

## A Deep Learning Approach to Automatic Call Routing

Rajiv Shah Director of Solution Architect and Professional Services

2019



We deliver the fastest big data analytics processing platform to run your analytics & machine learning in production, at scale



**300+** Direct customers

**50+ / 500+** Fortune / Organizations

## 5,000+

Large installations in production (OEM)

**25+** ISVs

## **GigaSpaces Select Customers**





## ABOUT THE USE CASE

This use case shows how to modernize existing software architecture for an efficient call center routing workflow

## **JSE CASE BENEFITS:**



Enhance Customer Experience with automatic routing that prevents customers from being buried in a hierarchical menu



Reduce Average Handle Time for optimized efficiency



Improve Customer Experience

Faster call routing to the correct agent means a more satisfied customer

Reduce Costs: lower AHT

Faster call resolution: Faster routing

+

Routing to correct agent

Enhanced System Agility

Higher agility when adding new categories or departments

# TECHNICAL CHALLENGES

Performance

Event Driven Architecture based on prediction criteria is required for optimal performance supporting peak events

#### Simplification

Leveraging existing opensource frameworks such as BigDL in a unified platform simplifies architectural complexity

#### Continuous ML Training

Continuous model training

based on previous transcribed calls + automatic training of alternative models ensure models with higher scoring



Automatic routing to the right agent for the perfect personalized experience



# Operationalizing AI Example – Automatic Call Routing



## LIVE DEMO: Instant Insight to Action

- Run Deep Learning with BigDL on transactional data in real-time for instant insights
- Trigger transactional workflows based on prediction criteria and scoring for real-time business impact
- Simplify architecture, eliminate GPU requirements & reduce component and cluster sprawl for optimal performance & TCO



#### Call session assistant



See the brain behind it!

Insightedge Web-UI

Spark Jobs

#### Model: TextClassification Training time: 10 min

Accuracy: 0.787037

BigDL v0.2.0

#### Call Center BigDL/InsightEdge module

Text

Click on the microphone icon and begin speaking.					
In-process calls (Powered by Intel BigDL) $_{\circ}$					

Search...

#### Call sessions 4

Id

Id	Category	Agent Id	Time (ms)	Text
4	comp.sys.mac.hardware	1	64	My Mac is just s***
3	comp.os.ms-windows.misc	7	67	My Microsoft computer sex
2	comp.sys.mac.hardware	0	94	Hey I have a problem with my Mac
1	comp.os.ms-windows.misc	4	114	Hello I have a problem with my windows





## GigaSpaces Coverage





## GigaSpaces Competitive Edge

**SPEED** 

Any Data

Live, **Transactional & Historical Data** 

**Deploy Anywhere** 



ANALYTICS

# Data Analytics: Undeniable Value to your Business

## Dynamic Pricing

Helps grow sales by **30% annually** 

## **Optimized Operations**

Saves **\$100sK** in annual savings (banking example)

### **Risk Analysis**

Reduces loan losses by 10-30%

## Call Center Automation

Increases efficiency by over 90%

### **Predictive Maintenance**

Reduces maintenance costs by up to **75%** per mile (transportation example)

### **Personalized Recommendation**

Increases conversions by **up to 20X** for brick & mortar stores via location-based promotions

### **Fraud Analytics**

Reduces losses by **3 to 5%** in mature environments and by **over 30%** in evolving contexts

# The Velocity of Business

"To prevent fraud, anomaly detection needs to happen against 500,000 txn/sec in less than 200 milliseconds" "A typical e-commerce website will experience 40% bounce if it loads in more than 3 seconds, including personalization offers" "A call center receives 450,000 calls/day, each call needs to be routed in less than 60 milliseconds"



FINANCIAL SERVICES



ECOMMERCE



**TELCO** 

## Use Cases Spanning Industries Benefit from Near Real-time AI **Decision Support Systems Built on GigaSpaces**



- Fraud
- Credit risk scoring
- Customer 360
- Customer churn

**FINANCIAL SERVICES** 



**INSURANCE** 

- Usage based insurance Customer 360
  - Customer churn
    - Claims management



- RETAIL ECOMMERCE
- Personal recommendations
- Intelligent inventory mgmt.
- Customer 360
- Locations-based promotions



- Predictive maintenance
- Fleet management
- Customer 360

TRANSPORTATION

- •
- Inventory planning
  - Customer 360  $\bullet$ 
    - Predictive maintenance





MEDIA/ **TELCO** 

- Customer 360 (incl. churn)
- Intelligent call center routing
- Data Center Infrastructure Monitoring (DCIM)
- Predictive maintenance

# InsightEdge: Unifying Real-Time Analytics, Al and Transactional Processing in One Platform

- Rich ML & DL support
- Extreme performance
- Fully Transactional
- ACID Compliance
- Enterprise-grade (Security, High Availability)
- Co-located Apps and Services
- Seamless integration with Big Data ecosystem
  - Data sources (Kafka/Nifi/Talend/etc.)
  - Data lakes (S3/Hadoop/etc.)
  - BI tools (Tableau/Looker/etc.)











## AnalyticsXtreme: Accelerating Your Data Lake by 100X for Real-time Analytics

Your data is immediately searchable, queryable, and available for analytics

- Single logical view for hot, warm and cold data
- Hot data resides on in-memory data grid and historical data on HDFS/Object Store
- Hot data is mutable and historical data is immutable (parquet)

## Fast Access

 Fast access to frequently used historical data

## Access any data through a unified layer

- Analytics (Spark ML)
- Query (Spark SQL)

## Automatic lifecycle management

• Automatically handles the underlying data movement, optimization and deletion





# Ultra-low latency and high throughput transactional processing IMDG



# Co-located Analytics and AI with Transactional Processing







- Persistent Memory +249% than SSD
- RAM (off-heap) +350% than SSD



- Persistent Memory +159% than SSD
- RAM (off-heap) +180% than SSD





- CAPEX reduction of up to 50% with RAM off-heap vs. on-heap
- CAPEX reduction of up to 75% with AEP vs. RAM on-heap
- OPEX reduction by X10





Define which data resides on which layer per class and per field

(\*) See vendor specifications

Figure 2: Memory-Storage Hierarchy with Persistent Memory Tier





https://builders.intel.com/persistent-memory-developer-challenge

#IntelDCISummit



## Kubernetes and Docker

VARIOUS DATA SOURCES







![](_page_29_Picture_0.jpeg)

## Leverage leading BI Platforms

#### Tableau

![](_page_29_Figure_3.jpeg)

#### Qlik

![](_page_29_Figure_5.jpeg)

![](_page_29_Picture_6.jpeg)

#### Looker

![](_page_29_Figure_8.jpeg)

#### Power BI

![](_page_29_Figure_10.jpeg)

#### NoSQL vs. GigaSpaces

![](_page_30_Figure_1.jpeg)

Per Node Replication Factor: 2 Record size:1KB RAM: 32GB CPU: 16 cores Disk: 1.2TB SSD

![](_page_30_Figure_3.jpeg)

■ HBase ■ IE SSD ■ IE PMEM ■ IE RAM

![](_page_30_Picture_5.jpeg)

![](_page_31_Picture_0.jpeg)

GigaSpaces is now focused on in-memory data processing... The combination of Spark and XAP will enable GigaSpaces to target the new breed of real-time analytics and hybrid operational and analytic workloads.

#### FROST 🔗

#### SULLIVAN

InsightEdge contains all the necessary SQL, Spark, Streaming, and Deep Learning toolkits for scalable data-driven solutions... our preferred solution components: the three-tier Kappa model, including Spark and Kafka, as implemented by GigaSpaces, in combination with its commercial InsightEdge platform.

> Everyone Wants "Real-time Analytic Insights" But Which Architecture Will Get You There?

# CASE STUDY: Fast Global Fabric for Risk, Trading and Market Data

#### **BUSINESS CHALLENGE:**

 Prior to executing a trade, a credit check needs to run and guarantee that the counterparty is not exceeding their limit

#### **TECHNICAL CHALLENGE:**

- Complete control over all eTrading platforms
- Regulatory enforcement set by RISK rules on all users trades on a daily basis
- Regulation analysis and checks
- Client onboarding
- Traversal framework
- Referential data for other apps

#### **IMPLEMENTATION:**

- All reservations, limits and client data is stored in the GigaSpaces in-memory platform
- All the requests are executed via the platform
- GigaSpaces is used in front of the database to speed up data access
- A worldwide deployment is done (Paris, NY and London) with GS asynchronous replication between each site to populate the data in NY and London

#### **RESULTS:**

- Three sites with 99.999 HA, replicated WW (Paris, London, New-York and Hong Kong)
- Reduced cluster and component sprawl
- Real-time risk analysis and credit checks complying with regulations

RISK

SOCIETE GENERALE

> SERVICES HUB

- Efficient scalable multi data-centre architecture
- Read: 700 K per day
- Write/Update/Remove : 20 K per day
- Next phase is to add ANOTHER site (TOKYO)

![](_page_32_Figure_23.jpeg)

![](_page_33_Picture_0.jpeg)

![](_page_33_Picture_1.jpeg)

#### **BUSINESS CHALLENGE:**

• Demand forecasting and price optimization in real-time based on threshold changes

#### TECHNICAL CHALLENGE:

- Ingest ~ billion of records in minutes
- Ability to query data from multiple geographies in real-time at low latency
- Ability to update with low latency multiple locations to adjust forecast and influence
- Cloud nativeness

#### **RESULTS:**

- Agility: Reduced forecasting ingestion from 3 hours to 8 minutes
- Live interactive querying and analytics through Spark SQL < 150ms latency

![](_page_33_Figure_12.jpeg)

![](_page_34_Picture_0.jpeg)

![](_page_34_Picture_1.jpeg)

#### **BUSINESS CHALLENGE:**

• Flight availability forecasting real-time based on various factors: date, city pair, #seats requested, marketing class, Point of Sale (PoS), quota limits, traffic restrictions, etc.

#### **TECHNICAL CHALLENGE:**

- Various internal systems (Reservation, Shopping, eCommerce Systems)
- Open API for external systems: Airlines, Global Distribution Systems (GDSs) and BOTs (automated searching).
- Auto scaling and sub-sec latency
- Multi tenancy (small/med/large airlines)

#### **RESULTS:**

- Querying and analytics response time < 50ms latency
- High Performance with up to 200K transaction/sec
- Scaling Near Linear (X100)
- Increase throughput by X& and reduce network overhead by 10%

LOOK TO BOOK 2000 1500 500 0 2014201520162017201820192020202120222023202420252026

> Ratio of Bookings per Availability Requests increases **by 100**

![](_page_35_Picture_0.jpeg)

## CASE STUDY: PriceRunner Compares Prices for Millions of Offers in Milliseconds

#### **BUSINESS CHALLENGE:**

• PriceRunner receives prices from 18,000 different merchants and has 4.4 million unique visitors per month, needed to ensure real-time comparisons for their customers at high peak periods such as the night before Black Friday where traffic increases between 10–20 times the normal traffic.

#### TECHNICAL CHALLENGE:

- Support scalability requirements at peaks without compromising performance
- No downtime
- Real-time analytics on transactional data
- Event-driven applications powering integrated applications
- Microservices architecture for rapid development and deployment

#### 18,000 different merchants

### 200 million prices updates

### 1 Billion requests a month

5-8 millisecond performance

![](_page_35_Figure_14.jpeg)

eCommerce

"Innovation is a key tenant of our strategy, and adoption of GigaSpaces InsightEdge real-time machine learning technology will highly differentiate our services by enabling us to run advanced analytics models on our hot data and instantly predict prices to improve the customer experience."

#### Roger Forsberg, CTO PriceRunner

![](_page_36_Picture_0.jpeg)

![](_page_36_Picture_1.jpeg)

#### **BUSINESS CHALLENGE:**

- Dynamic pricing engine based on CO2 tax regulations for B2B and B2C
- Many car configurations are unique, but all parts are not significant for CO2 calculation

#### **TECHNICAL CHALLENGE:**

- The current pricing engine workload is around 60 to 80 calculations/s, expected to increase to 2000
- Pricing calculations are obsolete after 24 hours.
- Each CO2 returned value must be exact
- All requests (both internal and external) must be equally treated

#### **RESULTS:**

- Querying and analytics response time < 100ms latency
- Reduce infrastructure footprint by a factor of X4-6
- Scaling up by X20

![](_page_36_Figure_14.jpeg)

Pricing Requests increases by 20x

![](_page_37_Picture_0.jpeg)

#### **BUSINESS CHALLENGE:**

- Enhance customer experience with quicker First Call Resolution
- Reduce Average Handle Time for optimized efficiency

#### TECHNICAL CHALLENGE:

- Ingestion of millions of CRM cases and data from other repositories into a unified analytics platform
- Leveraging ML models in real time
- Continuous model training

<b>DATABANK</b>	Q Search
Ticket ID #54367	DATA SO
Customer Name #54367	
Type Enterprise	International Payment declined
Support Level Bronze	
Last Contact Date 20.12.18	Case Description
	Read mo

#### URCES CUSTOMER CUSTOMER TICKET 95.32% Credit Limit exceeded #56409 93.05% Authentication required #33487 86.16% Beneficiary account unknown #180762 Beneficiary 77.98% account dormant #180762 71.53% Intermediary bank changes #60975 **Case Resolution Case Description** Check that credit limit Check here to email ment to is not exceeded instructions to customer Read more Email ore

**CONTACT CENTER** 

Reducing mean time to resolution by 5-10X Average time of 50ms to search and find similar cases

![](_page_38_Picture_0.jpeg)

## Fraud and Money Laundering Detection in Real-time

#### FINANCIAL SERVICES

#### **BUSINESS CHALLENGE:**

- Detecting fraud on mobile payment applications in real-time
- Detecting the deposit of the same check in multiple accounts at different banks in real-time
- User experience: application availability 24×7
- TCO reduction: reduce dependency on expensive RDBMS (Oracle)

#### TECHNICAL CHALLENGE:

- IMC Platform to ingest 4 TB of data daily
- Fully consistent transactional In-Memory Map-Reduce
- Millisecond response
- Analyze and validate against a large dataset of live (multiple TB) in memory and archived data (to Cassandra NoSQL and Hadoop)

#### **RESULTS:**

- Sub-second response for accurate fraud detection to stop the transaction
- TCO Reduction: RAM and SSD for runtime data compared to Oracle DB or SAN
- Fault-tolerant, highly available, scaling on demand

## Ingest **4 TB** daily Handle **1.5M** events per second

![](_page_38_Figure_18.jpeg)

![](_page_39_Picture_0.jpeg)

## Instant Payments for real-time transactions and high reliability to enhance the overall customer experience

## FINANCIAL SERVICES

#### **BUSINESS CHALLENGE:**

- Enable and accelerate instant payment solutions and meet regulatory requirements on a global scale
- Automatically track purchases and other server-to-server communication
   in real time
- Store payment transactions, order information and other sales internally

#### TECHNICAL CHALLENGE:

- Ability to handle added data volumes 15k payment/sec receipts introduced by management of new SEPA European payment regulation
  - Assure no-downtime for mission critical service

#### **RESULTS:**

- Running low-latency payment and business logic calculations
- No downtime assured
- Real-time analytics and Machine Learning preventing fraud and adherence to regulations
- Design to deployment in just a few months leveraging microservices architecture

![](_page_39_Figure_15.jpeg)

Payment transaction in **500** milliseconds End-to-end validation in seconds

![](_page_40_Picture_0.jpeg)

## CASE STUDY: PriceRunner Compares Prices for Millions of Offers in Milliseconds

#### **BUSINESS CHALLENGE:**

• PriceRunner receives prices from 18,000 different merchants and has 4.4 million unique visitors per month, needed to ensure real-time comparisons for their customers at high peak periods such as the night before Black Friday where traffic increases between 10–20 times the normal traffic.

#### TECHNICAL CHALLENGE:

- Support scalability requirements at peaks without compromising performance
- No downtime
- Real-time analytics on transactional data
- Event-driven applications powering integrated applications
- Microservices architecture for rapid development and deployment

#### 18,000 different merchants

### 200 million prices updates

### 1 Billion requests a month

5-8 millisecond performance

![](_page_40_Figure_14.jpeg)

eCommerce

"Innovation is a key tenant of our strategy, and adoption of GigaSpaces InsightEdge real-time machine learning technology will highly differentiate our services by enabling us to run advanced analytics models on our hot data and instantly predict prices to improve the customer experience."

#### Roger Forsberg, CTO PriceRunner

![](_page_41_Picture_0.jpeg)

#### **BUSINESS CHALLENGE:**

- Enhance customer experience with quicker First Call Resolution
- Reduce Average Handle Time for optimized efficiency

#### TECHNICAL CHALLENGE:

- Ingestion of millions of CRM cases and data from other repositories into a unified analytics platform
- Leveraging ML models in real time
- Continuous model training

<b>DATABANK</b>	Q Search
Ticket ID #54367	DATA SO
Customer Name #54367	
Type Enterprise	International Payment declined
Support Level Bronze	
Last Contact Date 20.12.18	Case Description
	Read mo

#### URCES CUSTOMER CUSTOMER TICKET 95.32% Credit Limit exceeded #56409 93.05% Authentication required #33487 86.16% Beneficiary account unknown #180762 Beneficiary 77.98% account dormant #180762 71.53% Intermediary bank changes #60975 **Case Resolution Case Description** Check that credit limit Check here to email ment to is not exceeded instructions to customer Read more Email ore

**CONTACT CENTER** 

Reducing mean time to resolution by 5-10X Average time of 50ms to search and find similar cases

![](_page_42_Picture_0.jpeg)

![](_page_42_Picture_1.jpeg)

![](_page_42_Picture_2.jpeg)

![](_page_42_Picture_3.jpeg)

![](_page_42_Picture_4.jpeg)

## EXTREME PERFORMANCE

![](_page_42_Picture_6.jpeg)

![](_page_42_Picture_7.jpeg)

MISSION CRITICAL AVAILABILITY

sec from data to insight to action

# millions

of IOPS

10X less expensive than only RAM with In-memory performance No Downtime at leading enterprise customers for

YEARS And still counting

![](_page_43_Picture_0.jpeg)

## WHY GIGASPACES?

Real-time insights
Boost your performance
Simplify your architecture
Lower TCO / Enhance ROI

![](_page_44_Picture_0.jpeg)

### Enterprise Grade System of Record

![](_page_44_Picture_2.jpeg)

#### **Optimized Data Replication:**

Field-proven, reliable, high performance replication mechanism to replicate data between peer nodes in the data grid

![](_page_44_Picture_5.jpeg)

#### Data Partitioning:

Transparent content-based data partitioning to evenly and intelligently distribute data across your cluster

![](_page_44_Picture_8.jpeg)

#### Transaction Support:

Full transaction support, including local, distributed and XA transactions

#### Write Behind:

Asynchronous and reliable propagation of data to any external data source

#### Locking Support:

RDBMS locking and transaction isolation for robust and hassle-free data access

![](_page_44_Picture_16.jpeg)

Apache ZooKeeper\*

0

Multi-Site Deployment: Replicate and share data between multiple, geographically-distributed, active clusters for global activity

**Network Segmentation Protection:** Ensure data remains consistent in case of network segmentations of all types

#### Security:

Role-based authentication for data and operations, Support for Kerberos, Spring, TLS and more

#### Change API: Update data by specifying only the required change instead of the entire updated object

♥ 3

🔶 👍

10

![](_page_44_Picture_22.jpeg)

![](_page_44_Picture_23.jpeg)

![](_page_44_Picture_24.jpeg)

![](_page_44_Picture_26.jpeg)

Querving:

Sophisticated query engine

Customize the query's result

set by defining which fields

with support for SQL and

example queries

Projection API

should be returned

#### Aggregations

Sum, Avg, Min, Max, GroupBy and more, or even your own user-defined aggregations

![](_page_44_Picture_30.jpeg)

Advanced Querying & Indexing

Indexing:

Predefined and add-hoc Property indexing for fast data access

![](_page_44_Figure_33.jpeg)

Geospatial:

Enhance your data model with shapes and use spatial operations to find matches

### Full Text Search:

Go beyond plain text with regular expressions, fuzzy search, proximity matching and more

![](_page_44_Picture_38.jpeg)

SOL Functions:

Abs, Round, Length, Upper, Lower and more, or even your own userdefined functions

![](_page_44_Figure_41.jpeg)

![](_page_45_Picture_0.jpeg)

## Data Model Flexibility & Interoperability

![](_page_45_Picture_2.jpeg)

Native: Highly optimized, POJO driven API which exposes all the unique capabilities of the platform

![](_page_45_Picture_4.jpeg)

#### JPA:

Support data grid access using the standard JPA API for seamlessly scaling your JEE data access layer

#### Document:

Completely schema-free data API that supports upgrading the application's data model on the fly

![](_page_45_Picture_10.jpeg)

Key-Value: Simple and intuitive Map-based interface for simple caching scenarios

# 

Microsoft

#### **REST API:** Standard REST endpoint provides access to the data grid from any app, Platform and programming language

Native C# interface that enables

any .NET application to access

#### **Cross Language Access:**

.Net:

the data grid

support for heterogeneous environments, with seamless interoperability among them all

## Messaging & Event Features

Publish/Subscribe Messaging:

takes place in the data grid to

publish/subscribe paradigm

Support for implementation of

triggering of processing logic

Point-to-Point Messaging:

complex workflows and

across the data grid

listeners using the

![](_page_45_Picture_17.jpeg)

![](_page_45_Picture_18.jpeg)

![](_page_45_Picture_19.jpeg)

![](_page_45_Picture_23.jpeg)

![](_page_45_Picture_25.jpeg)

Content Based Routing: Routing of events to relevant cluster members based on their content

# Propagation of any event that

![](_page_45_Picture_30.jpeg)

Fully durable pub/sub messaging for data consistency and reliability

![](_page_45_Picture_32.jpeg)

#### **FIFO Groups**

Ensure in-order and exclusive processing of events belonging to the same group, while parallelizing across groups

### Workflow Support:

![](_page_45_Picture_36.jpeg)

![](_page_45_Picture_37.jpeg)

![](_page_46_Picture_0.jpeg)

### Collocation of Data and Business Logic

![](_page_46_Picture_2.jpeg)

Spring on Steroids:

Deployment, provisioning and proactive management of any spring application, with or without a data arid

![](_page_46_Picture_5.jpeg)

Master-Worker Support: Intuitive and highly scalable master-worker implementation for distributing computationintensive tasks

#### **Dynamic Code Execution:** Dynamic code shipment and map/reduce-like execution across the grid for optimized processing and data access

## event 🧼

Event Tracking: Trace Cluster Events as they happen event 🛞 event 🥥 for improved visibility & easier event

troubleshooting (available through both admin API and UI)

![](_page_46_Picture_11.jpeg)

#### Alerts:

**UI Based Management** 

Web-based dashboard app for

single Space instance queries.

easy monitoring & management of

deployed app. Enhanced data grid

console for cluster wide queries or

Out-of-the-box identification & notification of risky situations (e.g., above-normal CPU utilization or data replication failure)

![](_page_46_Picture_14.jpeg)

#### **Application Dependencies:** Deploy modules as an application ensuring order of deployment

![](_page_46_Picture_17.jpeg)

#### Client side Cache Monitoring: Discover client-side cache and views connected to your spaces

![](_page_46_Figure_19.jpeg)

#### WAN Replication Monitoring: Discover client-side cache and views connected to your spaces

![](_page_46_Picture_21.jpeg)

Code and Data Collocation: Deployment of business logic and data as a single coherent unit for optimized performance

![](_page_46_Figure_23.jpeg)

**Robust Remoting Support:** Built on top of the data grid to provide fault tolerance, service auto discovery, cluster wide invocations and more

#### Security:

Customizable security policy to control who can run dynamic code on the arid

![](_page_46_Picture_28.jpeg)

![](_page_46_Picture_29.jpeg)

![](_page_46_Picture_31.jpeg)

**Event Containers Monitoring:** Trace embedded and remote event containers

![](_page_46_Picture_33.jpeg)

**Extensible Metrics Framework:** Measure both space and userdefined metrics, integrated with any tool (InfluxDB and Grafana out of the box)

![](_page_46_Picture_35.jpeg)

**REST Admin API:** Comprehensive and intuitive API for monitoring and controlling every aspect of your cluster and application

![](_page_46_Picture_37.jpeg)

Grid Health Transparency & Monitoring

![](_page_46_Picture_38.jpeg)

# THANK YOU

BUILD IT >

TRY IT >

![](_page_47_Picture_3.jpeg)

innovate with confidence