# Global Big Data Conference

# BIG DATA BOOTCAMP

Tampa

December 9th, 10th & 11th 2016

Tampa Convention Center, 333 S Franklin St, Tampa, FL 33602

www.globalbigdataconference.com

Twitter : @bigdataconf

# IRI, The CoSort Company

## Vendor Background

- ISV specializing in data management and data protection

- Known since 1978 for "big data" transformation speed

- 7 of 8 software products share 1 metadata and Eclipse GUI

- A 'top big data provider' (CIO Review & Insight Success)

- Headquartered 1 hour southeast of Orlando, FL

- Resellers in more than 40 international cities

- Customers in every industry with big and/or sensitive data

**IRI**
The CoSort Company

# Selected IRI Customers

IRI customers process and protect data off the mainframe, for DW ETL/ODS ops, and in PII protection (privacy law compliance) initiatives. Hadoop use is optional. Most work with big and/or sensitive financial, call/click, or healthcare data.

## IRI Data Manager Suite

### FACT
*IRI FAst extraCT*

*Speed DB unloads for archival, migration, reorg and ETL*

• Extract tables to flat files in parallel using SQL queries
• Convert and re-format to change data types and layouts
• Create the data definitions for IRI software and DB loads
• Pipe to CoSort and DB loaders for faster reorg and ETL

### CoSORT
THE OPEN SYSTEMS STANDARD

*Speed or replace batch, BI, ETL, sort, and SQL programs*

• Filter, sort, join, aggregate, pivot, cleanse, lookup, calc, etc.
• Map, migrate, federate and replicate data from 125 sources
• Segment data, capture changes, report details / summaries
• Analyze changing dimensions, support complex transforms

### NextForm
Data & Database Migration

*Unlock data and move between apps, DBs, and platforms*

• Convert, federate, remap, and replicate legacy data
• Migrate data between databases and create new tables
• Change file formats, data types, and endian conditions
• Search and structure data in "dark data" documents

---

## Voracity
IRI
*Total Data Management*

*Consolidate tools and tasks to process, protect, prototype, present*

• Discover, define, and govern data in legacy and new sources
• Combine data integration, migration, protection, and analytics
• Exploit CoSort and Hadoop engines for optimum throughput
• Leverage Eclipse familiarity, functionality, and extensibility

eclipse

---

## IRI Data Protector Suite

### RowGen
Safe Intelligent Test Data

*Prototype DBs and ETL, stress-test, outsource, benchmark*

• Use real data models and formats, not production data
• Combine generation and selection, create new formats
• Preserve referential integrity and frequency distributions
• Feed test DBs, files, and custom reports simultaneously

### FieldShield
Data Masking & Encryption Solutions

*Comply with privacy laws, nullify breaches, govern data*

• Select shields for each field per business rules
• De-ID, encrypt, hash, mask, pseudonym, random, token
• Apply cross-table rules to save time and referential integrity
• Create an XML audit log of each job to verify compliance

### CellShield
Data Masking Add-In for Excel

*Profile and protect PAN/PHI/PII in Excel spreadsheets*

• Search and save patterns to discover sensitive data
• Locate, report, and open all found ranges in the LAN
• Click to encrypt, mask or pseudonymize data directly
• Auto-log protections to verify privacy law compliance

### Chakra Max
*Smart Data-Centric Audit & Protection*

*Define, monitor, block, and audit DB access*

• High-volume, data-centric audit and protection (DCAP)
• Monitor, block, alert, and log users in real-time
• Low-impact on DB performance and availability
• Classify and dynamically mask sensitive data with RBAC

---

Embedded or callable analytics:
BIRT, JupiterOne, NextCoder, R

### IRI
The CoSort Company

# Address the Challenges of Big Data

| Volume | Variety | Velocity | Veracity | Value |
|--------|---------|----------|----------|-------|
| BI and analytic tools choke on high volumes; they drag, hang or crash | The myriad of structured and unstructured sources is beyond most tools | IOT logs, dark data, CDRs, etc. are generated too fast for analysis | Garbage in=garbage out: low quality data jeopardizes analytic value | Without tackling the above, you won't get analytic value from big data |
| *Voracity blends and prepares data for analytic tools via **fast, combinatory transforms** like: filter, sort, join, aggregate and segment. Programs built on the CoSort SortCL language hand off digestible data chunks or cubes to BIRT, Qlik, R, SAS, Splunk, Tableau, etc.* | *Voracity either natively, or through partner drivers, connects to and integrates **>125 data sources** on premise or in the cloud. They can be structured, semi-structured, or unstructured, and static and streaming.* | *Voracity processes **streaming data** from: web services and brokers (MQTT, Kafka); pipes; in Hadoop Spark or Storm; SQL; and, through memory via input procedure calls to CoSort. Voracity's built-in task launcher can also run jobs in near-real-time.* | *Voracity's data **discovery and quality** features let you: search for strings and patterns, do fuzzy matching, validate, scrub, enrich, and unify data for DW/BI, MDM, and analytics.* | *Voracity runs with or without Hadoop on commodity hardware under an **affordable subscription** model based only on the number (not size) of servers. Its **Eclipse GUI** is free, familiar, and flexible, to speed learning and time-to-solution.* |

# Supported Data Sources/Targets:

| | | | |
|---|---|---|---|
| Amazon EMR Hive | FinancialForce | Marketo | Pivotal Greenplum |
| Apache Cassandra | Force.com apps | MongoDB | Pivotal HD Hive |
| Apache Hadoop Hive | Hortonworks Hive | MS Dynamics CRM | Salesforce.com |
| Cloudera CDH Hive | Hubspot | MS SQL Azure | ServiceMAX |
| Cloudera Impala | Lightning Connect | Oracle Eloqua | Spark SQL |
| Database.com | MapR Hive | Oracle Service Cloud | Veeva CRM |

*… plus 'legacy list' on next 2 pages >>*

**IRI** The CoSort Company

# Global Big Data Conference

| | | | |
|---|---|---|---|
| Acucobol Vision | Delimited | MaxDB | SQL Server |
| Altibase (FACT) | Derby (WB) | Mongo (WB) | SQLite |
| ASN.1 TAP3 | ESDS | MF-ISAM | Sybase ASA/E & IQ |
| BIRT DB (WB) | Excel (WB) | WF Var. Length | Tibero (WB) |
| BIRT Hive (WB) | ELF web logs | MySQL | Teradata (WB) |
| BIRT JDBC (WB) | Fixed | Oracle | Text |
| BIRT POJO (WB) | Heap / print | Outlook (WB) | UTF-8 & 16 |
| C-ISAM | HSQLDB (WB) | PDF (WB) | Variable Block |
| CLF web logs | IDX 3, 4 & 8 | PostgreSQL | Variable Sequential |
| CSV | Informix | Powerpoint (WB) | VSAM MVS (UniKix) |
| DB2 (UDB) | Ingres | Record Sequential | Web Services (WB) |
| DB2 for i5/OS (WB) | LDIF | RTF (WB) | Word (WB) |
| DB2 for z/OS (WB) | Line Sequential | SQL Anywhere | XML |

**IRI**
The CoSort Company

| Access | D3 | GA-Power 95, R91 | K-ISAM | Pathway | RMS |
|--------|-----|------------------|--------|---------|-----|
| Adabas | Datacom | Gemstone | Knowledgeman | PDS | Reality/X |
| Advanced Pick | Dataflex | GENESIS | KSDS | PervasiveSQL | RRDS |
| ALLBASE | Db4o | Gigabase | Lotus | Pick/Pick64+ | SAP HANA |
| Alpha5 | dBase | H2 | Manman | PI-Open | Sequoia |
| Amazon RDS | Desktop Adapter | IDMS | Mentor / pro | Powerflex | Sharebase |
| Azure | DL/1 | IDS | MO | Powerhouse | Supra |
| BizTalk | DSM | Image | Model 204 | Progress | Terracotta |
| Cache | Enscribe | IMS | Mumps | QueryObject | Total |
| Clipper | Enterprise Adapter | Interbase | MyBase | rBase | Ultimate |
| Codasyl | FileMaker | Intersystems | Netezza | R83 | UltPlus |
| CorVision | Firebird | ISM | NonStop SQL | Rdb | Unidata |
| ConceptBase | Focus | Jasmine | ObjectStore | REALITY | Universe |
| D-ISAM | FoxPro | JBase | Paradox | Red Brick | VSAM VSE |

# Global Big Data Conference

## Sources

**Big Data**
cloudera · MAPR · Hortonworks · Spark · HIVE · Pivotal · SAP HANA · Netezza

**Call Detail Records**
ASN.1 Formats

**Cloud & SaaS**
amazon web services · Microsoft Azure · salesforce · HubSpot · Marketo · eloqua

**Databases**
ORACLE DATABASE · IBM DB2 · MySQL · TERADATA · mongoDB · Microsoft SQL Server · TIBERO DataBase · SYBASE · DATASTAX · ALTIBASE

**Files & Pipes**
COBOL, CSV, LDIF, LS-RS-VS, MFVL, Text, VB, Vision, XML

**Mainframe**
Adabas, Datacom, IDMS, IMS, ISAM, Pick, Unidata, VSAM, etc.

**Semi & Unstructured**

**Other Sources**
Custom Apps, ETL/ELT Tools, Packaged Apps, Web Logs

## Targets

**Big Data**
cloudera · MAPR · Hortonworks · Spark · HIVE · Pivotal · SAP HANA · Netezza

**BI & Analytic Tools**
BIRT · SAP BusinessObjects · Excel · IBM Cognos · MicroStrategy · ORACLE BUSINESS INTELLIGENCE · Qlik Q · R · splunk> · Spotfire · tableau

**Cloud & SaaS**
amazon web services · Microsoft Azure · salesforce · HubSpot · Marketo · eloqua

**Custom Reports**
Detail & summary reports

**Databases**
ORACLE DATABASE · IBM DB2 · MySQL · TERADATA · mongoDB · Microsoft SQL Server · TIBERO DataBase · SYBASE · DATASTAX · ALTIBASE

**Files & Pipes**
COBOL, CSV, LDIF, LS-RS-VS, MFVL, Text, VB, Vision, XML

**Other Targets**
Custom Apps, Data & SpreadMarts, ETL/ELT Tools, Federated Views, Packaged Apps, Test Suites

## DISCOVER
Data Classification
Dark Data Search
DB & File Profiling
ER Diagramming
Metadata Definitions
Metadata Forensics
Multi-Method Search

## INTEGRATE
Public/Private Mashup
Change Data Capture
Bulk DB Un/Load
Data Federation
One Pass ETL

## MIGRATE
Data & File Types
Endianness
Databases
ETL Jobs
JCL Sorts

## GOVERN
Data Lineage
Data Masking
Data Quality
Metadata & MDM
Test Data Generation

## ANALYZE
Embedded BI
BIRT & Splunk Feeds
Clickstream Analytics
Customer Segmentation
Slowly Changing Dimensions

## IRI Voracity
Total Data Management

## DESIGN
ADS Mapping Manager
Form Editors
Graphical Dialogs
Outlines & Palettes
Script Editors
Visual Workflow
Wizards & Rules

## DEPLOY
CoSort CLI/API (SMP)
Eclipse & Other Job Launchers
Java, Paques, SQL
MapReduce (Grid)
Spark (In-Memory)
Storm (Streaming)
Tez (Batch)

PartnerNet SILVER PARTNER Novell · intel Software Partner · IBM Business Partner · Business Partner hp · SOLUSI 247 24 hours 7 days integrated ICT solution · AnalytiXDS · eclipse · msdn Microsoft Developer Network · MICRO FOCUS · TRILLIUM SOFTWARE A Harte Hanks Company · redhat · ORACLE PARTNERNETWORK

www.globalbigdataconference.com

Voracity includes PII discovery facilities for multi-source data **classification**, string (literal or in-dictionary), pattern, and fuzzy-match **searches**, statistical **reports**, and automatic **metadata** creation. Fit-for-purpose wizards in Voracity perform:



- Data classification, with rule matcher libraries
- DB profiling and E-R diagramming
- Dark data discovery and structuring, with forensic metadata display
- Flat-file statistical and value searching
- Metadata discovery and definition
- Metadata sharing, lineage tracking, etc.

Voracity combines fast ETL engines and task consolidation techniques with simple metadata in Eclipse that's shared by all IRI software and other products, like AnalytiX DS for ETL code conversion. You can use Voracity to speed or *re-platform* megavendor tools, and optimize:

- EDW, LDW, ODS, data lakes
- Data quality (cleansing)
- VLDB unload/reorg/load jobs
- SCD, CDC, pivoting, unification

*Job Design …*

In addition to GUI wizards, diagrams, and dialogs, you can also hand-code the underlying 4GL programs in Voracity's syntax-aware editor.

This job sorts and filters an employee CSV file into two target files, while also redacting ID #'s and commissions, and encrypting the salary.

IRI
The CoSort Company

*Job Deployment …*

Voracity's 4GL scripts run on the command line or in batch from the GUI or shell.

BIRT or Splunk can also run them as they report or index.

Voracity can also schedule and run them seamlessly in MR2, Spark, Spark Stream, Storm or Tez.



Map once, deploy anywhere

*Preparing a run configuration for* ***Hadoop*** *...*

Once our gateway is open, we can tell any job to run in Hadoop.

Here, we specify MR2 as the engine, and our working directory in HDFS.

**DISCOVER**   **INTEGRATE**   **MIGRATE**   **GOVERN**   **ANALYZE**

*Voracity*
*Total Data Management*

*The Job Manager view shows our Hadoop job running, plus the status of other jobs.*

**hadoop-demo - IRI Development - demo-hdfs/employee-dept.scl - IRI Workbench**

File   Edit   Navigate   Search   Project   Run   Window   Help

Quick Access

Console   Problems   Properties   Scheduler   HDFS Browser   **Job Manager** ✕

| ID | Name | Engine | Status | User | Start | End |
|---|---|---|---|---|---|---|
| 0000005-1611301139275-oozie-oozi-W | demo | MR2 | RUNNING | yava | Thu, 01 Dec 2016 19:58:18 GMT | null |
| 0000004-1611301139275-oozie-oozi-W | demo | MR2 | SUCCEEDED | yava | Thu, 01 Dec 2016 19:43:28 GMT | Thu, 01 Dec 2016 19:44:38 GMT |
| 0000003-1611301139275-oozie-oozi-W | demo | MR2 | SUCCEEDED | yava | Thu, 01 Dec 2016 19:39:14 GMT | Thu, 01 Dec 2016 19:40:24 GMT |
| 0000002-1611301139275-oozie-oozi-W | demo | MR2 | SUCCEEDED | yava | Thu, 01 Dec 2016 19:30:12 GMT | Thu, 01 Dec 2016 19:31:22 GMT |
| 0000001-1611301139275-oozie-oozi-W | demo | MR2 | SUCCEEDED | yava | Thu, 01 Dec 2016 18:44:12 GMT | Thu, 01 Dec 2016 18:44:48 GMT |
| 0000000-1611301139275-oozie-oozi-W | demo | MR2 | SUCCEEDED | yava | Thu, 01 Dec 2016 18:41:20 GMT | Thu, 01 Dec 2016 18:42:08 GMT |
| 0000000-161129080054185-oozie-oozi-W | Hadoop | MR2 | SUCCEEDED | yava | Wed, 30 Nov 2016 14:17:18 GMT | Wed, 30 Nov 2016 14:34:06 GMT |
| 0000040-161116100421533-oozie-oozi-W | Ex11 | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 20:25:25 GMT | Tue, 29 Nov 2016 14:19:09 GMT |
| 0000039-161116100421533-oozie-oozi-W | Ex11 | MR2 | KILLED | yava | Mon, 28 Nov 2016 18:32:30 GMT | Tue, 29 Nov 2016 14:29:09 GMT |
| 0000038-161116100421533-oozie-oozi-W | chiefs | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 18:27:20 GMT | Tue, 29 Nov 2016 14:18:24 GMT |
| 0000037-161116100421533-oozie-oozi-W | Ex11 | MR2 | KILLED | yava | Mon, 28 Nov 2016 18:24:21 GMT | Tue, 29 Nov 2016 14:18:10 GMT |
| 0000036-161116100421533-oozie-oozi-W | Ex11 | MR2 | KILLED | yava | Mon, 28 Nov 2016 18:24:04 GMT | Tue, 29 Nov 2016 14:18:09 GMT |
| 0000035-161116100421533-oozie-oozi-W | Ex11 | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 17:31:58 GMT | Mon, 28 Nov 2016 17:32:12 GMT |
| 0000034-161116100421533-oozie-oozi-W | Ex11 | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 17:20:56 GMT | Mon, 28 Nov 2016 17:21:10 GMT |
| 0000033-161116100421533-oozie-oozi-W | Ex11 | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 17:19:35 GMT | Mon, 28 Nov 2016 17:19:49 GMT |
| 0000032-161116100421533-oozie-oozi-W | Ex11 | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 17:17:57 GMT | Mon, 28 Nov 2016 17:18:11 GMT |
| 0000031-161116100421533-oozie-oozi-W | Ex11 | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 16:51:29 GMT | Mon, 28 Nov 2016 16:52:04 GMT |
| 0000030-161116100421533-oozie-oozi-W | chiefs | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 16:32:06 GMT | Mon, 28 Nov 2016 16:32:41 GMT |
| 0000029-161116100421533-oozie-oozi-W | Hadoop | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 15:29:27 GMT | Mon, 28 Nov 2016 15:30:07 GMT |
| 0000028-161116100421533-oozie-oozi-W | Hadoop | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 15:11:38 GMT | Mon, 28 Nov 2016 15:12:17 GMT |
| 0000027-161116100421533-oozie-oozi-W | Hadoop | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 14:48:27 GMT | Mon, 28 Nov 2016 14:49:34 GMT |
| 0000026-161116100421533-oozie-oozi-W | Hadoop | MR2 | SUCCEEDED | yava | Mon, 28 Nov 2016 14:47:42 GMT | Mon, 28 Nov 2016 14:48:52 GMT |
| 0000025-161116100421533-oozie-oozi-W | Hadoop | MR2 | SUCCEEDED | yava | Mon, 21 Nov 2016 15:24:42 GMT | Mon, 21 Nov 2016 15:25:57 GMT |
| 0000024-161116100421533-oozie-oozi-W | tester | MR2 | SUCCEEDED | yava | Mon, 21 Nov 2016 03:54:02 GMT | Mon, 21 Nov 2016 03:55:17 GMT |
| 0000023-161116100421533-oozie-oozi-W | tester | MR2 | SUCCEEDED | yava | Fri, 18 Nov 2016 12:24:32 GMT | Fri, 18 Nov 2016 12:25:42 GMT |
| 0000022-161116100421533-oozie-oozi-W | tester | MR2 | SUCCEEDED | yava | Fri, 18 Nov 2016 12:13:38 GMT | Fri, 18 Nov 2016 12:14:48 GMT |
| 0000021-161116100421533-oozie-oozi-W | tester | MR2 | SUCCEEDED | yava | Fri, 18 Nov 2016 12:07:27 GMT | Fri, 18 Nov 2016 12:08:37 GMT |
| 0000020-161116100421533-oozie-oozi-W | encaes256 | MR2 | SUCCEEDED | yava | Fri, 18 Nov 2016 09:41:22 GMT | Fri, 18 Nov 2016 09:41:57 GMT |

**IRI**
The CoSort Company

*The HDFS Browser and Data Viewer show the target file and its contents ..*

You can also use the viewer window to manage all of your input and output data directly in HDFS..

*Wizards for ...*

**Slowly Changing Dimensions**

**Change Data Capture**
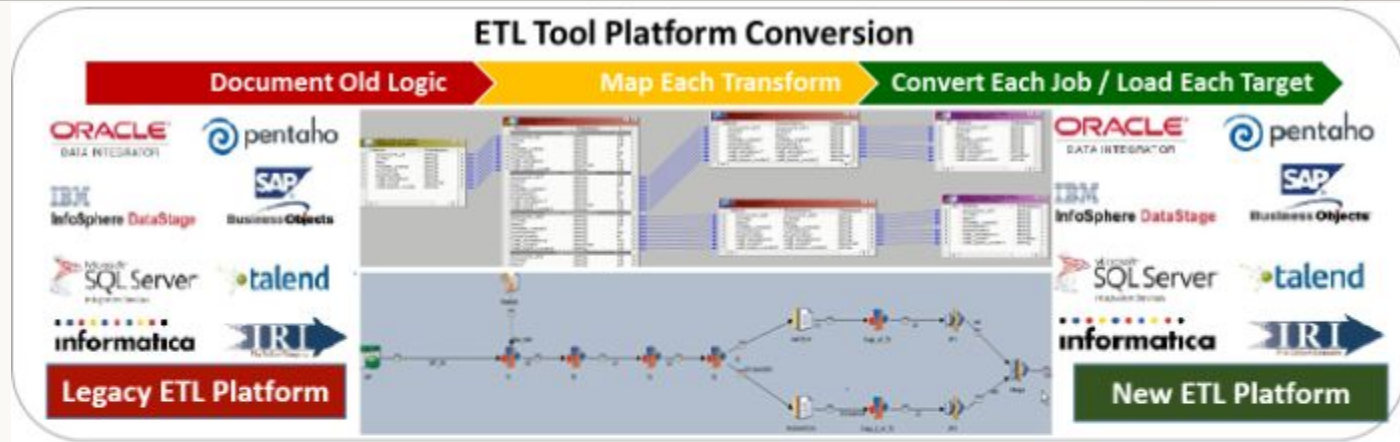
Global Big Data Conference

DISCOVER **INTEGRATE** MIGRATE GOVERN ANALYZE
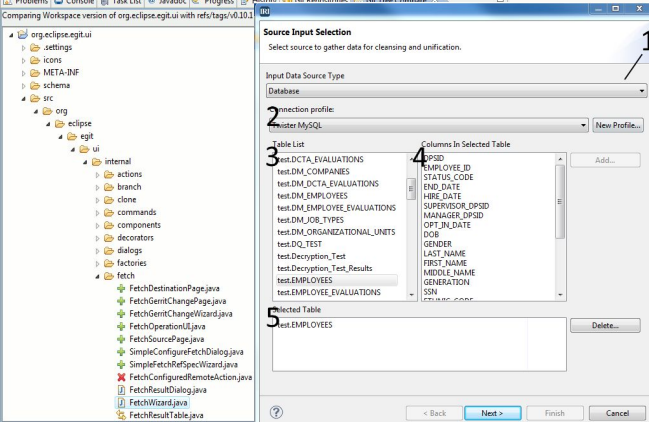
*With AnalytiX DS, ETL tool and SQL users can **convert** their existing data integration jobs to faster, simpler, and **far less expensive** Voracity workflows.*

**ETL Tool Platform Conversion**

Document Old Logic → Map Each Transform → Convert Each Job / Load Each Target

Legacy ETL Platform → New ETL Platform

Performance (like Ab Initio or Teradata)
Capability (like Informatica or DataStage)
DB affinity (like SSIS or ODI)
Eclipse ergonomics (like Talend)
Affordability (like Pentaho)

Voracity
Total Data Management

IRI
The CoSort Company

**DISCOVER    INTEGRATE    MIGRATE    GOVERN    ANALYZE**

Voracity **converts, replicates, and reformats** data from mainframe datasets, relational and NoSQL databases, index and sequential files, dark data documents, and cloud apps.



• Change data types, record layouts, file formats, and endianness

• Migrate column values, layouts, and relationships (constraints) between DBs

• Copy or update data from one or more sources to one or more targets

• Federate, or virtualize, data by mashing up data from disparate sources and creating custom, ad hoc views

IRI
The CoSort Company

Voracity's data governance and information stewardship features include:

- **Master data** management
- Data **class** and **rule libraries**
- Enterprise **metadata management**
- Static and dynamic **data masking**
- **Test data** generation & management
- **DB firewall** (via IRI Chakra Max)

- Connect and interact with **multiple sources** and targets, on-prem or cloud
- **Discover and classify** data in DB, flat-file, and dark-data (document) sources
- Mask **static or streaming** inputs, NoSQL DBs, and files in LUW, HDFS and S3
- Select from **12 masking categories** (e.g., encrypt, hash, pseudonymize, redact)
- **Address multiple** protections, targets and recipients all in one job, one I/O
- Apply consistent, cross-table masking rules for **referential integrity**
- Support **conditional security**, based on patterns, values, or ranges
- Specify target protections and formats in **Eclipse or portable** job scripts
- Integrate with **DB apps** via ODBC. Use .NET and Java SDK for dynamic masking
- Retain data **realism via FPE** and pseudonymization for testing, outsourcing
- **Mask during** big data ETL, migration, sub-setting, and BI/analytic jobs
- Log job and system runtime detail to XML audit files to **verify compliance**

Masking Features

IRI
The CoSort Company

Masking Complex XML

- Create synthetic but realistic **random and random-real** test data simultaneously

- Improve **DB prototypes**, application quality, benchmarking, and devops

- Leverage DDL, production file, and/or custom metadata

- Preserve structural and **referential integrity**

- Produce data in any type, structure, volume, value range, and "if" condition

- Synthesize **composite values** and custom (master) data formats

- Generate computationally valid and invalid NID, SSN, or CC#

- Set and graph test data **value distributions** (linear, normal, random, etc.)

- Apply common attribute rules (e.g., lookups) for pattern-matched field names

- **Filter, transform, and pre-sort** test data as you generate it

- Write loader metadata, and perform the loading, automatically

- Build test flat-file and custom detail and summary reports

- **Subset and mask** databases automatically as an alternative approach

- Use Java SDK functions to generate test data in apps and Hadoop
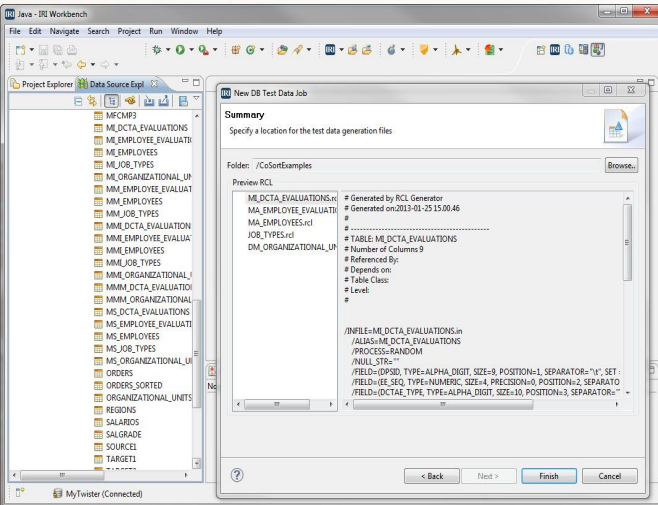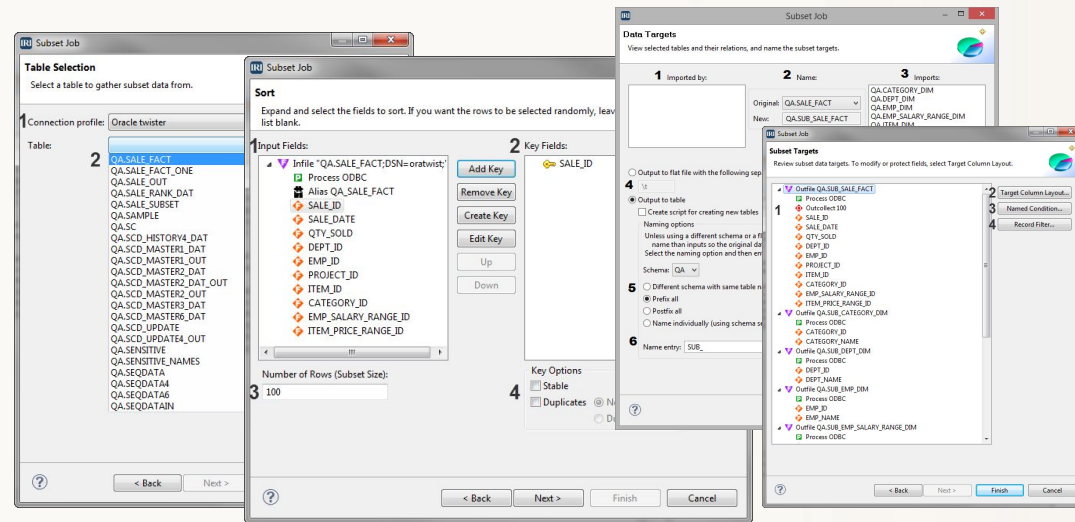
TDM Features

IRI
The CoSort Company

## Synthetic Data for:

- Flat files
- EDW ETL tools
- RDB & NoSQL
- Data lakes
- Mainframe jobs
- SAP, Teradata
- Cloud/SaaS apps

TDM Features

Both test data generation/population and DB subsetting wizards with built-in data masking are included in Voracity to facilitate DB and EDW prototyping. Either way, the test data is realistic, referentially-correct, and privacy-law compliant.

From its one Eclipse IDE (IRI Workbench) Voracity supports multiple analytic approaches …

*Voracity Analytic
Option 1:*
Embedded BI

Unlimited 2D reporting
in custom-formatted,
detail and summary
files, XML, HTML, etc.

*Voracity Analytic Option 2:*
BIRT Integration

Prepare and present data simultaneously from an "IRI Data Source" in BIRT

*Voracity Analytic
Option 3:*
Cloud Dashboard

Leverage drill-down,
browser-based
dashboard applications
like this one from
NextCoder

*Voracity Analytic Option 4:*
Splunk Add-On

Prepare data you need to index ad hoc, with a Voracity job launched from Splunk

*Voracity Analytic Option 5:*
Data Blending

Prepare CSV, XML or table subsets to reduce time-to-display 2-20X, along with data quality, privacy, and storage

# Global Big Data Conference

DISCOVER    INTEGRATE    MIGRATE    GOVERN    **ANALYZE**

*Data Preparation for R ...*

On a PC with 6GB of RAM, R could only process 30MB of data in 3MB chunks. Rt needed 11 jobs or nodes to break down the data and merge the results…

… The same data prep in Voracity happens in *just one* sort-join-aggregate program (and I/O pass) that runs 45% faster than R in this small case.

IRI
The CoSort Company

*Voracity Analytic Option 6:*
Big SM Streams

Leverage advanced text and social media analytic engines with NLP and Kafka support



📈 Facebook sentiment by minute

J1 JUPITERONE

| | | Clinton | Trump |
|---|---|---|---|
| Facebook | Fox News | 9,036K | 8,188K |
| | Donald J. Trump | 2,575K | 3,784K |
| | CNN | 1,934K | 3,024K |
| | Hillary Clinton | 1,677K | 2,685K |
| | NBC News | 917K | 1,535K |
| | Yahoo News | 798K | 1,146K |
| | The New York Times | 720K | 954K |
| | Washington Post | 486K | 809K |
| | The Huffington Post | 367K | 910K |
| Twitter | | 48,720K | 99,317K |

value
367K    99,317K

Powered by **JupiterOne**

— Clinton  — Trump

**DISCOVER**   **INTEGRATE**   **MIGRATE**   **GOVERN**   **ANALYZE**

## Data Curation

### Profile & Acquire

Discover and extract data and metadata in disparate sources. Define custom structures, mask formats, and build test data.

### Cleanse & Unify

Filter, enrich, scrub and standardize data in multiple sources. Select, fuzzy-search, and merge reference data into master tables and values.

### Process & Provide

Integrate, migrate, govern, and analyze data in the same job and I/O pass. Visualize and feed test or real targets in any format.

### Protect & Audit

De-ID data at the field level as you acquire, transform, report, or franchise. Encrypt, hash, pseudonymize, redact, tokenize, etc.

### Express & Predict

Aggregate, cross-calc, and format data in detail, summary and trend reports, or, hand-off results to your analytic tool or BIRT charts in memory.

### Convert & Replicate

Migrate legacy databases, or files and data types -- or specify new target record layouts -- in copies, or subsets, of data in any format or schema.

### Publish & Share

Federate, save, or populate multiple targets at once, and connect to them and their metadata in secure repositories for change tracking, etc.

IRI
The CoSort Company

# Retail

**Voracity Uses**

## Micro-target customers

Use Voracity to <u>segment</u> purchase groups for targeted marketing, and to create holistic, <u>unified</u> views of each customer that help you customize service and build loyalty.

## Leverage Consumer Psychology

Use Voracity to <u>integrate</u> consumer behavior and sentiment data against seasonal, regional, weather, and other factors, and mine it with <u>regression analyses</u> that reveal trends.

## Price Smarter

Use Voracity to integrate preference and pricing data from retail data brokers, public data, your own pricing history, and competitive research.

**IRI**
The CoSort Company

# BFSI

**Voracity Uses**

## Assess Credit Risk

Use CoSort and Hadoop engines in Voracity to blend traditional credit data with sources like utility bill and rental payments to improve score accuracy, facilitate lending, marketing, etc.

## Optimize Loan Performance

Use Voracity to blend and prepare internal and external data points (borrower history, industry repayment stats, social/market forces, etc.) for visual analytics on risk factors vs. loan rates.

## Expose Insurance Fraud

Use Voracity to rapidly sort, filter, and expose claim data outside normal parameters to identify suspicious behavior, and feed it to visualization and notification apps in the same IDE.

**IRI**
The CoSort Company

# Healthcare

## Voracity Uses

### Improve Treatment Outcomes

Flow IoT data through slowly changing dimension or change data capture processes in Voracity to compare patient data with diagnostic values to spot, alert, and correct for abnormalities.

### Individualize Drug Therapies

Rapidly integrate genetic data into single-node-type networks, gene-set libraries, and bi-partite graphs to help reveal new relationships between patient genes, drugs and phenotypes.

### See the Whole Patient

Use Voracity' search, join, consolidate, and masking features to unify and de-identify patient information from family, provider, demographic, diagnostic and treatment data silos.

# Energy & Transport

Voracity Uses

## Conserve & Troubleshoot

Use the IoT edge aggregation and hub analytics in Voracity on smart meter and thermostat data to identify peak uses, or on grid sensor and weather data to re-route power, inspect, repair, etc.

## Improve Traffic Flow

Combine data from street cameras and sensors, cell phone apps and weather data in Voracity and feed it directly into BIRT or BIRT-connected Integeo geospatial reports to warn drivers.

## Optimize Fleet Performance

Use IoT analytic and alerting features in Voracity to predict and prevent equipment failures, and its DW/BI prowess against historic O&D and pricing data to maximize passenger revenues.

IRI
The CoSort Company

# Telco & Media

## Voracity Uses

### Monetize Calls & Clicks

Use Voracity to process CDRs and clickstream data for billing and analytics, and to sell that data to marketing affiliates and others who can permissibly use it.

### Anticipate Spending Trends

Use Voracity to extract string and pattern-matching values from social data from Hubspot, etc., and munge it with transaction and demographic data to identify and predict content preferences.

### Throttling & Enforcement

Use Voracity to identify excessive bandwidth usage or illegal behavior from network traffic and web logs, and tie it to analytic and notification mechanisms in the same IDE.

**Voracity Uses**

Reliance Communications (RC) is broadband and telco company in india with 110M subscribers. To meet daily SLAs in billing and analytics for wireless (mobile) and global (landline) segments, RC must process and report on hundreds of millions of call detail records (CDRs) every day.

RC uses 64-bit Solaris servers and Oracle. The CDRs come from binary switch data mediated into flat files that the CoSort engine in Voracity transforms *before* DataStage ETL & BOBJ reports.

*"Prior pilots failed from slow and inaccurate results, and SLAs were missed as call volume grew. After Voracity jobs transformed flat files in the 60GB range, the processing bottleneck disappeared, and our analytic results were always accurate."*

**RELIANCE**
**Communications**
Anil Dhirubhai Ambani Group

**IRI**
The CoSort Company

## Voracity Uses

DataBase Technologies (DBT) in Parsippany, NJ builds and maintains VLDB CRMs for ADP, Verizon, Merrill Lynch, Seagrams, and Universal Studios.

DBT integrates 350M transaction records per day, joining them to files up to 100M rows each, and accumulating the data over time for analysis. Their first 350GB dataset took over two days to load, so it had to be pre-sorted.

*"It's fun to watch the system performance monitor and see all those processors working in the high 90 percentages and the disks utilizing the fast data rates you pay for."*

Voracity filter, sort, and join operations, were 10x faster than those in MS SQL Server …. 9.5 minutes versus 98 @350GB.
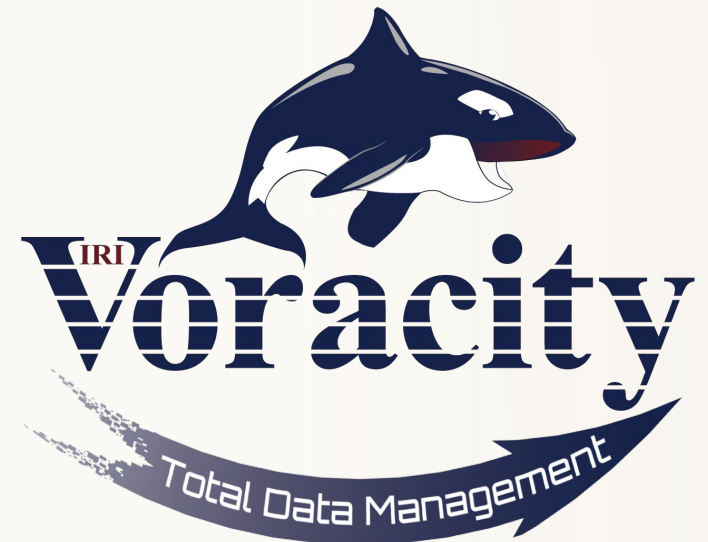
**DATABASE**
TECHNOLOGIES

**IRI**
The CoSort Company

Global Big Data Conference

Learn and Share

IRI.com    IRI blog

IRI Voracity Data Management Group on LinkedIn

www.globalbigdataconference.com