ORACLE®

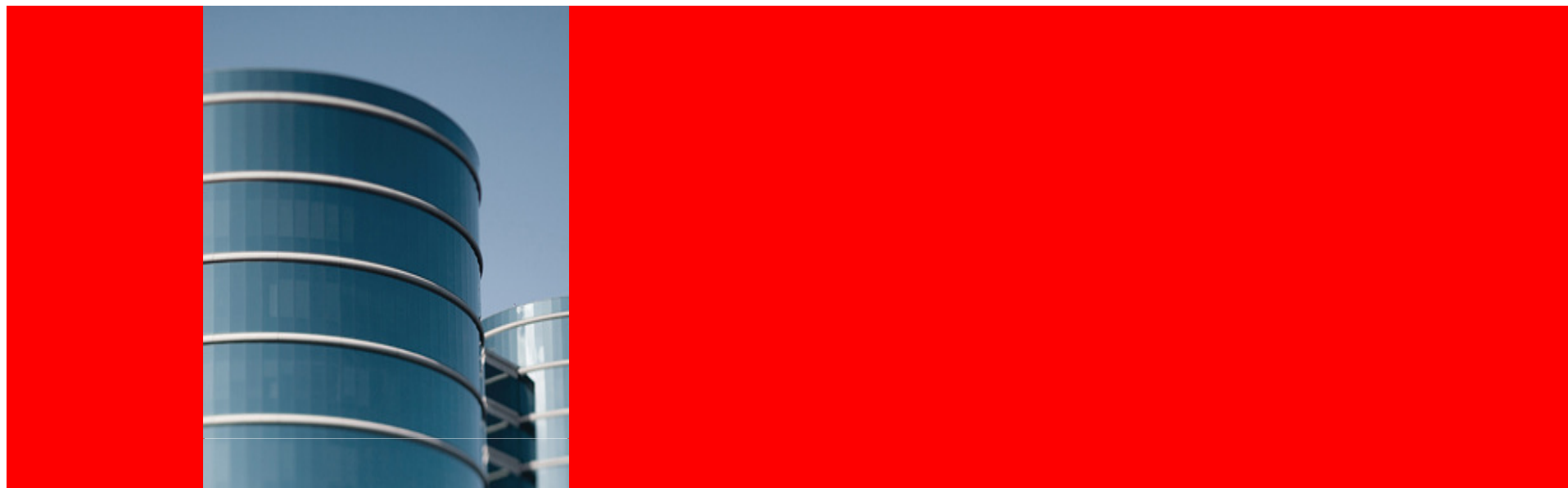**ORACLE**®

# Get more value from your DWH

Ugo Pollio – Sales Consulting and Business Development
Oracle EMEA Data Integration Solutions

# Program Agenda

- The challenge of Quality of Data
- Oracle Data Integration Solution
- Get More Value from your DWH
- Q&A

# Data Quality issues are on your daily business

# An example in DWHs…

# Oracle Data Integration Solution
## Best-in-class Heterogeneous Platform for Data Integration

| Oracle Applications | Custom Applications | MDM Applications | Business Intelligence | Activity Monitoring | SOA Platforms |
|---|---|---|---|---|---|

### Comprehensive Data Integration Solution

**SOA Abstraction Layer**

| Process Manager | Service Bus | Data Services | Data Federation |
|---|---|---|---|

| Oracle Data Integrator | Oracle GoldenGate | Oracle Enterprise Data Quality |
|---|---|---|
| ELT/ETL | Real-time Data | Data Profiling |
| Data Transformation | Log-based CDC | Data Parsing |
| Bulk Data Movement | Bi-directional Replication | Data Cleansing |
| Data Lineage | Data Verification | Match and Merge |

Storage | Data Warehouse/ Data Mart | OLTP System | OLAP Cube | Flat Files | Web 2.0 | Web and Event Services, SOA

**ORACLE**

# ODI Saves Money

## E-LT Runs on Existing Servers with Shared Administration

### Typical: Separate ETL Server

- Proprietary ETL Engine
- Expensive Manual Parallel Tuning
- High Costs for Standalone Server

### ODI: No New Servers

- **Lower Cost:** Leverage Compute Resources & Partition Workload efficiently
- **Efficient:** Exploits Database Optimizer
- **Fast:** Exploits Native Bulk Load & Other Database Interfaces
- **Scalable:** Scales as you add Processors to Source or Target
- **Manageability:** unified Enterprise Manager

### Benefits

- Better Hardware Leverage
- Easier to Manage & Lower Cost
- Simple Tuning & Linear Scalability

**Next Generation Architecture**

# E-LT

Transform

**E**xtract   **L**oad   **T**ransform

**Conventional ETL Architecture**

**E**xtract   **T**ransform   **L**oad

# ODI is Simpler

## Speed Project Delivery and Time to Market with ODI

- Development Productivity
  - 40% Efficiency Gains

- Environment Setup (ex: BI Apps)
  - 33-50% Less Complex



**ODI Declarative Design**

Define **What** You Want → Automatically Generate Dataflow

Define **How**: Built-in Templates

**Conventional ETL Design**

| Number of Setup Steps | 7 |
|---|---|
| Number of Servers | 1 |
| Number of connections | 3 |

| Number of Setup Steps | 10 |
|---|---|
| Number of Servers | 3 |
| Number of connections | 7 |

# ODI is Faster

## Up to 7TB per hour of real world data loading and complex transformations

### Over 7 TB(60 Billion Rows) per Hour

Legend: GB / Hour, Rows / Hour

Categories (x-axis):
Exadata V1 Qtr Rack, Exadata V1 Half Rack*, Exadata V1 Full Rack*, Exadata V2 X2-2 Qtr Rack, Exadata V2 X2-2 Half Rack*, Exadata V2 X2-2 Full Rack*, Exadata V2 X2-8*

### Exadata X2-2 Qtr Rack - 1 TB in 39 Minutes

| | 100 GB | 300 GB | 1,000 GB |
|---|---|---|---|
| ODI ELT w/ Exadata | 03:57 | 11:50 | 39:27 |
| Conventional ETL | 09:26 | 27:09 | 01:30:17 |

## ODI ELT (on Exadata/any DW)

- ODI scales with Exadata
  - Loads increase linearly as Exadata scales
- ODI runs on Exadata – no ETL hardware required
  - No new hardware required as data sets grow
  - ODI processes used only during integration runs
  - Exadata continually available for OLTP, BI, DW, etc
- Common administration, monitoring and management
- All the benefits of rapid tools-based ETL development

## Conventional ETL

- As data sets grow, more hardware (**$$**) needed to scale
- ETL parallel optimization and design (**$$$**) is heavily dependent on resources available to the ETL environment
  - Sources, integrations, targets must be designed to match processing power of ETL environment
  - Source flat files split to match # of ETL engine CPU's
  - Integration grid setup appropriately to match # of ETL engine CPU's
  - Target partitions, table spaces to match # of ETL engine CPU's
- ETL engine hardware resources only used for ETL
  - Cannot be utilized for OLTP, BI, DW, etc.
- Hardware not co located, multiple vendors
- Different management, monitoring and administration from database and BI infrastructure (**$$**)

# Product Architecture
## for Enterprise Scale Deployment

**Desktop**

**JVM**

**ODI Studio**

Java EE Application

ODI SDK

- Designer
- Operator
- Topology
- Security

**WebLogic 11g / Application Server**

FMW Console

ODI Plug-in

Servlet Container

Java EE Application

ODI SDK

ODI Console

Runtime WS

Java EE Agent

Web Service Container

Public WS

Data Services

Data Sources Connection Pool

**Repositories**

ODI Master Repository

ODI Work Repository

**Sources and Targets**

Legacy

Files / XML

DBMS

Applications
ERP/CRM/PLM/SCM

DW / BI / EPM

**JVM**

Runtime WS

Standalone Agent

# Unique Qualities of Enterprise Data Quality

**Integrated DQ Solution**

- Seamless integration of all core DQ capabilities
- Engineered for business users
- Integrated team collaboration and management

**Modern Architecture**

- Easy to configure and integrate 'DQ Services'
- Modern, open architecture (Java, SOA, etc.)
- Collaborative, multi-user project support

**Designed as a Platform**

- Enables innovation and reuse
- Enables delivery of complete business solutions
- Allows partners opportunity to showcase their expertise and specialization

ORACLE®

# Oracle EDQ for Customer Data

**Fastest Time To Value**

- Easy to use - Single screen combines all key DQ functions
- Simple graphical UI, no coding required
- Collaborative environment with integrated issue management
- Designed for business user

**Most Tunable Solution**

- Fully tunable rules
- Closed-loop rule building with instant feedback
- No 'black boxes'
- Easy to extend and re-use

**Lowest Total Cost of Ownership**

- Integrated solution
- Rapid integration and deployment
- Low maintenance overheads
- Scalable performance

ORACLE®

# Oracle EDQ for Customer Data

## DQ-Based Solutions

Business Solutions

Domain Knowledge

- Customer-delivered
- Partner-delivered
- Oracle-delivered

## Enterprise Data Quality

| | |
|---|---|
| Profiling | Analysis |
| Parsing | Standardization |
| Match/Merge | Reporting |
| Case Management | |

### Solutions
- Pre-configured business solutions and domain knowledge
- Can be delivered by
  - Customers
  - Partners
  - Oracle
- Solutions can be components, processes, domain knowledge, up to complete applications

### Data Quality Platform
- End-to-end enterprise data quality solution spanning both customer and product data quality
- Highly configurable to match specific business needs
- Case management tools for tracking and web-based KPI reporting for increased productivity

ORACLE®

# Product Architecture



- All Java Server (Stateless)
- Java Webstart Client Applications
- Fully integrated with a single repository and UI
- Batch and Real-time Execution
- Connects to virtually any source/target of data
- Platform Independent

# ODI & OEDQ – ELT+CM

**Source**

**Target**

**Staging**

**ODI** **E**xtract Agent

**ODI** **L**oad

**ODI** **T**ransform

**OEDQ** **C**leanse & **M**atch

**E**xtract > **L**oad > **T**ransform > **C**leanse > **M**atch

ORACLE

# Integrated Data Quality with ODI
## Oracle Enterprise Data Quality Runtime with Data Integrator



- **Best of breed Quality**
  - Proven, scalable DQ engine
  - Rich capabilities for cleansing, standardization, validation, match and merge
  - Extensible by customers

- **Out-of-box integration**
  - ODI integrates with Quality functions via pre-built ODI OpenTool
  - Drag and drop graphical icon for inserting DQ flows into ODI

# Oracle Enterprise Data Quality – Standardization

Name: Dr Ellen Van Der Heijde

Name: Mr RJ & Mrs FB MacDonald

Name: Jalila Abdul-Alim (Do Not Call)

Title: Dr
First: Ellen
Last: Van Der Heijde
Gender: Female

Title: Mr
First: R
Middle: J
Last: MacDonald
Gender: Male

Title: Mrs
First: F
Middle: B
Last: MacDonald
Gender: Female

First: Jalila
Last: Abdul-Alim
Gender: Female
Note: Do Not Call

- Standardize, Transform and Parse
- Split names and name elements
- Identify individuals and businesses
- Derive additional attributes

ORACLE®

# Oracle Enterprise Data Quality - Matching

First: Bob
Last: Fulmar
Gender: Male
Email: chem291_rjf@barker.edu

Title: Mr
First: Robert
Last: Fulmar
Gender: Male
DoB: 12/05/1978
Phone: 555-120-1329
Address:
9405 Main St
Fairfax
Virginia
22030

Title: Dr
First: R
Last: Fulmer
DoB: 01/01/1978
Email: chem291_rjf@barker.edu
Address:
9407 Main Street
Fairfax
VA
22031-4001

Title: Dr
First: Robert
Last: Fulmar
Gender: Male
DoB: 12/05/1978
Email: chem291_rjf@barker.edu
Phone: 555-120-1329
Address:
9407 Main St
Fairfax
VA
22031-4001

- Match & Merge data from disparate sources
- Create 'best' record based on survivorship rules

ORACLE®

# Built for the Business User

- Short learning curve & time-to-value

- Solution for owners of the business problem

- Integrated team collaboration

# Oracle EDQ for Product Data

**Built from the ground up for product data**

- Handles the variability of product data – structure, standards, categories
- Use Customer's data to build references

**Ability to govern a largely un-governed process**

- Stewardship, oversight and remediation combined in a single interface
- Optimal combination of automation and remediation

ORACLE®

# The Product Data Problem – Unstructured & Non-Standard

## What is this?



10hp motor 115V Yoke mount

MOT-10,115V, 48YZ,YOKE

mtr, ac(115) 10 horsepower 115volts

This 10hp yoke mounted motor is rated for 115V with a 5 year warranty

10 Caballos, Motor, 115 Voltios

TEAO HP = 10.0 1725RPM 115V 48YZ YOKE  MTR

Motor, TEAO, 1725 RPM, 48YZ, 15 Voltios, Montaje de Yugo, hp = 10

| Item | Motor |
|---|---|
| Classification | 26101600 |
| Power | 10 horsepower |
| Voltage | 115 |
| Mounting | Yoke |

Product data is much more variable and unpredictable than other data types

ORACLE®

# An example in DWHs…

## Data Quality Control

Asked by Jean-Pierre Paisley | posted 10 days ago | Replies (4)

Hie Guys,

We have set up a cognos platform for our reporting for my executives, which are run overnight. We get data from various points totaling more than 130, each sending more than 30000 records daily. The problem however that is now arising is that usually reports are with some points or modules being offline or down, thus will not have sent their data. We are thus having a situation where we will have inaccurate data usually at time of reporting and sometimes if a key point is offline this results in big variances. They are thus now losing confidence in the BI platform.

Does anyone have any ideas on how we can ensure data quality control, probably ways we can flag to show that there the data is incomplete?

ORACLE®

- Get Information from data is always a challenge:
  - Because of lack of standards, or different standards across different sources
  - Because of missing values and typos
  - Robust Matching features are the only that solve duplicate issues
- Business executors always need trustable data:
  - For taking the right decision
  - To get business insights and further develop the business
  - Discover potential gaps
- Customer believe in Oracle DQ strategies because:
  - Best of breed technology
  - Most flexible and customizable as per customer's needs

ORACLE

# Next Gen Data Warehouse
## Change data into valuable Information

| Acct Name | Closed Rev | Profitability | Share of Total Cust Rev |
|---|---|---|---|
| Berkeley Asset Management | 5,346,500 | 4,233,584 | 18% |
| First Bank of CA | 2,450,000 | 1,887,857 | 8% |
| A. K. Parker Distribution | 2,404,000 | 938,716 | 8% |
| Columbia Bank | 1,564,000 | 1,564,000 | 5% |
| Collins Pharmaceutical | 1,300,000 | 954,979 | 4% |
| A. K. Parker Distrib | 1,006,000 | 500,242 | 3% |
| AK Parker Distribution | 592,000 | 240,585 | 2% |
| Parker Distribution | 150,000 | 54,320 | 1% |
| A. K. Parker Dist | 25,000 | 25,000 | 0% |
| Grand Total | 14,837,500 | 10,399,282 | 50% |

⚠ CAUTION

Do not trust this information!

- **Profiling**: • Investigate, Analyze, Audit
- **Cleansing**: • Standardize, Enrich, Deduplicate
- **Control**: • Govern over time

| Acct Name | Closed Rev | Profitability | Share of Total Cust Rev |
|---|---|---|---|
| Berkeley Asset Management | 5,346,500 | 4,233,584 | 18% |
| A. K. Parker Distribution | 4,177,000 | 1,758,863 | 14% |
| First Bank of CA | 2,450,000 | 1,887,857 | 8% |
| Columbia Bank | 1,564,000 | 1,564,000 | 5% |
| Collins Pharmaceutical | 1,300,000 | 954,979 | 4% |
| Grand Total | 14,837,500 | 10,399,282 | 50% |

- **The Business Issue**
  - BI Reports are not trustable, because of the state of source data

- **Reduce risks**
  - Improve data quality by integrating cleansing as part of the process
  - Eliminate data redundancies

- **Improve Business Insights**
  - Improved business insight with improved data quality
  - Better profiling of data to eliminate gaps in insight

ORACLE

# Mission-Critical Systems and Batch Processing

Too Much Data, Not Enough Time



**What time of your day is your business *NOT* at it's peak?**

# E-LTQ In-line Predictive Quality
## What if you go longer than the batch window?

**Sources**

**ODS / Stage Area**

Emp

Dept

Geo

Sales

**Data Mart**

DIM

DIM

FACT

DIM

DIM

**REPORTING**

11:00pm    **B A T C H   T I M E   W I N D O W**    8:00am

Data are loaded, transformed during a batch time window, before users and applications get access.
This windows easily becomes a challenge because of:
- data volume increasing
- only static controls can be applied on the flow, eventually discarding bad data. Discarded data potentially generate inconsistencies on the final target.

ORACLE

# E-LTQ In-line Predictive Quality
## What if you discard or, worse, you load, bad data?



**Sources**

**ODS / Stage Area**

Emp

Dept

**Data Mart**

DIM

DIM

FACT

DIM

DIM

**REPORTING**

ACCESS DENIED

11:00pm    **BATCH TIME WINDOW**    8:00am

If bad data are processed without cleansing or discarded , the Data mart cannot be accessed until bad data are fixed. This usually is done by IT operations, with lot of efforts, Users & Manager pressing for restoring data ASAP.

ORACLE

# E-LTQ In-line Predictive Quality
## Check your data dynamically, ensure Quality

**Sources**

**ODS / Stage Area**

**Data Mart**

**E-LT ODI flow in parallel**

EMP  DEPT

DIM     DIM

FACT

DIM     DIM

**REPORTING**

**?**

**GO!**

A > ?
Sum (B) = ?
$C_{tn-1} - C_{tn} < ?$

**EDQ in parallel**

A > 87%
Sum (B) = 1002
$C_{tn-1} - C_{tn} < 2.3\%$

11:00pm          **B A T C H   T I M E   W I N D O W**          8:00am

**Check while you transform & Load your
Business Thresholds A > 80%, sum(B) = 1000, $C_{tn-1} - C_{tn} < 3\%$**

ORACLE

# Data Quality **Firewall**
## profile, repair, check, alert, report



**Source**

Cobol copybooks

Databases

TXT files

**EDQ+ ODI**

**Database**

**OBI-EE**

Oracle Data Quality for Data Integrator

Global Data Router — Transformer — Parser — Postal Matcher — Relationship Linker

**Discarded Record**

**Human Workflow**

- DWH, improving reliability and quality
- New ERP/CRM installation, and legacy data integration
- Master Data Management projects
- Data synchronization projects

ORACLE

# Business Intelligence
## Real Time Data Warehouse



Source 1
EMP  DEPT
On-Disk Logs

Low Impact Real-Time CDC

Source 2
EMP  DEPT
On-Disk Logs

Continuous ELT&DQ

DIM  DIM
FACT
EMP  DEPT  DIM  DIM

ODS Schema  DW Schema

- **Solution**
  - Log-based capture of database transactions from source systems
  - Load to target with sub-second latency
  - Transformation performed on the database using E-LT in continuous mini-batches

- **Benefits**
  - No resource / performance impact to OLTP
  - Fresh data available for better decision making
  - Get double-duty from database investment by using it for transformations
  - No batch windows necessary – key for global businesses

**ORACLE**

# Key Capabilities

Interactive exploration of data, identifying distribution and outlying values with drill-downs

Identify and quantify issues in data

| TITLE | Count | % |
|---|---|---|
| Mr | 816 | 40.8 |
| Ms | 468 | 23.4 |
| Mrs | 309 | 15.4 |
| Miss | 251 | 12.5 |
| Dr | 15 | 0.7 |
| Rev | 1 | <0.1 |
| Prof. | 1 | <0.1 |
| Col. | 1 | <0.1 |

| Input field | Without data | Singleton | Duplications | Distinct values | Comment |
|---|---|---|---|---|---|
| CU_NO | 1 | 1997 | 3 | 1998 | Potentially damaged key; Investigate nulls; Investigate duplicates |
| CU_ACCOUNT | 1 | 2000 | 0 | 2000 | Potentially damaged key; Investigate nulls |
| TITLE | 139 | 3 | 1859 | 8 | |
| NAME | 1 | 1980 | 20 | 1990 | Potentially damaged key; Investigate nulls; Investigate duplicates |
| GENDER | 148 | 0 | 1853 | 2 | |
| BUSINESS | 331 | 1629 | 41 | 1649 | Investigate duplicates |
| ADDRESS1 | 2 | 1926 | 73 | 1954 | Potentially damaged key; Investigate nulls; Investigate duplicates |
| ADDRESS2 | 80 | 554 | 1367 | 839 | Investigate nulls |
| ADDRESS3 | 969 | 278 | 754 | 379 | |
| POSTCODE | 239 | 1604 | 158 | 1672 | |
| AREA_CODE | 117 | 64 | 1820 | 270 | |
| TEL_NO | 7 | 1875 | 119 | 1934 | Potentially damaged key; Investigate nulls |
| EMAIL | 65 | 1904 | 32 | 1920 | Potentially damaged key; Investigate nulls; Investigate duplicates |
| ACC_MGR | 5 | 0 | 1996 | 30 | Investigate nulls |
| DT_PURCHASED | 3 | 1090 | 908 | 1499 | Investigate nulls |
| DT_ACC_OPEN | 3 | 1093 | 905 | 1500 | Investigate nulls |
| DT_LAST_PAYMENT | 4 | 1026 | 971 | 1425 | Investigate nulls |
| DT_LAST_PO_RAISED | 3 | 1003 | 995 | 1433 | Investigate nulls |
| BALANCE | 2 | 7 | 1992 | 10 | Investigate nulls |

# Oracle Enterprise Data Quality – Audit



- Validate data against business rules
- Publish results to data quality dashboard

# Key Features

- Fully Unicode Compliant

| Origin_nat | Name_nat |
|---|---|
| 日本 | テクス・テクサン |
| 조선 민주주의 인민 공화국 | 이설희 |
| 대한민국 | 안성기 |
| 대한민국 | 심은하 |
| 조선 | 솅ㅈㅗ▵ (세종대왕) |
| Repubblika ta' Malta | Trevor Żahra |
| Norge | Tor Åge Bringsværd |
| Noreg | Herbjørn Sørebø |
| پاکستان | نصرت فتح علی خان |
| Perú | Nicómedes Santa Cruz |
| Polska | Lech Wałęsa |
| Portugal | Amália Rodrigues |
| Puerto Rico | Olga Tañón |
| Rōma | Pūblius Cornēlius Scīpiō Africānus |
| Россия | Михаил Горбачёв |
| Россия | Борис Гребенщиков |
| רוסלאַנד | שלום עליכם |
| Sápmi | Áillohaš |
| Slovensko | Ľudovít Štúr |
| Slovenija | Frane Milčinski - Ježek |
| Sverige | Björn Borg |
| Συρακούσα | Ἀρχιμήδης |
| تاجیکستان | صدر الدین عینی |

| Character | Decimal | Hex | Total |
|---|---|---|---|
| 原 | #21407 | #x539f | 1 |
| č | #269 | #x10d | 1 |
| ث | #1579 | #x62b | 1 |
| š | #352 | #x160 | 1 |
| 林 | #26519 | #x6797 | 1 |
| ū | #363 | #x16b | 1 |
| µ | #956 | #x3bc | 1 |
| ד | #1492 | #x5d4 | 1 |
| ú | #250 | #xfa | 1 |
| ο | #959 | #x3bf | 1 |
| ג | #1490 | #x5d2 | 1 |
| ľ | #317 | #x13d | 1 |
| 기 | #44592 | #xae30 | 1 |
| א | #1488 | #x5d0 | 1 |
| λ | #955 | #x3bb | 1 |
| δ | #948 | #x3b4 | 1 |
| ן | #1503 | #x5df | 1 |
| ぐ | #12368 | #x3050 | 1 |
| ô | #244 | #xf4 | 1 |
| η | #951 | #x3b7 | 1 |
| 鷗 | #40407 | #x9dd7 | 1 |
| כ | #1499 | #x5db | 1 |
| ī | #299 | #x12b | 1 |

ORACLE

# Key Features

- Comprehensive DQ Functionality with a single UI and Repository

# Key Features

- Provided Extensions for Customer Data and Locales
- Highly Extensible



Tool Palette - Customer Data

- Add Gender
- Advanced Email Check
- Country Code from Country
- Country from City
- Country from Nationality
- Create Identifier
- Extract Building Identifier
- Extract Town
- Extract Zipcode
- GeoNames Check City
- GeoNames Country Codes from City
- Match Entities
- Match Households
- Match Individuals (Name, Address, DoB)
- Standardize Entity Names

Tool Palette - United States

- US Address Parser
- US Derive State from City Name
- US Standardize State
- US Standardize Street Words
- US Standardize Towns

ORACLE®

# **Transformation** – **Data Improvement**

- Fully configurable data transformation rules
- Operates in both Batch and Real-Time
- Full control over data updates
- Original data always preserved (and all steps in between)
- Source data may either be staged and processed or 'streamed' through the process



| | |
|---|---|
| **Use profiling results to create your own data improvement rules** | **Use provided processors for common tasks such as address standardization** |

- Designed for business users
- Flexible matching engine for any data with many comparison algorithms
- Provided template match processors for individual, entity and address matching
- Easy reuse of configured match processors
- Fully configurable outputs (Links, Groups, Master and Slaves, Best Record)
- Operates in both Batch and Real-Time
- See Match Essentials deck for more information on Matching

Advanced Options | Assign Relationship Review
Review Results | Assign Merged Review
Configure Bulk Review Rules | View Match Statistics
Delete Realtime Review Results | Delete Manual Decisions

**Match Entities**

Input    Identify    Cluster    Match    Merge

| Rule | Priority | Name WMP | Name contains | Name CMP | Name word match | Name first word | Name initials | Add1 s/w | Town ed | Postcode e | Decision |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Rule out records with no name words matching | 0 | * | * | * | None | * | * | * | * | * | NOMATCH |
| Name stand, Address | 95 | 80+ | * | * | * | * | * | true | * | Exact | MATCH |
| Name stand close, Address | 93 | 60-79 | * | * | * | * | * | true | * | Exact | MATCH |
| Name contains, Address | 92 | * | true | * | * | * | * | true | * | Exact | MATCH |
| Name chars, Address | 91 | * | * | 80+ | * | * | * | true | * | Exact | MATCH |
| Name stand possible, Address | 90 | 50-59 | * | * | * | * | * | true | * | Exact | MATCH |
| Name stand, Postcode | 88 | 80+ | * | * | * | * | * | * | * | Exact | MATCH |
| Name stand close, Postcode | 86 | 60-79 | * | * | * | * | * | * | * | Exact | MATCH |
| Name contains, Postcode | 84 | * | true | * | * | * | * | * | * | Exact | MATCH |
| Name chars, Postcode | 82 | * | * | 80+ | * | * | * | * | * | Exact | MATCH |
| Name first word, Address | 81 | * | * | * | * | Exact | * | true | * | Exact | REVIEW |
| Name stand possible, Postcode | 80 | 50-59 | * | * | * | * | * | * | * | Exact | REVIEW |
| Name first word, Postcode | 78 | * | * | * | * | Exact | * | * | * | Exact | REVIEW |
| Name, Address1 and Town | 77 | 80+ | * | * | * | * | * | true | 0-1 | * | REVIEW |
| Name, no Address | 75 | 80+ | * | * | * | * | * | no data | * | no data | REVIEW |
| Name, Address1 only | 70 | 80+ | * | * | * | * | * | true | * | * | REVIEW |
| Address only | 65 | * | * | * | * | * | * | true | * | Exact | REVIEW |
| Name only | 50 | 80+ | * | * | * | * | * | * | * | * | REVIEW |

ORACLE

# Reporting

- Highly flexible reporting interface
  - Export any Results views automatically to database/file
  - 1-click export of results to Excel from the Director client
- Dashboard reporting provides stakeholder view of Data Quality KPIs with trend analysis
- Example reports
  - Automatic Matches / Review Matches / Non-Matching Records
  - Match Group Size Report
  - Match Rule Report
  - Data Validity Report
  - Profiling Report
  - Etc.

ORACLE

# Reporting - Immediate drilldown reporting in Director

- Results Browser

# **Reporting** - Collect important results into Results Books



**Results Browser - Key Results**

Viewing all 17 records

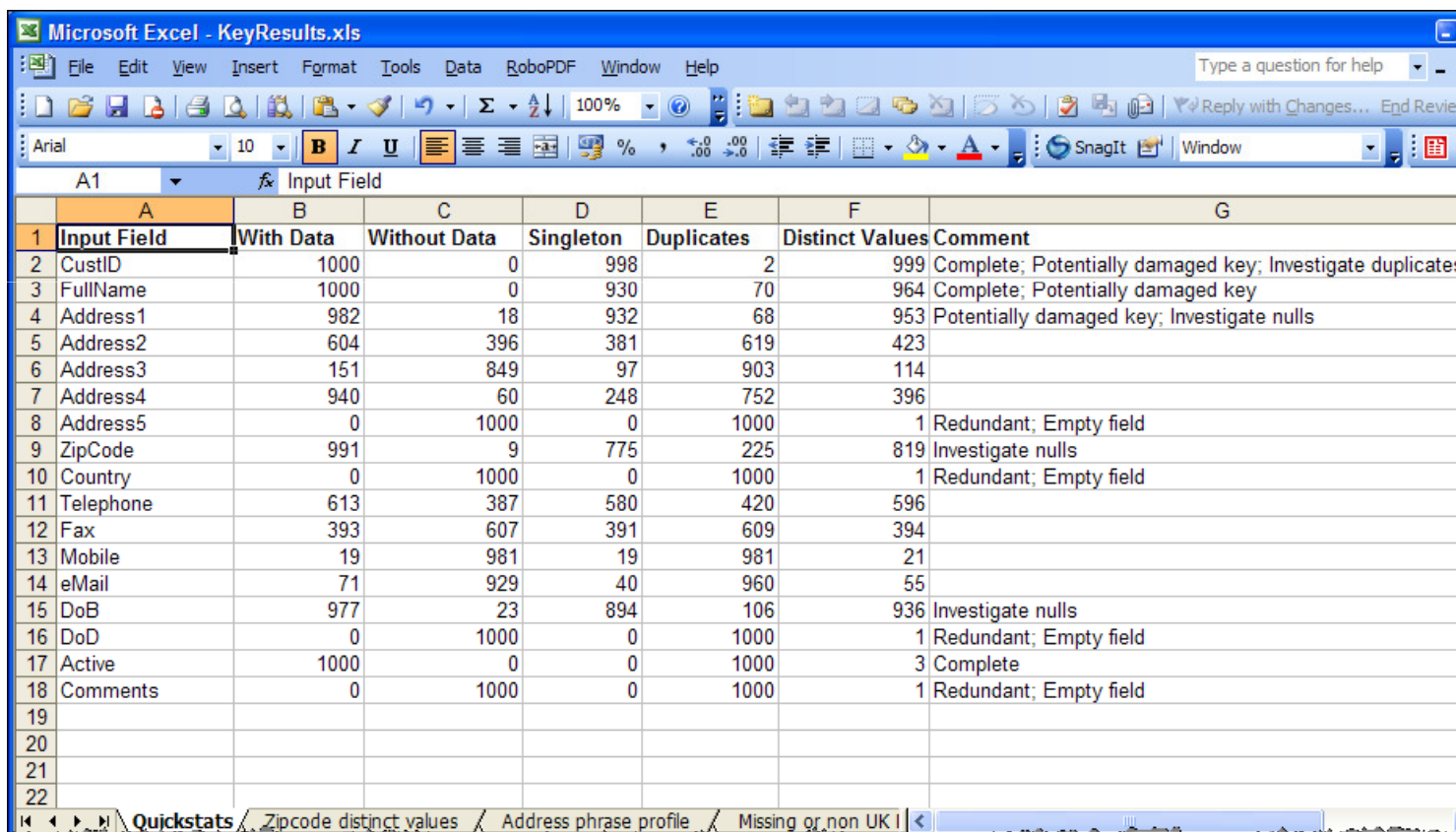| Input Field | With Data | | Without Data | | Singleton | | Duplicates | | Distinct Values | | Comment |
|---|---|---|---|---|---|---|---|---|---|---|---|
| CustID | 1000 | 100% | 0 | 0% | 998 | 99.8% | 2 | 0.2% | 999 | 99.9% | Complete; Potentially damaged key; Investigate duplicates |
| FullName | 1000 | 100% | 0 | 0% | 930 | 93.0% | 70 | 7.0% | 964 | 96.4% | Complete; Potentially damaged key |
| Address1 | 982 | 98.2% | 18 | 1.8% | 932 | 93.2% | 68 | 6.8% | 953 | 95.3% | Potentially damaged key; Investigate nulls |
| Address2 | 604 | 60.4% | 396 | 39.6% | 381 | 38.1% | 619 | 61.9% | 423 | 42.3% | |
| Address3 | 151 | 15.1% | 849 | 84.9% | 97 | 9.7% | 903 | 90.3% | 114 | 11.4% | |
| Address4 | 940 | 94.0% | 60 | 6.0% | 248 | 24.8% | 752 | 75.2% | 396 | 39.6% | |
| Address5 | 0 | 0% | 1000 | 100% | 0 | 0% | 1000 | 100% | 1 | 0.1% | Redundant; Empty field |
| ZipCode | 991 | 99.1% | 9 | 0.9% | 775 | 77.5% | 225 | 22.5% | 819 | 81.9% | Investigate nulls |
| Country | 0 | 0% | 1000 | 100% | 0 | 0% | 1000 | 100% | 1 | 0.1% | Redundant; Empty field |
| Telephone | 613 | 61.3% | 387 | 38.7% | 580 | 58.0% | 420 | 42.0% | 596 | 59.6% | |
| Fax | 393 | 39.3% | 607 | 60.7% | 391 | 39.1% | 609 | 60.9% | 394 | 39.4% | |
| Mobile | 19 | 1.9% | 981 | 98.1% | 19 | 1.9% | 981 | 98.1% | 21 | 2.1% | |
| eMail | 71 | 7.1% | 929 | 92.9% | 40 | 4.0% | 960 | 96.0% | 55 | 5.5% | |
| DoB | 977 | 97.7% | 23 | 2.3% | 894 | 89.4% | 106 | 10.6% | 936 | 93.6% | Investigate nulls |
| DoD | 0 | 0% | 1000 | 100% | 0 | 0% | 1000 | 100% | 1 | 0.1% | Redundant; Empty field |
| Active | 1000 | 100% | 0 | 0% | 0 | 0% | 1000 | 100% | 3 | 0.3% | Complete |
| Comments | 0 | 0% | 1000 | 100% | 0 | 0% | 1000 | 100% | 1 | 0.1% | Redundant; Empty field |

Quickstats | Zipcode distinct values | Address phrase profile | Missing or non UK Postcodes | Telephone number patterns | Building identifier extraction

Records where building id could not be extracted | Customer deduplication | Customer deduplication rules

ORACLE

# Reporting - Output Results Books

Export Results Books as part of an automated job
1-click Export of a Results Book to Excel

# Arabic – Name Transcription Approach

- For 98+% of individual names, transcription occurs directly to the IC form of the name using a large database of Arabic names

- For the remaining names, a custom dictionary is used

- If a name is still unrecognized (<1%) it is transliterated using character-based transliteration and flagged as an exception

- Easy to add transcriptions for exceptions to the custom dictionary

- It is also possible to override the transcription for a specific names

| namesurname | dnGivenNames | dnFamilyName | dnFullName |
|---|---|---|---|
| محمد أحمد إبراهيم | MUHAMMAD AHMED | IBRAHIM | MUHAMMAD AHMED IBRAHIM |
| محمد الزنيطى | MUHAMMAD | AL ZANATI | MUHAMMAD AL ZANATI |
| محمد بيت المال | MUHAMMAD BAYT | AL MAL | MUHAMMAD BAYT AL MAL |
| محمد حسين فضل الله | MUHAMMAD HUSSEIN FADL | | |
| عيسى عبدالكافى | ISSA ABDUL | | |
| محمد عبد الجواد | MUHAMMAD ABDUL | | |
| محمد على الحويج | MUHAMMAD ALI | | |
| محمد على الحويز | MUHAMMAD ALI | | |
| محمد محمود الحجازي | MUHAMMAD MAHMOUD | | |
| محمد مطوق مطوق | MUHAMMAD MATUQ | | |
| مشعان الجبوري | MUSHAN | | |
| محمد طاهر حموده سعيله | MUHAMMAD TAHIR HAMMUDAH | | |
| معتوق محمد معتوق | MATUQ MUHAMMAD | | |
| مفتلح محمد كوبح | MFTLH MUHAMMAD | | |
| هادي كوبر | HADI | | |
| معمر محمد القذافى | MAMAR MUHAMMAD | | |
| وئام وهاب | WIAM | | |

**Results Browser**

Job: Greek to Latin master

| Original Script Name | Original Script Name.Transliterated |
|---|---|
| Ευαγγελος Αντωνιου | Evangelos Antoniou |
| Ροης Σπυρος Πογιαντζῆς | Rois Spyros Pogiantzis |
| Λυκουργος Κυπριανου | Lykourgos Kyprianou |
| Ιωαννης Κατελουζος | Ioannis Katelouzos |
| Παναγιωτης Χαικαλης | Panagiotis Chaikalis |
| Δημητριος Κονδυλιος | Dimitrios Kondylios |
| Χρηστος Δημητριος Χατζοπουλος | Christos Dimitrios Chatzopoulos |
| Χρηστος Δημητριος Χατζοπουλος | Christos Dimitrios Chatzopoulos |
| Γεωργιος Αλεξανδρης | Georgios Alexandris |
| Γεωργιος Αλεξανδρης | Georgios Alexandris |
| Αναστασιος Βαβατσικλης | Anastasios Vavatsiklis |
| Κωνσταντινος Καπολλας | Konstantinos Kapollas |
| Κωνσταντινος Παπαναγιωτου | Konstantinos Papanagiotou |
| Κωνσταντινος Π... .νγιω.. .υ | Kon...antinos P... ...ist... |

# **Arabic** – **Name Matching Approach**

- All known Latin variant representations of an Arabic name are recognized in matching using a dictionary of 5m variants
- The 5m may be filtered to the most frequent representations only if required
- High confidence matching even where transcription standards may be very different
- Can match both Arabic to Arabic and Arabic to English
- Can match Arabic to other languages via comprehensive transliteration capabilities for other languages
- Wide variety of additional matching algorithms and transformation capability, for example to cope with:
  - Missing names
  - Out of order names
  - Typos
  - Etc.
- Complete control over matching