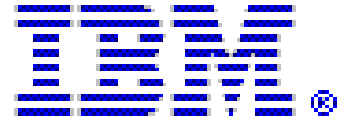


PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.



**IBM Information Server Suite**  
**Course: DX741**  
**QualityStage 8 Essentials**

---

Copyright, Disclaimer of Warranties and Limitation of Liability

© Copyright IBM Corporation February 2007

IBM Software Group  
One Rogers Street  
Cambridge, MA 02142

All rights reserved. Printed in the United States.

IBM and the IBM logo are registered trademarks of International Business Machines Corporation.

The following are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both:

AnswersOnLine	DynamicServer, WorkgroupEdition	RedBrick Decision Server
AIX	Enterprise Storage Server	RedBrickMineBuilder
APPN	FFST/2	RedBrickDecisionscape
AS/400	Foundation.2000	RedBrickReady
BookMaster	Illustra	RedBrickSystems
C-ISAM	Informix	RelyonRedBrick
Client SDK	Informix4GL	S/390
Cloudscape	InformixExtendedParallelServer	
Connection Services	InformixInternet Foundation.2000	Sequent
Database Architecture	Informix RedBrick Decision Server	
DataBlade	J/Foundation	SP
DataJoiner	MaxConnect	System View
DataPropagator	MVS	Tivoli
DB2	MVS/ESA	TME
DB2 Connect	Net.Data	UniData
DB2 Extenders	NUMA-Q	UniData&Design
DB2 Universal Database	ON-Bar	UniversalDataWarehouseBlueprint
Distributed Database	OnLineDynamicServer	UniversalDatabaseComponents
Distributed Relational	OS/2	UniversalWebConnect
DPI	OS/2 WARP	UniVerse
DRDA	OS/390	VirtualTableInterface
DynamicScalableArchitecture	OS/400	Visionary
DynamicServer	PTX	VisualAge
DynamicServer.2000	QBIC	WebIntegrationSuite
DynamicServer with Advanced DecisionSupportOption	QMF	WebSphere
DynamicServer with Extended ParallelOption	RAMAC	
DynamicServer with UniversalDataOption	RedBrickDesign	
DynamicServer with WebIntegrationOption	RedBrickDataMine	

Microsoft, Windows, Window NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

Java, JDBC, and all Java-based trademarks are trademarks or registered trademarks of Sun Microsystems, Inc. in the United States, other countries, or both.

UNIX is a registered trademark of The Open Group in the United States and other countries.

All other product or brand names may be trademarks of their respective companies.

All information contained in this document has not been submitted to any formal IBM test and is distributed on an "as is" basis without any warranty either express or implied. The use of this information or the implementation of any of these techniques is a customer responsibility and depends on the customer's ability to evaluate and integrate them into the customer's operational environment. While each item may have been reviewed by IBM for accuracy in a specific situation, there is no guarantee that the same or similar results will result elsewhere. Customers attempting to adapt these techniques to their own environments do so at their own risk. The original repository material for this course has been certified as being Year 2000 compliant.

This document may not be reproduced in whole or in part without the prior written permission of IBM.

Note to U.S. Government Users – Documentation related to restricted rights – Use, duplication, or disclosure is subject to restrictions set forth in GSA ADP Schedule Contract with IBM Corp.

# Table of Contents

Topic	Page
Data Quality Issues	5
QualityStage 8 Architecture	40
Developing with QualityStage	56
Investigation	80
Standardize	124
Match	198
Survivorship	283
Special Topics Data Quality Methodology Full Run Migration Tool Utility Business Glossary	303

## Course contents

---

- Data quality issues
- Information Server purpose and architecture
- Introduction to DataStage and QualityStage
- Investigation
- Standardization
- Match
- Survivorship
- Special Topics
  - Data quality methodology
  - QualityStage Migration Tool

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



## Data Quality Issues

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

4.0.3

## Unit objectives

- After completing this unit, you should be able to:
  - Describe data quality issues
  - Describe where QualityStage fits into a data cleansing project

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Data quality challenges

- Different or inconsistent standards in structure, format or values
- Missing data, default values
- Spelling errors, data in wrong fields
- Buried information
- Data anomalies

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Data quality – why do we care?

- Accurate reports
- Accurate information for support operations
- Support development of applications that go beyond original scope for which data was designed
  - Master Data Management
  - Data Warehouse
  - Analytical applications

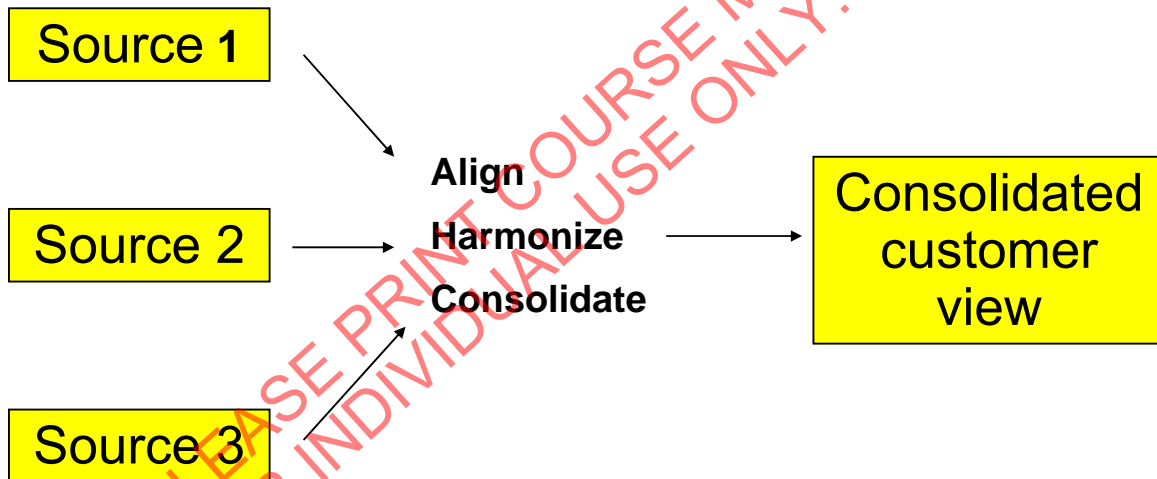
PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



# Master Data Management



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Examples of source systems :

- Orders
- Human resources
- Customer support

The consolidated customer view can represent a data warehouse or a customer relationship management system (CRM).

## Different or inconsistent standards

	Name Field	Location
Source 1	MARK DI LORENZO	MA93
	DENIS E. MARIO	CT15
	TOM & MARY ROBERTS	IL21
Source 2	DILORENZO, MARK	6793
	MARIO, DENISE	0215
	ROBERTS, TOM & MARY	8721
Source 3	MARC DILORENZO ESQ	BOSTON
	MRS DENNIS MARIO	HARTFORD
	MR & MRS T. ROBERTS	CHICAGO

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

- Three sources representing the same four customers, MARK DILORENZO, and DENISE MARIO, TOM ROBERTS and MARY ROBERTS.
- .The & in the third record defines a relationship between two of the customers
- .The ESQ in the first record from source 3 is a title
- .The locations are the same location but coded differently
- .MA93 = 6793 = BOSTON
- .Different structure for storing data: first name lastname and lastname, firstname
- .Sometimes there is a middle initial, sometimes there isn't
- .Not all records have a title word: Mr. Mrs, Esq

## Missing data & default values

Do the field values match the meta data labels?

NAME	SOC. SEC. #	TELEPHONE
Denise Mario DBA	228-02-1975	6173380300
Marc Di Lorenzo ETAL	999999999	3380321
Tom & Mary Roberts	025-37-1888	
First Natl Provident	34-2671434	415-392-2000
	101010101	508-466-1200
Astorial Fedrl Savings	LN#12-756	212-235-1000
Kevin Cooke, Receiver	18-7534216	FAX 528-9825
John Doe Trustee for K	111111111	5436

© Copyright IBM Corporation 2007

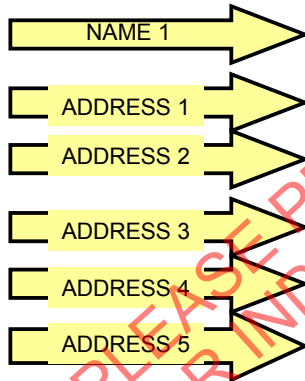
The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

▲Default SSN of 999999999, and 111111111

▲Missing phone numbers

## Buried information

### Legacy Meta Desc.



### Legacy Record Values

Robert A. Jones (TTE) Robert Jones Jr.  
First Natl Provident  
(FBO) Elaine & Michael Lincoln UTA  
(DTD 3-30-89) 59 Via Hermosa  
c/o Colleen Mailer (Esq)  
Seattle, WA 98101-2345

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# The anomalies nightmare

CUSNUM	NAME	ADDRESS	SALES \$
90328574	IBM	187 N.Pk. Str. Salem NH 01456	0
90328575	I.B.M. Inc.	187 N.Pk. St. Sarem NH 01456	0
90238495	International Bus. M.	187 No. Park StSalem NH 04156	0
90233479	Int. Bus. Machines	187 Park Ave Salem NH 04156	0
90233489	Inter-Nation Consults	15 Main St. Andover MA 02341	0
90234889	Int. Bus. Consultants	PO Box 9 Boston MA 02210	00
90345672	I.B. Manufacturing	Park Blvd. Boston MA 04106	00

No common key

Anomalies

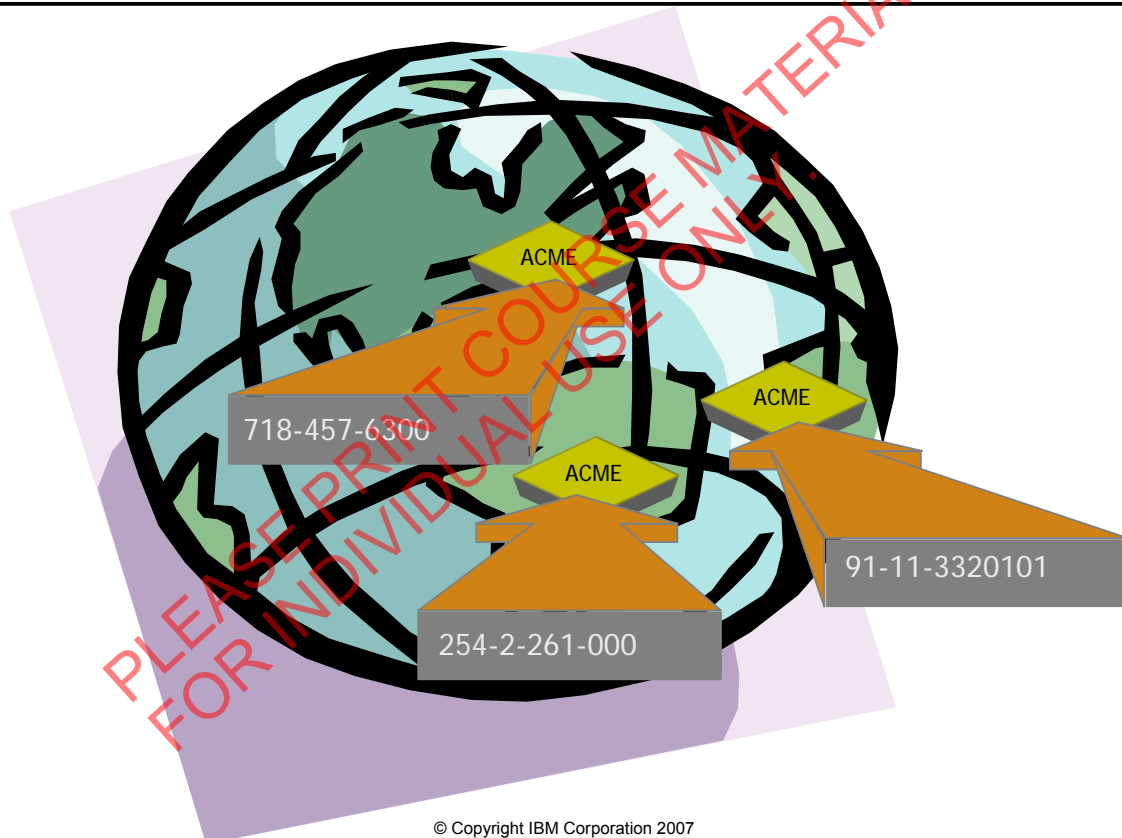
Lack of Standards

Spelling Errors

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## No consolidated of a single entity



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Errors due to:

- Data Entry Errors
- System field limitations
- Mergers and acquisitions
- Feeding legacy data into new systems

## What data challenges do you face?

- No consistent naming convention
- Business terms and spillover text
- Missing values or data in the wrong fields
- Buried information
- Misspelling
- No unique key linking records together

Acct #	Name	Address	City	State	Zip	Note
5154155	Peter J. Lalonde	40 Beacon St.	Melrose, Mass		02176	ODP
5152335	LaLonde, Peter	76 George 617-210-0824	Boston	YES	MA	02111
5146261	Lalonde, Sofie	40 Bacon Street	Melrose		MA	CHK ID
87121	Pete & Soph Lalond	76 George Road	Boston	MASS		FR Alert
87458	P. Lalonde FBO	S.Lalonde40 Becon Rd.	Melrose	MA	02	176

© Copyright IBM Corporation 2007  
 The course materials are provided for internal use only by the individual receiving the materials.  
 The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

The five common data Contaminants or DQ Issues are:

1. Lack of Standards – data coming from disparate systems
2. Spillover and Lack of Domain QualityStage
3. Misspellings and Multiple representations, Missing and invalid data
4. No consolidated key
5. Buried information

## Why investigate?

- Discover trends and potential anomalies in the data
- 100% visibility of single domain and free-form fields
- Identify invalid and default values
- Reveal undocumented business rules and common terminology
- Verify the reliability of the data in the fields to be used as matching criteria
- Gain complete understanding of data within context

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Why is Analysis and Assessment important. What does it provide?



## Investigate – single domain report

- Single domain

Field



Sample  
source  
data

Frequency



% of Total



qsInvColumnName	qsInvPattern	qsInvSample	qsInvCount	qsInvPercent
Zip			61	2.10345
Zip	01385	01385	1	0.0344828
Zip	13104	13104	3	0.103448
Zip	00000	00000	2	0.0689655
Zip	000000000	000000000	1	0.0344828
Zip	00630	00630	1	0.0344828
Zip	00654	00654	1	0.0344828
Zip	007731449	007731449	1	0.0344828
Zip	00926	00926	1	0.0344828
Zip	00927	00927	1	0.0344828
Zip	012389347	012389347	1	0.0344828
Zip	015411134	015411134	2	0.0689655
Zip	017011902	017011902	1	0.0344828
Zip	017462547	017462547	2	0.0689655

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

qsInv – QualityStage Investigate

## Investigate – word pattern report

- Freeform text (Word)

Field

Pattern

Sample  
source  
data

qsInvColumnName	qsInvPattern	qsInvSample	qsInvCount	qsInvPercent
FullName	?,?	TOWNSEND , PANDA	36	1.24138
FullName	?,??	ODISH , SAMER HERMIZ	5	0.172414
FullName	?,?F	COBB , EARLUST ANN	1	0.0344828
FullName	?,?I	BANCARTER , TRENTYN A	25	0.862069
FullName	?,?I.	GRIDLEY , LEVERITT S.	11	0.37931
FullName	?,F	ANTENNA , SALVATORE	332	11.4483
FullName	?,F&F	KAVANACH, MARY & OSCAR	2	0.0689655
FullName	?,F-FI	GIACOBBE , SHERRY-LYNN N	2	0.0689655
FullName	?,F?	BANDINI , DENISE WALSH	31	1.06897
FullName	?,FF	HANLY , PETER MICHAEL	56	1.93103
FullName	?,FFF	WYRICK , JIMMY RAY OLIVER	2	0.0689655
FullName	?,FFI	HOLT , JO ANN C	2	0.0689655
FullName	?,FI	BEERMANN , VERNON E	668	23.0345
FullName	?,FI.	BENAVIDES , KALA R.	74	2.55172

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Frequency

% of Total

## What is standardize?

- Applying business logic to data chaos.
  - Pattern manipulation
- Enforcing business standards on data elements.
  - Standards definition
- Transforming the input to an output which meets the business requirement.
  - Field structuring

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## How to standardize

- Parse specific data fields into smaller, atomic data elements
  - Atomic data elements are called tokens
  - Categorize identified elements
    - Separate Name, Address, and Area from freeform Name & Address lines
    - Identification of Distinct Material Categories (e.g. Sutures vs. Orthopedic Equipment)
- Refine data elements
  - Example 1
    - Name = 'DR PAUL E JONES' becomes:
      - > Title = 'DR'
      - > First Name = 'PAUL'
      - > Middle Name = 'E'
      - > Last Name = 'JONES'
  - Example 2
    - Part Description = 'BLK LATEX GLOVE' becomes:
      - > Color = 'BLACK'
      - > Type = 'LATEX'
      - > Part = 'GLOVE'

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

### Why

- Invoke enterprise wide standards
  - Common usage of abbreviations across the enterprise
  - Single entity definition/metadata across the enterprise
- Improve searching ability
  - Search by single domain entity: FName, LName, PartName,...
- Improve matching ability
  - Entity level matching NOT mixed domain matching:
    - FName → FName ('JOHN' → 'JOHN'),
    - LName → LName ('SMITH' → 'SMITH'),
    - NOT Name → Name ('MR JON P SMITH' → 'SMITH, JON')
  - Match consistent standardized values not free form variations:
    - 'JAMES' to 'JAMES' not 'JIM' to 'JAMES',
    - 'ST' to 'ST' not 'STREET' to 'STR',
    - 'BLACK' to 'BLACK' not 'BLACK' to 'BLK'
- Enable Categorization through standardized single domain entities
  - OFFICE EQUIP = {PC, FAX, COPIER, PRINTER, PHONE,...}

## Why standardize?

- Normalize values in data fields to standard values
  - Transform First Name = 'MIKE' → 'MICHAEL'
  - Transform Title = 'Doctor' → 'Dr'
  - Transform Address = 'ST. Michael Street' → 'Saint Michael St.'
  - Transform Color = 'BLK' → 'BLACK'
- Apply phonetic coding to key words - facilitates record linkage
  - NYSIIS
  - Soundex
  - Typically applied to Name fields (first, last, street, city)

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007  
The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

NYSIIS and Soundex are both phonetic coding algorithms.

## QualityStage standardize

- Uses a highly flexible pattern recognition language
- Can employ field or domain specific standardization (i.e. unique rules for names vs. addresses vs. dates, etc.)
- Contains customizable classification and standardization tables
- Utilizes results from data investigation

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# QualityStage standardize report example

Original data

Ind./Org. flag

FullName	NameType	Gender	FirstName	Middle	PrimaryName_USNAME
ECKER , JACOB E	I	M	JACOB	E	ECKER
KNIAT , KENNETH S	I	M	KENNETH	S	KNIAT
MENARD , LYNETTE H	I	F	LYNETTE	H	MENARD
STARBUCK , F DIANE	I	NULL	F	DIANE	STARBUCK
FIRST UNITED	O	NULL	NULL	NULL	FIRST UNITED
FIRSTAR BK	O	NULL	NULL	NULL	FIRSTAR BANK
APPERT LTD LIAB	O	NULL	NULL	NULL	APPERT LTD LIAB
BERTHA L KARRER	I	F	BERTHA	L	KARRER
J. BERNARD	I	NULL	J	NULL	BERNARD
JOHN F WIBLE TRST	O	NULL	NULL	NULL	JOHN F WIBLE TRUST
NELLIE HEALD	I	F	NELLIE	NULL	HEALD
BOROWITZ FAM TRUST	O	NULL	NULL	NULL	BOROWITZ FAM TRUST
FRANCIS BALLMAN TRUST	O	NULL	NULL	NULL	FRANCIS BALLMAN TRUS
OLGA DUEMELAND	I	F	OLGA	NULL	DUEMELAND
EUGENE B BOROWITZ	I	M	EUGENE	B	BOROWITZ
DONALD R HALL	I	M	DONALD	R	HALL

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Match

“Conditioned data and QualityStage’s matching engine link the previously unlinkable.”

- Match Construction:
  - Reliability of input data defines a match result.
- Statistical Analysis & Match Scoring:
  - Linkage probability determined on a sliding scale by field level comparison.
- Report Generation:
  - All business rules applied have easy to understand report structure.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



## What is match?

- Identifying all records on one file that correspond to similar records on another file
- Identifying duplicate records in one file
- Building relationships between records in multiple files
- Performing statistical and probabilistic matching
- Calculating a score based on the probability of a match

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Why match?

---

- Identify duplicate entities within one or more files
- Perform householding
- Create consolidated view of customer
- Establish cross-reference linkage
- Enrich existing data with new attributes from external sources

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## How to match

- Single file (Unduplication) or two file (Reference or Geomatch)
- Different match comparisons for different types of data (e.g. exact character, uncertainty/fuzzy match, keystroke errors, multiple word comparison ...)
- Generation of composite weights from multiple fields
- Use of probabilistic or statistical algorithms
- Application of match cutoffs or thresholds to identify automatic and clerical match levels
- Incorporation of override weights to assess particular data conditions (e.g. default values, discriminatory elements)

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## QualityStage match

- A wide variety of match comparison algorithms providing a full spectrum of fuzzy matching functions
- Statistically-based method for determining matches (Probabilistic Record Linkage Theory)
- Field-by-field comparisons for agreement or disagreement
- Assignment of weights or penalties
- Overrides for unique data conditions
- Score results to determine the probability of matched records
- Thresholds for final match determination
- Ability to measure informational content of data

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# QualityStage match examples

SetID	Record Type	FullName	AddressLine1	City
2126	XA	GEROSA, FRAN X	C/O NANCY C GEROSA	RIDGEFIELD
2126	DA	GEROSA, FRANCIS X	C/O NANCY C GEROSA	RIDGEFIELD
2126	DA	GEROSA, FRANK X	C/O NANCY C GEROSA	RIDGEFIELD
2126	DA	GEROSA, FRANCIS X	C/O NANCY C GEROSA	RIDGEFIELD
62	XA	BIONDI, KATHERINE A.	3142 CENTRAL ST	EVANSTON
62	DA	BIONDI, KATHERINE A.	3142 CENTRAL ST	EVANSTON
254	XA	STEFAN, JOHN R.	11009 AZALEA DR	PITTSBURGH
254	DA	STEFAN, JOHN R.	11009 AZALEA DR	PITTSBURGH
750	XA	RUMMEL, JACK R	640 SUMMERGREEN DRIVE	FRANKENMUTH
750	DA	RUMMEL, JACK R	640 SUMMERGREEN DR	FRANKENMUTH
15	XA	BANGERTER, EDWARD L	2060 CANDLE TREE CV	SANDY
15	DA	BANGERTER, EDWARD L	2060 CANDLE TREE CV	SANDY
389	XA	GOLDBLATT, RICHARD J	6410 TARREGA ST	CORAL GABLES
389	DA	GOLDBLATT, RICHARD J	6410 TARREGA ST	CORAL GABLES
431	XA	COLLINS, THERESA A	3699 CLAY ST APT 2	SAN FRANCISCO
431	DA	COLLINS, THERESA A	3699 CLAY ST APT 2	SAN FRANCISCO
134	XA	GRANT MORROW III	253 N COLUMBIA AVE	COLUMBUS
134	DA	GRANT MORROW III	253 N COLUMBIA AVENUE	COLUMBUS
1954	XA	ELKODSI, SUSAN L	8 SUNBEAM DR	TRUMBULL
1954	DA	ELKODSI, SUE	8 SUNBEAM DRIVE	TRUMBULL

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## What is survive?

- Creation of best-of-breed “surviving” data based on record or field level information
- Development of cross-reference file of related keys
- Production of load exception reports
- Creating output formats:
  - Relational table with primary and foreign keys
  - Transactions to update databases
  - Cross-reference files, synonym tables
  - Audit trails, exception reports

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Why survive?

- Provide consolidated view of data
- Provide consolidated view containing the “best-of-breed” data
- Resolve conflicting values and fill missing values
- Cross-populate best available data
- Implement business and mapping rules
- Create cross-reference keys

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## How to survive

- Highly flexible rules
- Record or field level survivorship decisions
- Rules can be based upon data frequency, data recency (i.e. date), data source, value presence or length
- Rules can incorporate multiple tests
- QualityStage features
  - Point-and-click (GUI-based) creation of business rules to determine best-of-breed “surviving” data
  - Performed at record or field level

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Survive has its own rules (not the same as standardization).



# QualityStage survive examples

## Example 1:

- The longest populated Middle and Last Name

Matched			Survived		
First Name	Middle Name	Last Name	First Name	Middle Name	Last Name
MARI		LEMELSON-LAPPNER	MARI	S	LEMELSON-LAPPNER
MARI	S	LEMELSON			

## Example 2:

- The longest populated Middle Name, Date of Birth, and SSN

Matched					Survived				
First Name	Middle Name	Last Name	DOB	SSN	First Name	Middle Name	Last Name	DOB	SSN
DENISE		TRIANO	19580211	98524173	DENISE	F	TRIANO	19580211	98524173
DENISE	F	TRIANO							

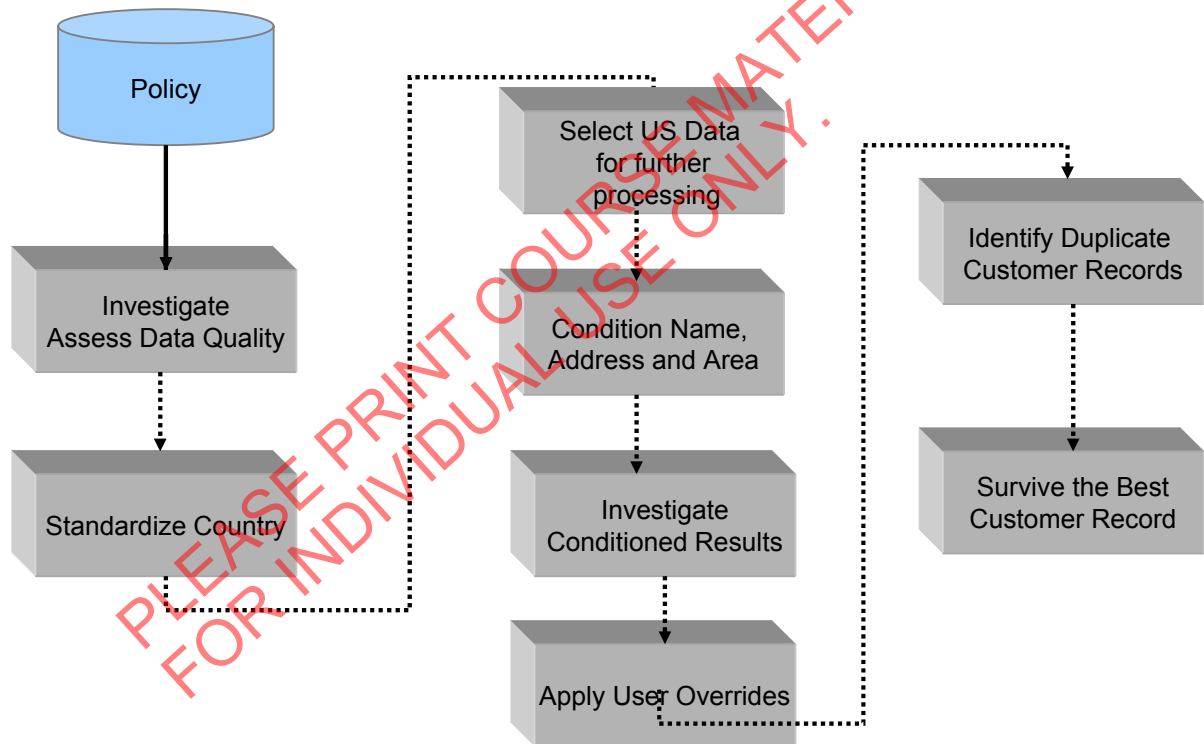
© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

These examples are field-level survivorship. Individual fields are mixed to form an output record.

Record-level survivorship would choose one record over another, perhaps based on source system or date.

# Course exercise project design



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint

1. (T/F) Data quality investigation cleans the source data.
2. (T/F) Standardization modifies the source data so that it can be loaded into the target system.
3. (T/F) Survivorship data can be either record based or field based.

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint solutions

1. (T/F) (T/F) Data quality investigation cleans the source data.  
*Answer: False*
2. (T/F) Standardization modifies the source data so that it can be loaded into the target system.  
*Answer: False*
3. Survivorship data can be either record based or field based.  
*Answer: True*

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Unit summary

Having completed this unit, you should be able to:

- Describe the five common data quality contaminants
  - Different standards
  - Missing and default values
  - Spillover and buried information
  - Anomalies
  - No consolidated view
- Describe each of the following processes:
  - Investigation
  - Standardization
  - Match
  - Survivorship

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 1: Review course project

---

- Course business case: WINN Insurance CRM project
- See QualityStage Essentials Exercises

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 2: copy student files

- Copy student files to disk
  - Use C: drive as root for folder

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



## QualityStage 8 Architecture

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

4.0.3



## Unit objectives

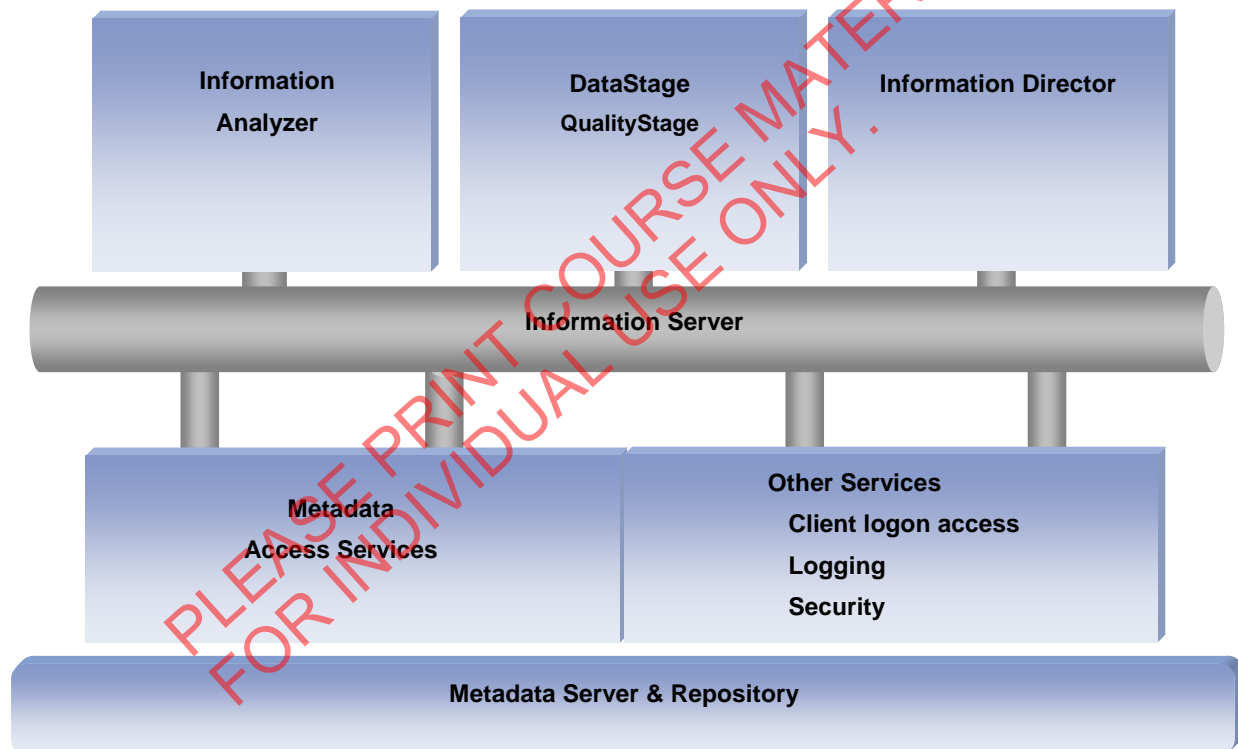
- After completing this unit, you should be able to:
  - Describe the Data Quality architecture
  - Identify data quality server and client components
  - Describe the methods of client/server communication

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Information Server conceptual architecture



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Information Server is really a suite of applications:

- ▲DataStage
- ▲QualityStage
- ▲Information Analyzer
- ▲Business Glossary
- ▲Information Director

QualityStage is an optional, add-on component to DataStage.

## QualityStage technical highlights

- Uses Enterprise level DataStage
  - DataStage design environment
  - Parallel execution engine
  - Stages are native enterprise operators
  - Match designer is embedded in DataStage Designer Client
  - Get DataStage data connectivity by default
    - No need for meta brokers, plug-ins
    - Common meta data
- Legacy (pre-version 8) QS job execution
  - Migration utility available to aid conversion from QS 7.x to QS 8
  - Converted jobs can be compiled and executed in the QS 8 environment

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

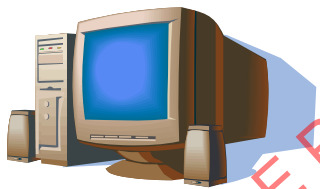
Data Quality stages are only available on Enterprise jobs.

# DataStage/QualityStage physical architecture

## Clients

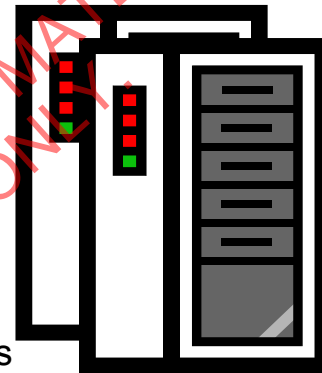
DataStage/QualityStage

Designer  
Director  
Administrator



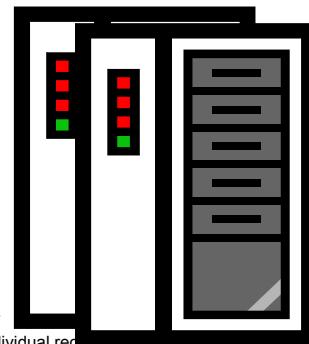
Windows

## Information Server



UNIX

Projects



Windows

Connect to projects  
via TCP/IP

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## DataStage clients

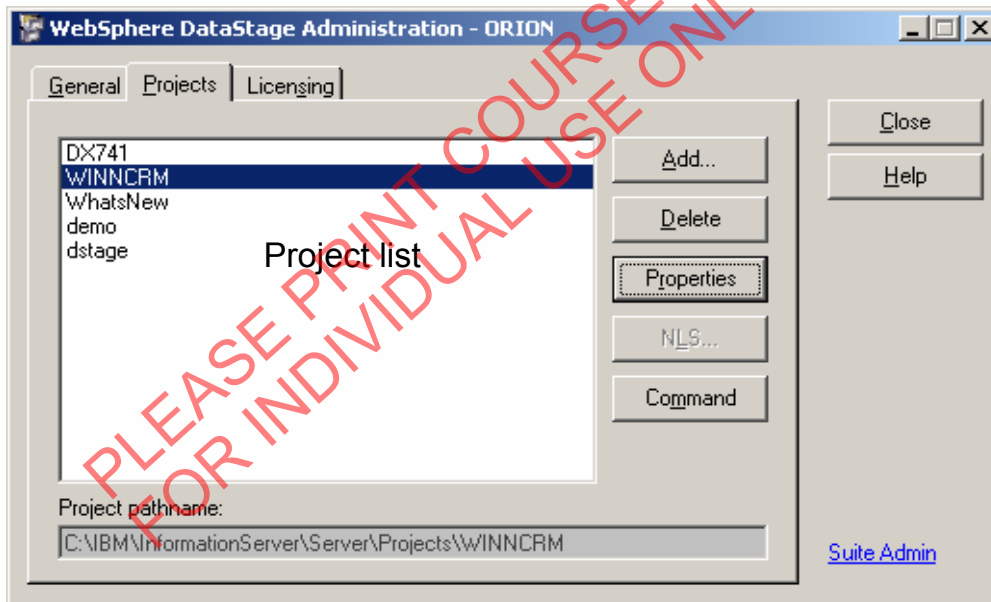
- Administrator
  - Add and delete projects
  - Set project defaults
  - Set project environment parameters
- Designer
  - Maintain data definitions
  - Add, modify, and delete jobs
  - Add, modify, and delete match specifications
  - Manage rule sets
  - Compile jobs
  - Run jobs
  - Provision rule sets and match specifications
- Director
  - Run jobs
  - Review job log
  - Schedule jobs

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# DataStage Administrator

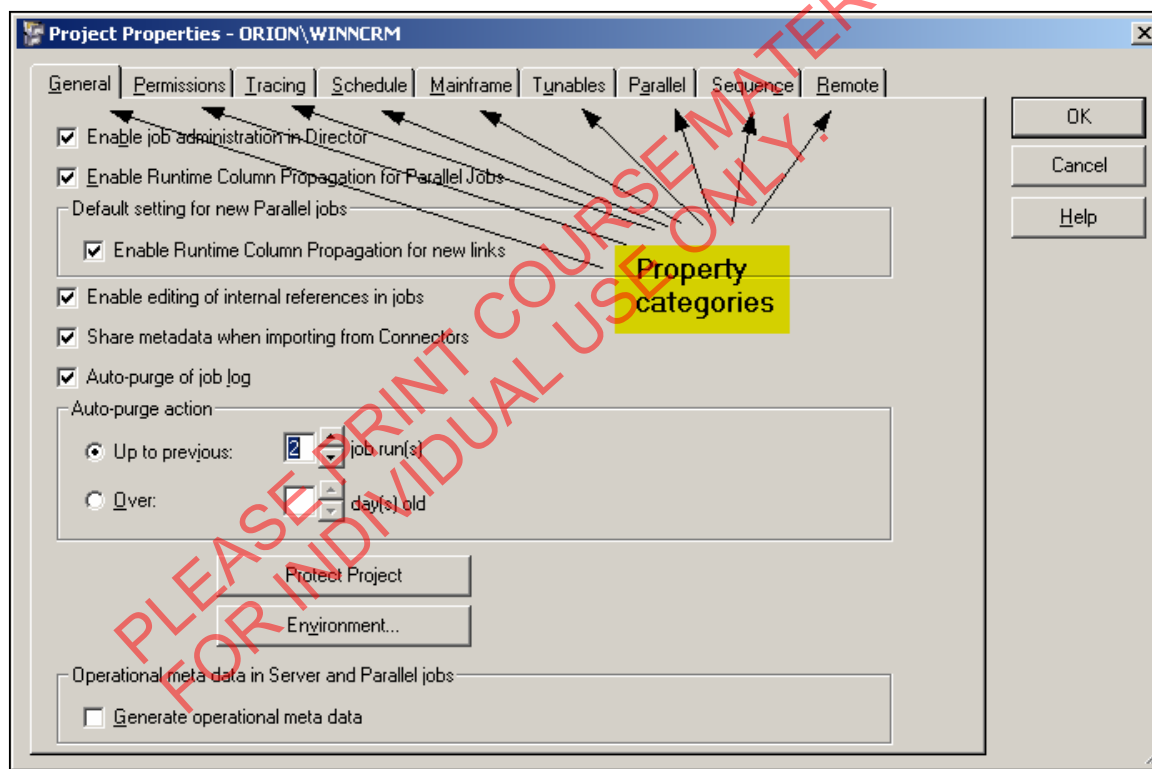
- Administrator
  - Create or delete projects
  - Set project defaults
  - Apply security



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# Project property defaults



© Copyright IBM Corporation 2007

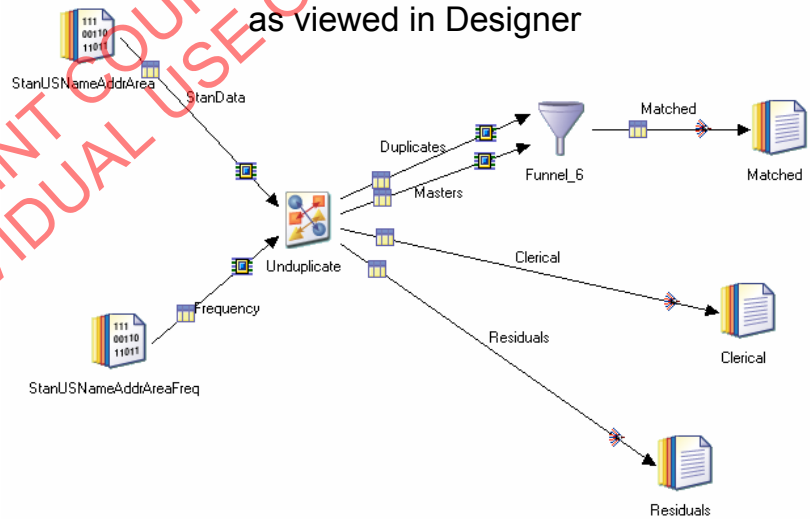
The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Project properties are divided into categories and accessed with the GUI tabs.

# DataStage Designer

- Designer
  - Client GUI for designing jobs
    - Windows 2000+, XP
    - Build meta data
    - Build Jobs
    - Modify Standardization Rules
    - Build match specifications
  - Designer Repository
    - Database

Sample QualityStage job  
as viewed in Designer



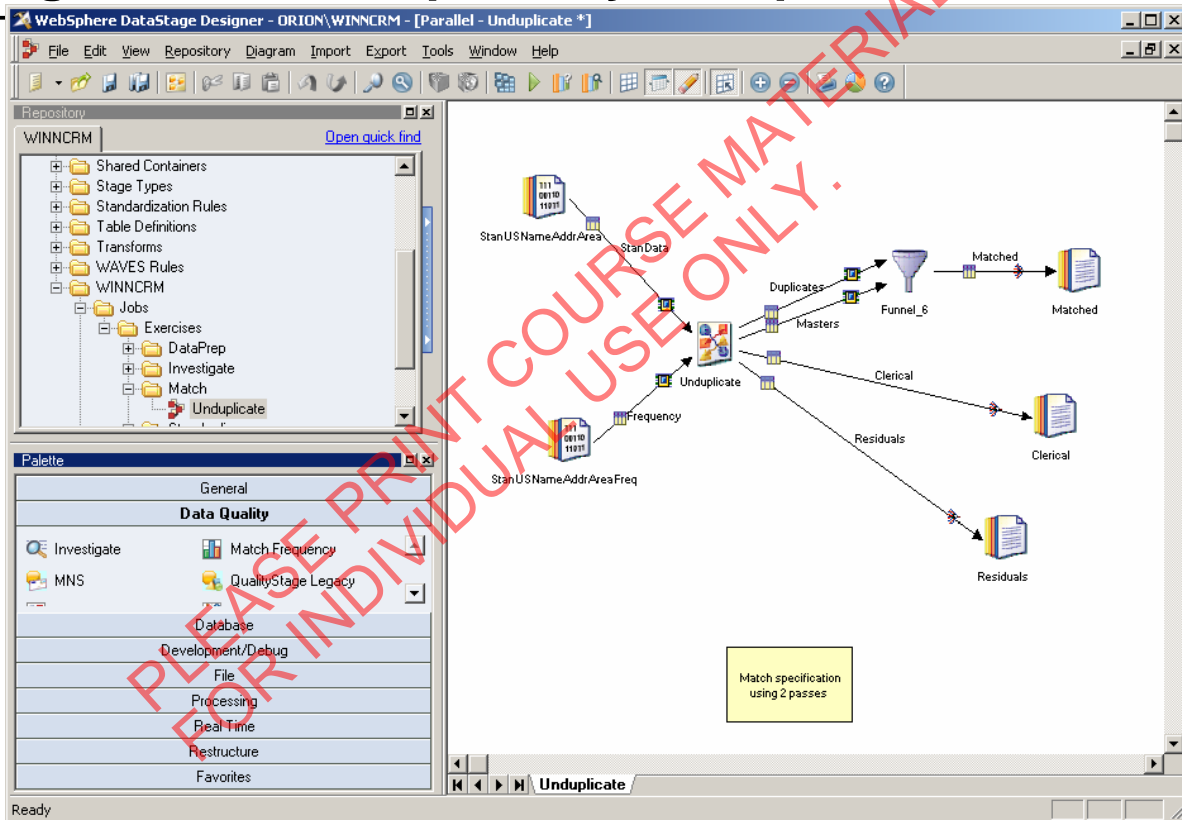
© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

QualityStage jobs are data flow oriented; the direction of the arrows shows the flow of data. Stages process the data.



# Designer canvas, repository, and palette

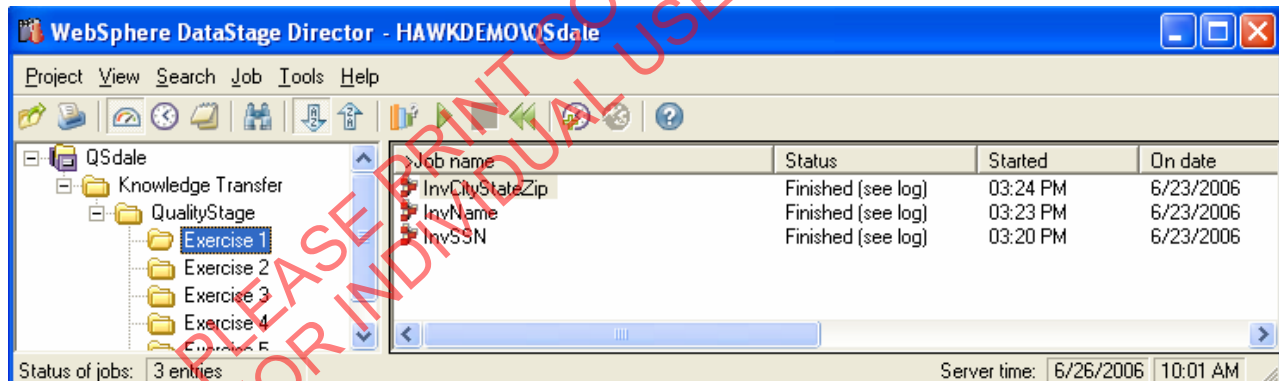


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# DataStage Director

- Director
  - Client GUI for managing job execution
  - Windows 2000+, XP
  - Run jobs – set job options and parameters
  - View job log
  - Schedule job execution



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

QualityStage jobs now produce standard DataStage job log. The log should be viewed for statistics.

# Job log viewed in Director

12:23:06 PM	12/7/2006	Control	Starting Job Unduplicate.
12:23:09 PM	12/7/2006	Info	Environment variable settings: (...)
12:23:09 PM	12/7/2006	Info	Parallel job initiated
12:23:09 PM	12/7/2006	Info	QSH script (...)
12:23:10 PM	12/7/2006	Info	main_program: IBM WebSphere DataStage Enterprise Edition 8.0.0 (...)
12:23:12 PM	12/7/2006	Info	main_program: orchgeneral: loaded (...)
12:23:13 PM	12/7/2006	Info	Unduplicate: Creating sub-operator: <QSmats -workDir ./.RT_QS23/V057-f Unduplicate.txt> (...)
12:23:13 PM	12/7/2006	Info	main_program: APT configuration file: C:/IBM/InformationServer/Server/Configurations/default.apt (...)
12:23:13 PM	12/7/2006	Warning	StanUSNameAddAreaFreq: When checking operator: When binding output schema variable "outRec": When binding ...
12:23:13 PM	12/7/2006	Warning	StanUSNameAddAreaFreq: When checking operator: When binding output schema variable "outRec": When binding ...
12:23:13 PM	12/7/2006	Warning	Residuals: When checking operator: When validating export schema: At field "StanZip3": "null_field" length (4) must ma...
12:23:13 PM	12/7/2006	Warning	Residuals: When checking operator: A sequential operator cannot preserve the partitioning (...)
12:23:13 PM	12/7/2006	Warning	Clerical: When checking operator: When validating export schema: At field "StanZip3": "null_field" length (4) must matc...
12:23:13 PM	12/7/2006	Warning	Clerical: When checking operator: A sequential operator cannot preserve the partitioning (...)
12:23:13 PM	12/7/2006	Warning	Matched: When checking operator: When validating export schema: At field "StanZip3": "null_field" length (4) must mat...
12:24:05 PM	12/7/2006	Warning	Matched: When checking operator: A sequential operator cannot preserve the partitioning (...)
12:24:06 PM	12/7/2006	Info	Unduplicate,0: Variable: GenderCode_USNAM (...)
12:24:06 PM	12/7/2006	Info	Unduplicate,0: 0126366747 3 0 0 0.90 0.00 D 9.73 -3.32 (...)
12:24:06 PM	12/7/2006	Info	Unduplicate,0: Frequency table(s) will be used
12:24:06 PM	12/7/2006	Info	Unduplicate,0: Default weights calculated for values OUTSIDE table (...)
12:24:06 PM	12/7/2006	Info	Unduplicate,0: <Pass 1> Blocks processed: 1275 (...)
12:24:07 PM	12/7/2006	Info	Unduplicate,0: Variable: GenderCode_USNAM (...)
12:24:07 PM	12/7/2006	Info	Unduplicate,0: 0126366747 3 0 0 0.90 0.00 D 9.73 -3.32 (...)
12:24:07 PM	12/7/2006	Info	Unduplicate,0: Frequency table(s) will be used
12:24:07 PM	12/7/2006	Info	Unduplicate,0: Default weights calculated for values OUTSIDE table (...)
12:24:07 PM	12/7/2006	Info	Unduplicate,0: <Pass 2> Blocks processed: 121 (...)
12:24:09 PM	12/7/2006	Info	Unduplicate,0: "" Output Statistics For UNDUPLICATE "" (...)
12:24:09 PM	12/7/2006	Info	Unduplicate,0: 2843 data records & 1599 match records joined
12:24:09 PM	12/7/2006	Info	Residuals,0: Export complete; 1244 records exported successfully, 0 rejected.
12:24:09 PM	12/7/2006	Info	Clerical,0: Export complete; 0 records exported successfully, 0 rejected.
12:24:10 PM	12/7/2006	Info	Matched,0: Export complete; 1599 records exported successfully, 0 rejected.
12:24:10 PM	12/7/2006	Info	main_program: Step execution finished with status = OK.
12:24:12 PM	12/7/2006	Info	main_program: Startup time, 0:34; production run time, 0:26.
12:24:16 PM	12/7/2006	Info	Contents of phantom output file (...)
12:24:17 PM	12/7/2006	Info	Contents of phantom output file (...)
12:24:17 PM	12/7/2006	Info	Parallel job reports successful completion
12:24:18 PM	12/7/2006	Control	Finished Job Unduplicate.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint

---

1. (T/F) DataStage Administrator executes jobs.
2. (T/F) DataStage Designer configures projects.
3. Which DataStage component displays objects in the designer database?

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint solutions

1. (T/F) DataStage Administrator executes jobs.  
*Answer: False*
2. (T/F) DataStage Designer configures projects.  
*Answer: False*
3. Which DataStage component displays objects in the designer database.  
*Answer: the repository view*

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Unit summary

Having completed this unit, you should be able to:

- Describe the Information Server components
- List the DataStage clients
- Describe a typical DataStage configuration

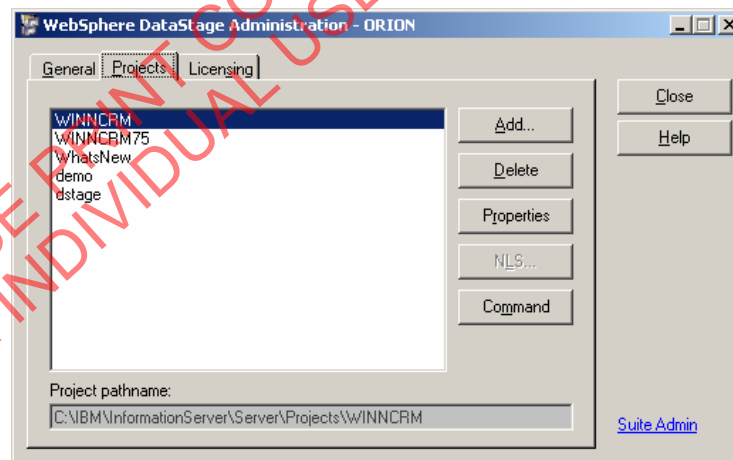
PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 3: configure QualityStage project

- Create a project using Administrator
- Set project properties
  - General defaults
  - Environment variables
  - Security groups and roles



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



## Developing with QualityStage

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

4.0.3



## Unit objectives

- After completing this unit, you should be able to:
  - Define data files and field definitions
  - Build DataStage Jobs
  - Deploy and run jobs
  - Locate and review results

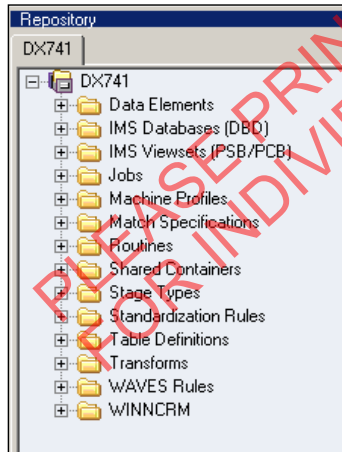
PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## QualityStage application

- Could comprise one or more projects
- Project components
  - Jobs
  - Stages
  - Data File Definitions
    - Meta data
  - Designer repository view shows project components



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Project components are displayed in the repository view as a hierarchy of folders.

## Job definition

---

- A job is an executable DataStage/QualityStage program
- Created by job compilation
- Jobs can be run in batch or in real time

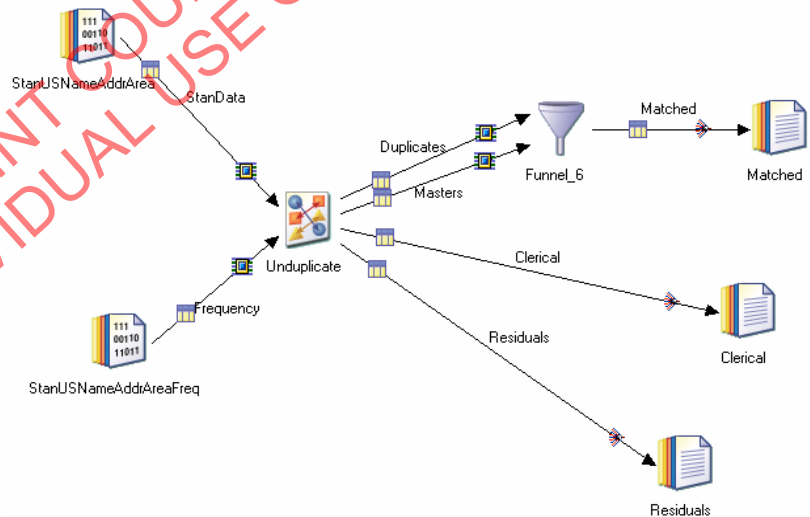
PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Job development overview

- Designer
  - Import or enter file definitions and meta data defining your sources and targets
  - Add stages and links defining the process or task
  - Compile the job
  - Run the job
  - Review results files
- Server
  - Runs the job
  - View job log



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Log onto project in Designer or Director

**Attach to Project**

Domain:  
orion:9080

User name:  
admin

Password:  
\*\*\*\*\*

Project:  
ORION/DX741  
ORION/demo  
ORION/dstage  
ORION/DX741  
ORION/whatsNew  
ORION/WINNCRM

OK  
Cancel  
Help

User name and Password controlled by Information Server

List of valid projects

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Logon service is provided by the Information Server backbone.

## Designer repository components

- Database which stores
  - Data file definitions
  - Job designs
  - Standardization rules
  - Callable components such as buildops
  - Containers
  - Data connection objects

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

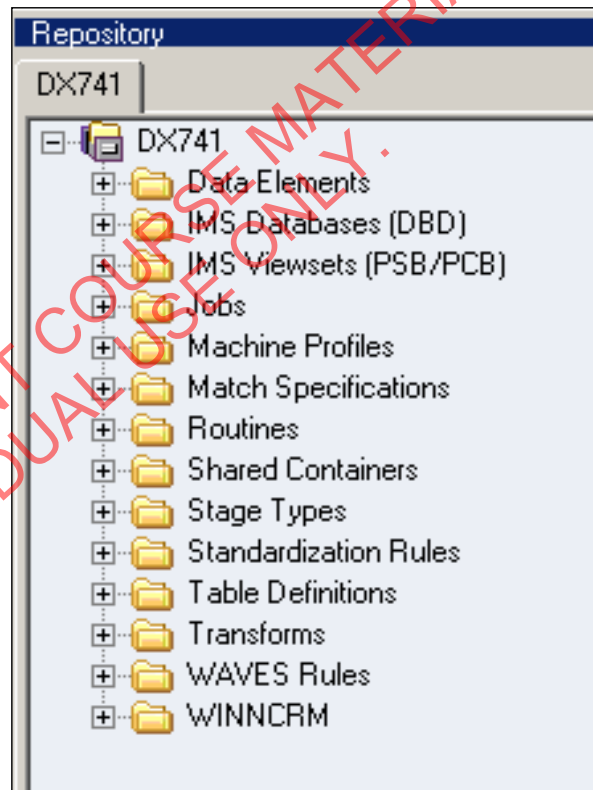
© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Project structure

Repository view

In Designer



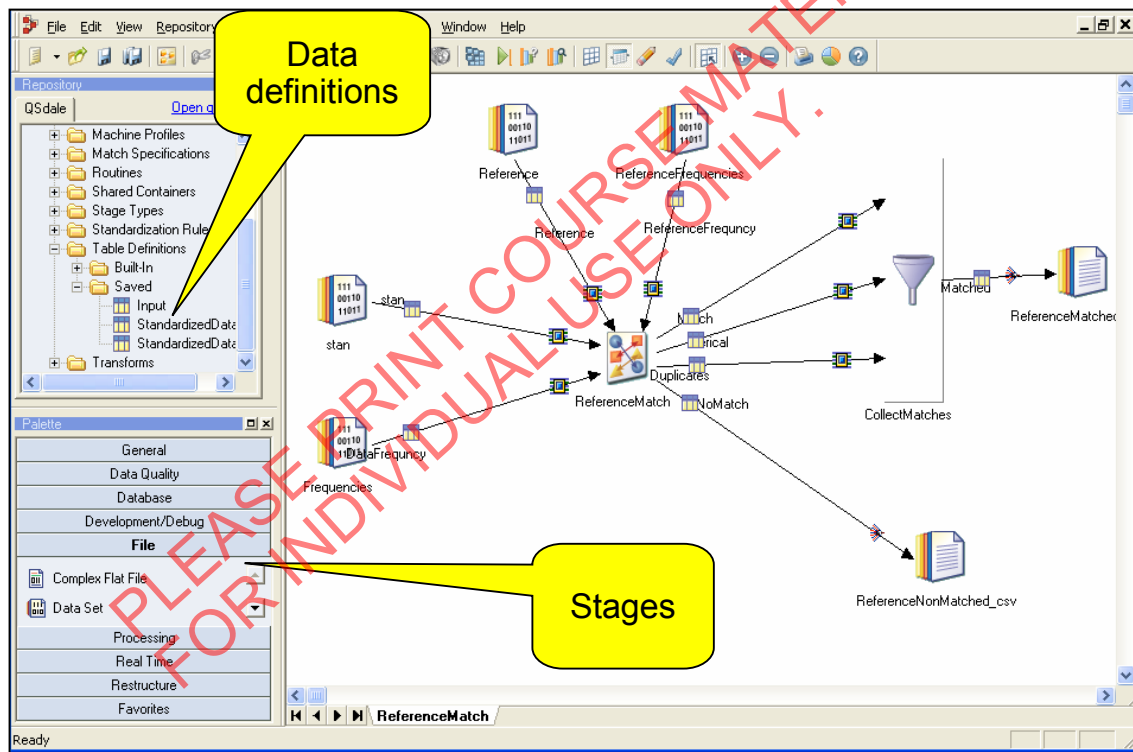
© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

DataStage components are arranged in a folder structure.

User can create new folders.

# DataStage design environment



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

The DataStage Designer has full import/export features – both meta data and DataStage components – as well as job designs.

QualityStage jobs have a DataStage visual representation. The above job is a two file reference match.

Jobs are comprised of stages and links.

Stages are functional units and links indicate the flow of data.

To build a job:

- Drag stages from palette

- Draw links between the stages



## Data definitions

- Entered or loaded via DataStage import mechanisms
  - Sequential file
  - ODBC
  - Plug-ins
  - MetaBrokers
- New and redefined columns can be added on the data flow via Transformer stage

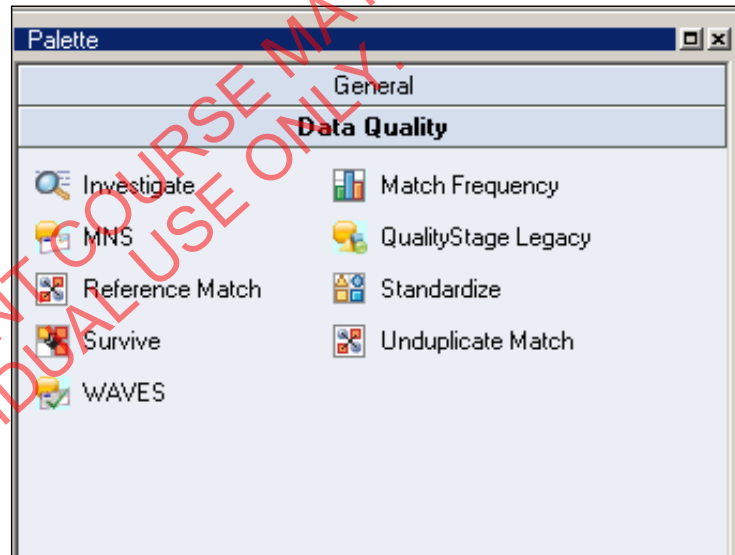
PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Data Quality folder

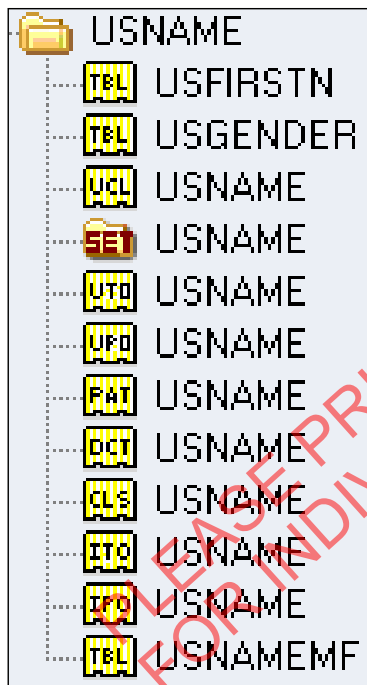
- Stages are the building blocks
- Focused in function
- All phases of data quality:
  - Investigate
  - Standardize
  - Match Frequency
  - Match
    - Unduplicate Match
    - Reference Match
  - Survive



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Standardization rule sets



Rule set for USERNAME

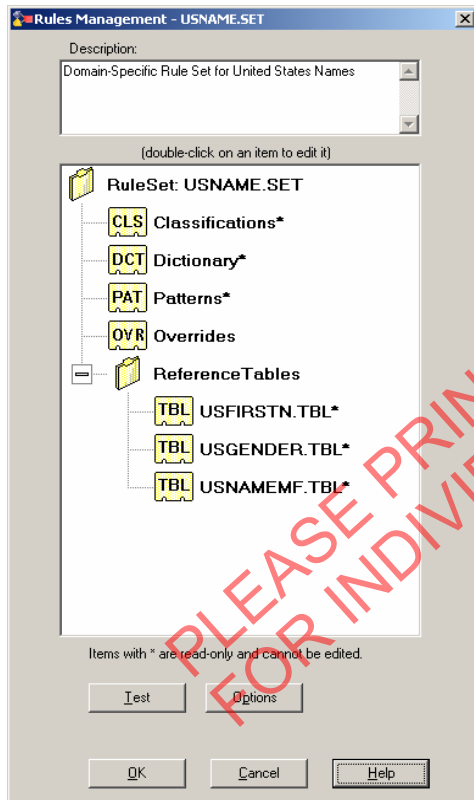
- Pre-defined rules for parsing and standardizing:
  - Name
  - Address
  - Area (City, State and Zip)
- Multi-national address processing
- Validate structure:
  - Tax ID
  - US Phone
  - Date
  - Email
- Append ISO country codes
- Rule sets are stored in the repository and provisioned to the job execution area

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Provisioning copies the rule set from the repository to the job execution area.

## Rule set components



- Can modify some rule set components
- Test rule sets
- Copy rule sets

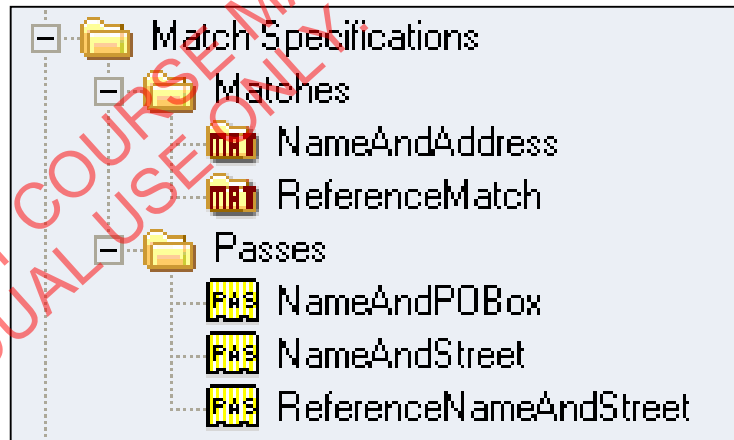
© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Rule sets are a group of control files that determine the standardization process.

## Match Specifications in the DataStage Repository

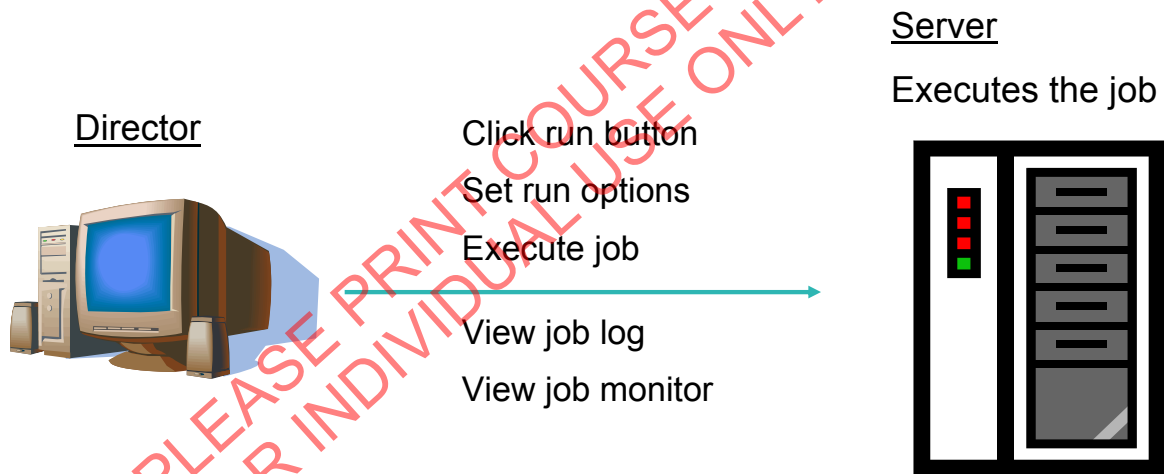
- Created using the Match Designer
- Allows online testing of match criteria



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Executing a job via Director



© Copyright IBM Corporation 2007

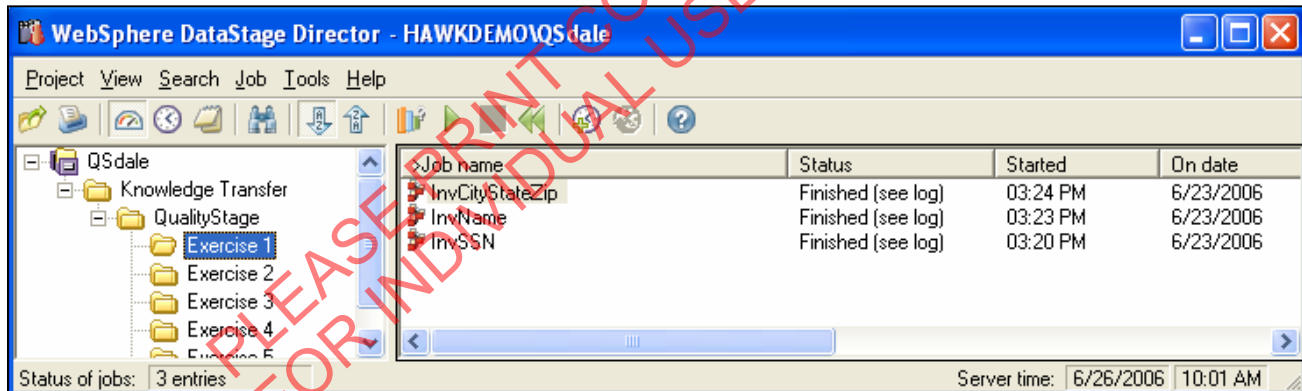
The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

A job can also execute via a shell script independent of the Director.

# Running a job in Director

- Director
  - Client GUI for running jobs
    - Windows 2000+, XP
    - View job logs and monitor
    - Job scheduling

Job status view



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

QualityStage jobs are data flow oriented; the direction of the arrows shows the flow of data. Stages process the data.

# Execution environment

**WebSphere DataStage Director - HAWKDEMO\QSDale**

Project View Search Job Tools Help

Project: QSDale

Knowledge Transfer

QualityStage

Exercise 1

Exercise 2

Exercise 3

Exercise 4

Exercise 5

Job name	Status	Started	On date
InvCityStateZip	Finished (see log)	03:24 PM	6/23/2006
InvName	Finished (see log)	03:23 PM	6/23/2006
InvSSN	Finished (see log)	03:20 PM	6/23/2006

Status of jobs: 3 entries

Server time: 6/26/2006 10:01 AM

**Data Quality Job Log**

3:24:58 PM 6/23/2006 Info Parallel job initiated

3:24:58 PM 6/23/2006 Info OSH script (...)

3:24:59 PM 6/23/2006 Info main\_program: IBM WebSphere DataStage Enterprise Edition 8.0.0 (...)

3:24:59 PM 6/23/2006 Info main\_program: orchgeneral: loaded (...)

3:24:59 PM 6/23/2006 Info InvestigateCityStateZip: Creating sub-operator: <field\_export-field qslnvSample -type string -schema '

3:24:59 PM 6/23/2006 Info main\_program: APT configuration file: D:\IBM\WDIS\Server\Configurations\default.apf (...)

3:24:59 PM 6/23/2006 Warning CityStateZipTokenReport: When checking operator: A sequential operator cannot preserve the partit

3:25:02 PM 6/23/2006 Warning CityStateZipPatternReport: When checking operator: A sequential operator cannot preserve the parti

3:25:02 PM 6/23/2006 Info InvestigateCityStateZip,0: User Classification Override file: ./RT\_QS1\V0S22/Controls\USAREA.ucl f

3:25:02 PM 6/23/2006 Info CustomerFile,0: Import complete; 600 records imported successfully, 0 rejected.

3:25:02 PM 6/23/2006 Info InvestigateCityStateZip,0: Field export complete. 600 records converted successfully, 0 rejected.

3:25:02 PM 6/23/2006 Info InvestigateCityStateZip,0: 600 input records processed

3:25:02 PM 6/23/2006 Info CityStateZipTokenReport,0: Export complete; 58 records exported successfully, 0 rejected.

3:25:02 PM 6/23/2006 Info InvestigateCityStateZip,0: 600 input records read; 4 kept

3:25:02 PM 6/23/2006 Info CityStateZipPatternReport,0: Export complete; 4 records exported successfully, 0 rejected.

3:25:02 PM 6/23/2006 Info main\_program: Step execution finished with status = OK.

3:25:02 PM 6/23/2006 Info main\_program: Startup time, 0:03; production run time, 0:00.

3:25:03 PM 6/23/2006 Info Parallel job reports successful completion

3:25:04 PM 6/23/2006 Control Finished Job InvCityStateZip.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

QualityStage jobs now produce the standard DataStage job log. The log should be viewed for statistics.



## Job Monitor statistics

WebSphere DataStage Director Monitor - InvCityStateZip

Stage/Link name	Link type	Status	Num rows	Started at	Elapsed time	Rows/sec
<b>CustomerFile</b>		Finished	600	:24:58 PM	00:00:04	150
Customer	>Out		600			150
<b>Copy</b>		Finished	600	:24:58 PM	00:00:04	150
Customer	<<Pri		600			150
ToCityStateZip	>Out		600			150
<b>InvestigateCityStateZip</b>		Finished	1200	:24:58 PM	00:00:04	300
ToCityStateZip	<<Pri		1200			300
PatternReport	>Out		0			
TokenReport	>Out		58			14
<b>CityStateZipPatternReport</b>		Finished	4	:24:58 PM	00:00:04	1
PatternReport	<<Pri		4			1
<b>CityStateZipTokenReport</b>		Finished	58	:24:58 PM	00:00:04	14
TokenReport	<<Pri		58			14

Job: InvCityStateZip Status: Finished (see log) Project: QSDale (HAWKDEMO) Server time: 10:07 AM

Interval: 10

Close

Help

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Job development process

- Define data files
  - Enter or import meta data
- Define job
  - Draw stages and links
  - Set stage properties
  - Compile
- Run the job
- Review results

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint

---

1. (T/F) The job monitor displays link statistics.
2. (T/F) The job log is viewed in DataStage Designer.
3. What protocol is used for communication between the DataStage clients and server?

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint solutions

1. (T/F) The job monitor displays link statistics.  
*Answer: True*
2. (T/F) The job log is viewed in DataStage Designer.  
*Answer: False*
3. What protocol is used for communication between the DataStage clients and server?  
*Answer: TCPIP*

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Unit summary

Having completed this unit, you should be able to:

- Define data files and field definitions
- Design jobs
- Deploy and run jobs
- Locate and review results

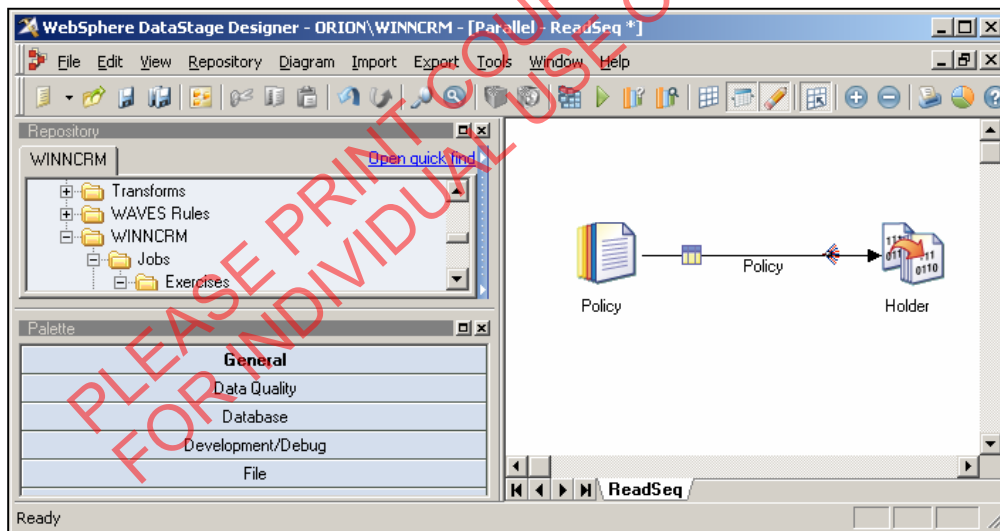
PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 4: Import meta data

- DataStage import mechanisms
  - DataStage components
    - Any object built in DataStage, such as jobs, table definitions, match specifications

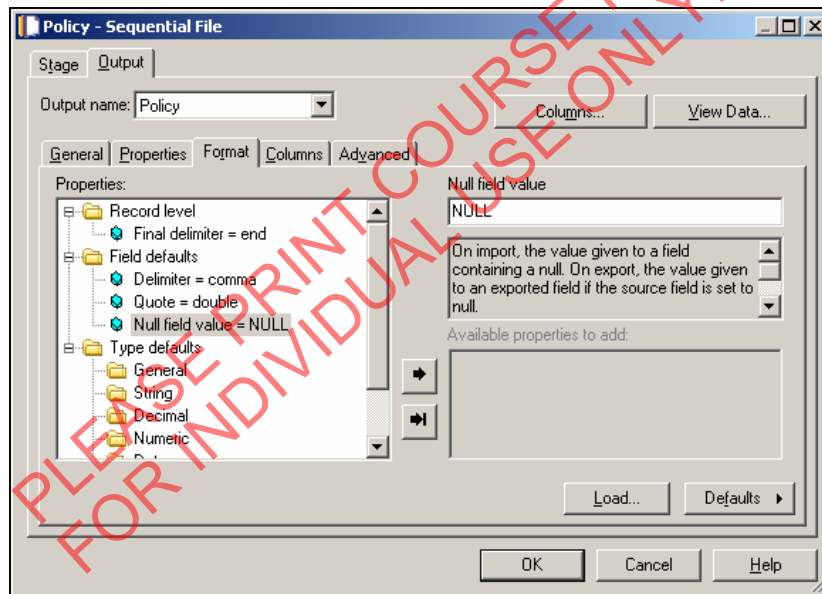


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 5: Build and run DataStage job

- Read sequential file
  - Must use format tab to handle nulls



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



# Investigation

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.



© Copyright IBM Corporation 2007  
The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

4.0.3



## Unit objectives

- After completing this unit, you should be able to:
  - Define data investigation
  - Build Investigate stages
  - Use character discrete, concatenate, and word investigations to analyze data fields
  - Locate and review results

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Investigation

- Verify the domain
  - Review each field and verify the data matches the meta data
- Identify data formats, missing and default values
- Identify data anomalies
  - Format
  - Structure
  - Content
- Discover “unwritten” business rules
- Identify data preparation requirements

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Investigate stage

- Features
  - Analyze free-form and single domain fields
  - Provide frequency distributions of distinct values and patterns
- Investigate methods
  - Character Discrete
  - Character Concatenate
  - Word

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Four Common Methods:

Character discrete – Inv multiple single-domain fields independently

Type C – View the character values

Type T – View the field format or “Template”

Type X – Ignore characters

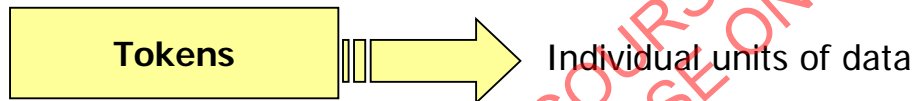
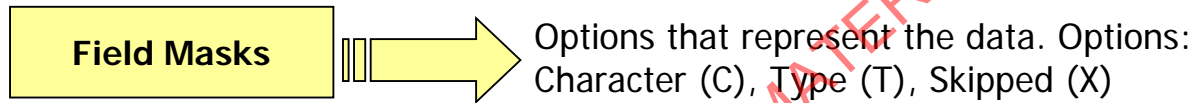
## Investigate methods

Method	Why
Character Discrete	Analyzing field values, formats, and domains
Character Concatenate	Cross-field correlation, checking logic relationships between fields
Word Investigation	Identifying free-form fields that may require parsing and discovery of key words for classification

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Investigate terminology



Character Mask	Usage
<b>C</b>	<b>For viewing the actual character values of the data</b>
<b>T</b>	<b>For viewing the pattern of the data</b>
<b>X</b>	<b>For ignoring characters</b>

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the  
The materials may not be modified, copied, distributed or transferred without the express prior written permission of IBM Corporation.

Examples of tokens:

string – 914 Josephine St

Token 1 – 914

Token 2 – Josephine

Token 3 - St

## Field mask examples

<b>Token</b>	<b>Mask</b>	<b>Result</b>
02116	<b>CCCCC</b>	02116
02116	<b>CCCXX</b>	021
01832-4480	<b>TTTTTTT</b>	nnnnn-nnnn
XJ2 6EM	<b>TTTTTT</b>	aanbnaa
(617) 338-0300	<b>CCCCCCCCCCCCC</b>	(617) 338-0300
617-338-0300	<b>TTTTTTTTTTTTT</b>	nnn-nnn-nnnn
6173380300	<b>CCCXXXXXXXXXX</b>	617
(617)3380300	<b>CCCXXXXXXXXXX</b>	(61

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Character discrete: field mask (C)haracter

- Usage: Domain quality
  - View the contents of each field to verify that the data values match the field labels
- Mechanism: Investigate stage
  - Generates Reports for frequency and pattern references

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

These reports provide the quantitative understanding of data values that will permit correlation of the various spellings, misspellings, abbreviations or other representation of data values

- ▲Also note any anomalies (anything suspect: out of range or defaults values), and how often each anomaly occurs?
- ▲Percent Populated per field: Note how often the field is populated
- ▲How many formats “templates” exist for the data?
- ▲The cardinality of the field: The number of distinct values
- ▲The frequency distribution: How often does each format occur?
- ▲How often does “data in the wrong domain” occur?

## Character discrete - character results

qsInvColumnName	qsInvPattern	qsInvSample	qsInvCount	qsInvPercent
SourceSystem	A	A	1534	52.8966
SourceSystem	H	H	366	12.6207
SourceSystem	L	L	1000	34.4828
PolicyNumber	003668461	003668461	1	0.0344828
PolicyNumber	003775219	003775219	2	0.0689655
PolicyNumber	004281148	004281148	1	0.0344828
PolicyNumber	004793986	004793986	1	0.0344828
PolicyNumber	004804210	004804210	1	0.0344828

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Investigates one or more single-domain fields

Each field is treated independently for frequency count and pattern reporting

Report names

- jobp.FRQ - sorted by frequency in descending order
- jobp.SRT - sorted alphabetical in ascending order
- job.PAT - reference file



## Character discrete: field mask (T)ype

- Usage: Data formats (patterns):
  - View the format of field which contain that you suspect may follow or conform to a specific format, e.g., dates, PIN, Tax ID, account numbers.
- Generates reports for frequency and pattern references

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

These reports provide the quantitative understanding of data values that will permit correlation of the various spellings, misspellings, abbreviations or other representation of data values

- ▲Also note any anomalies (anything suspect: out of range or defaults values), and how often each anomaly occurs?
- ▲Percent Populated per field: Note how often the field is populated
- ▲How many formats “templates” exist for the data?
- ▲The cardinality of the field: The number of distinct values
- ▲The frequency distribution: How often does each format occur?
- ▲How often does “data in the wrong domain” occur?

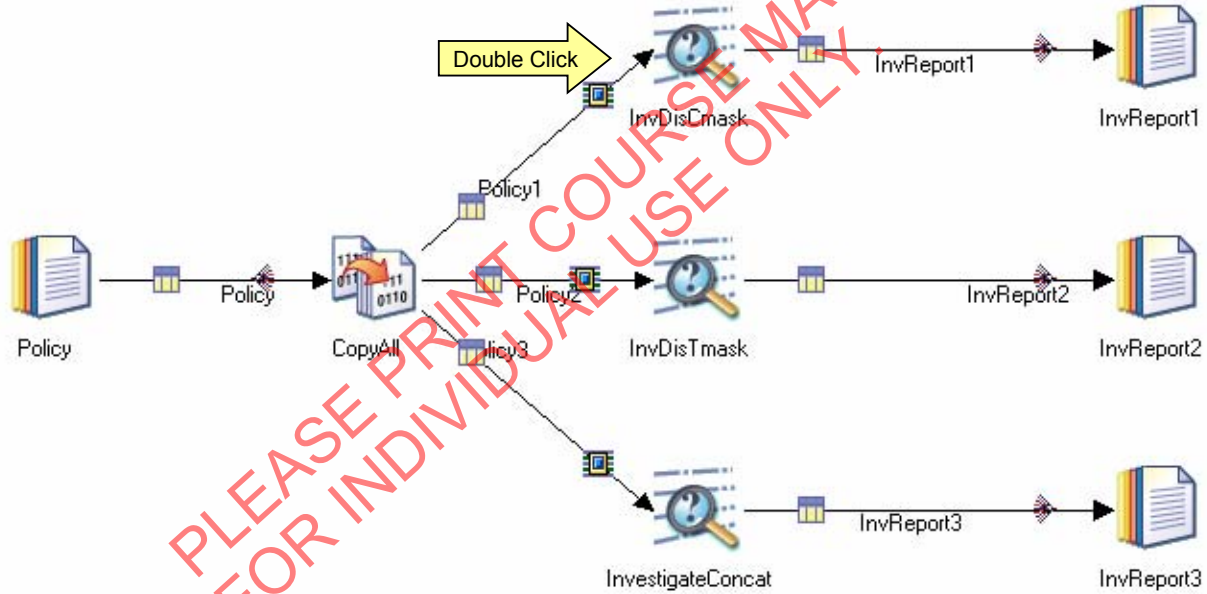
# Investigation Implementation

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

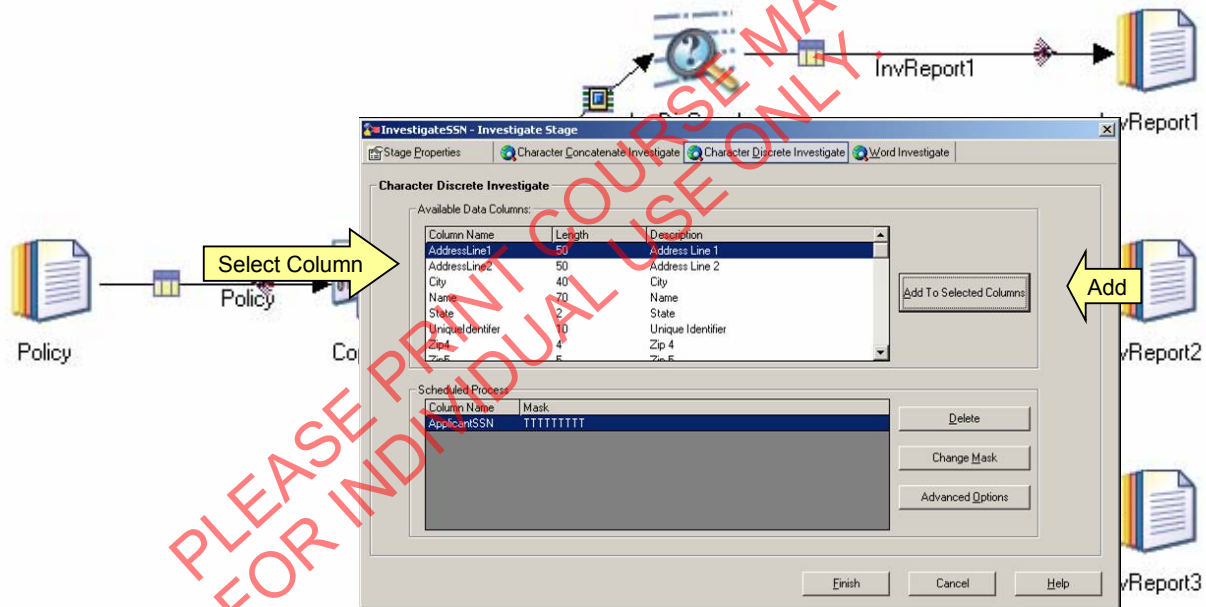
# QualityStage Investigation job – character



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

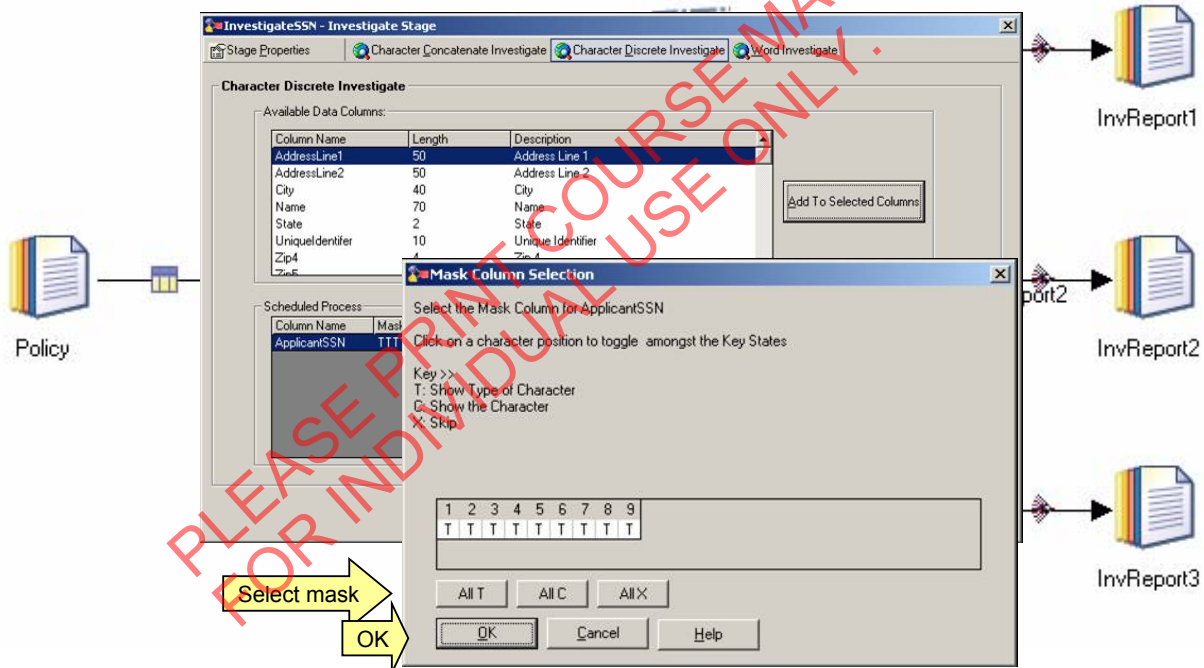
# Investigation - Character



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

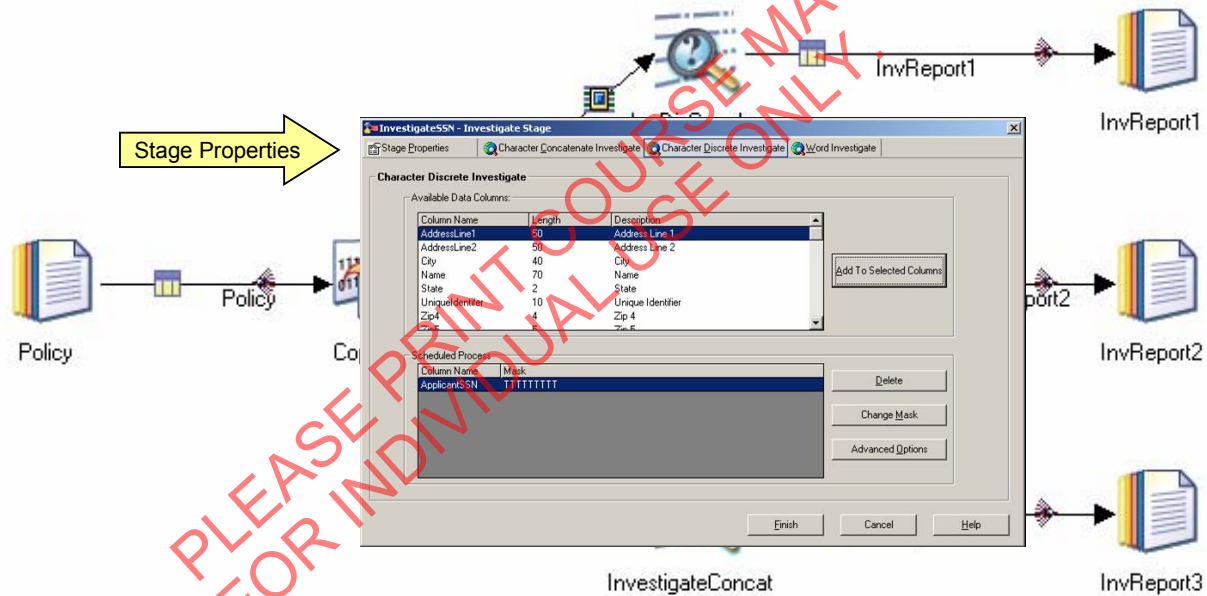
# Investigation - Character



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
 The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

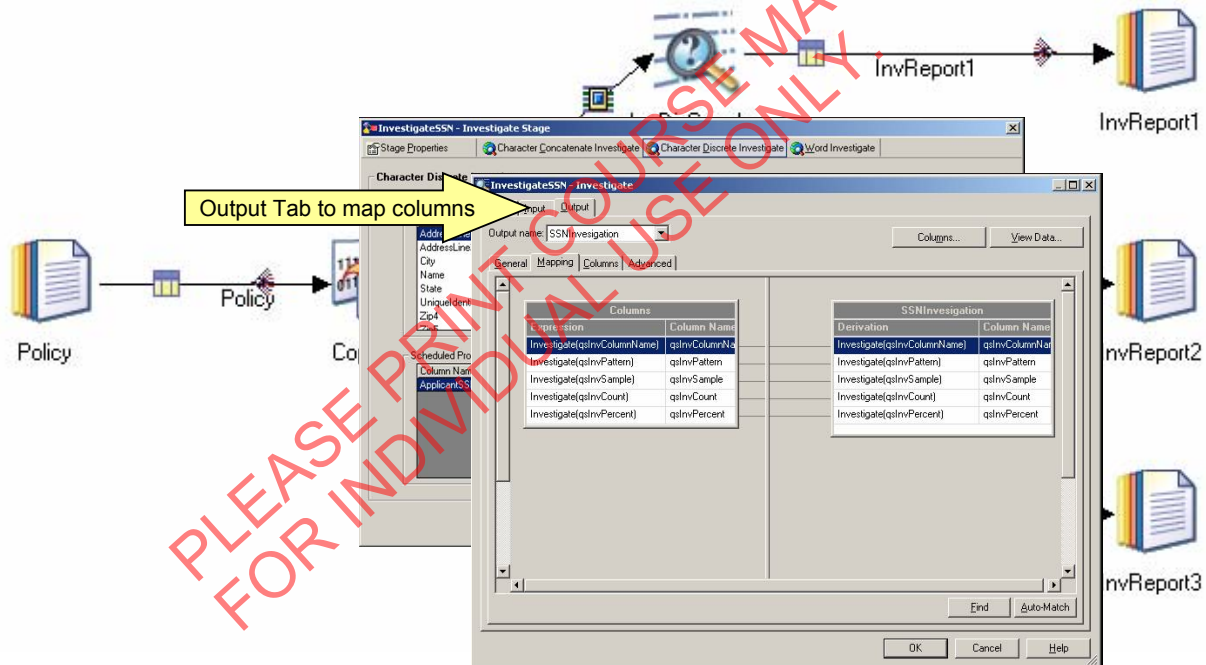
# Investigation - Character



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

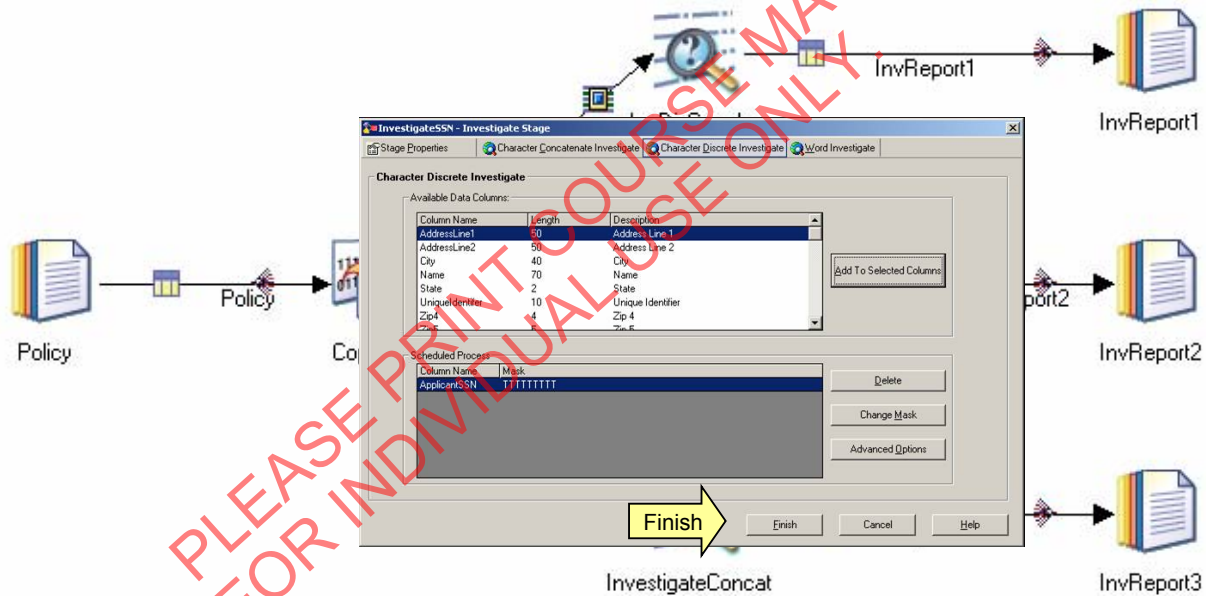
# Investigation - Character



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# Investigation - Character

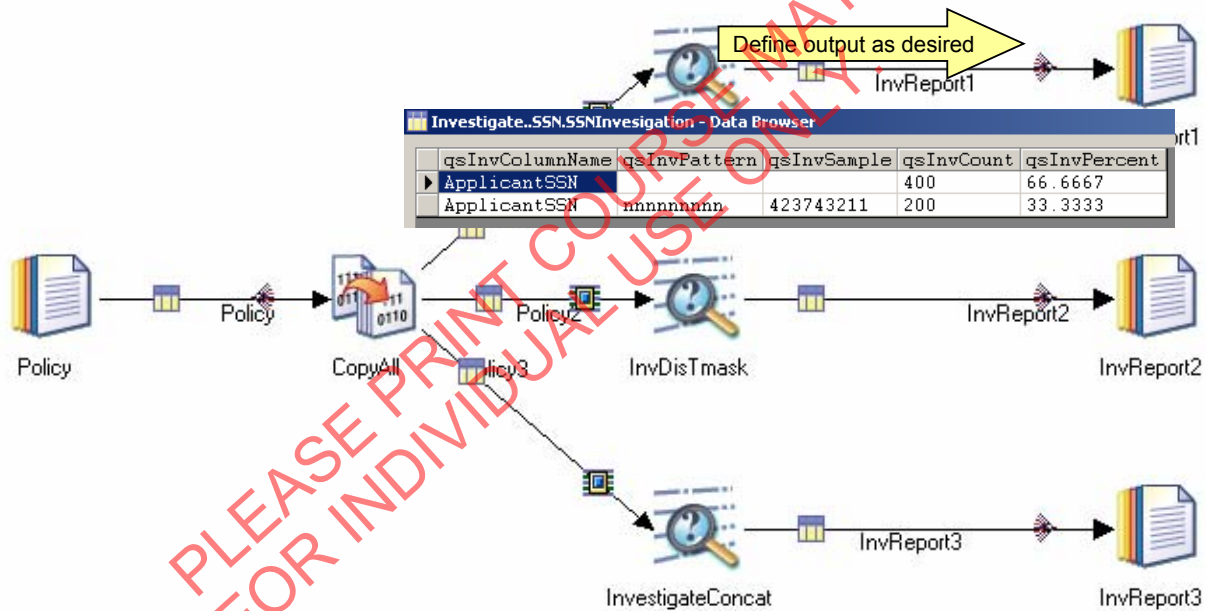


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



# Investigation - Character



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Character concatenate

- Identify Field Relationships
  - Investigate one or more fields to uncover any relationship between the field values.
  - Uses combinations of character masks
  - Generates Reports for frequency and pattern references

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

These reports provide the quantitative understanding of data values prevalence that will permit correlation of the various spellings, misspellings, abbreviations or other representation of data values

- ▲Also note any anomalies (anything suspect: out of range or defaults values), and how often each anomaly occurs?
- ▲Percent Populated per field: Note how often the field is populated
- ▲How many formats “templates” exist for the data?
- ▲The cardinality of the field: The number of distinct values
- ▲The frequency distribution: How often does each format occur?
- ▲How often does “data in the wrong domain” occur?

## Character concatenate results

### DOB and DOD Fields

qsInvColumnName	qsInvSample	qsInvCount	qsInvPercent
DOB+DOD		1184	40.8417
DOB+DOD	000000000000000000	2	0.0689893
DOB+DOD	190812150000000000	1	0.0344947
DOB+DOD	190901010000000000	1	0.0344947
DOB+DOD	191406090000000000	3	0.103484
DOB+DOD	191503300000000000	2	0.0689893
DOB+DOD	191507160000000000	1	0.0344947
DOB+DOD	191702250000000000	2	0.0689893
DOB+DOD	191703310000000000	2	0.0689893

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Word investigate

- Usage: Pattern free-form fields and lexical analysis
  - To view the pattern of the data within a freeform text field and parse it into individual tokens
- QualityStage process
  - Apply rule sets to free-form fields
  - Discover parsing requirements
  - Discover patterns in data
  - Generate reports for word frequency, pattern frequency distributions, and word classification

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

These reports provide the quantitative understanding of data values prevalence that will permit correlation of the various spellings, misspellings, abbreviations or other representation of data values

- ▲Also note any anomalies (anything suspect: out of range or defaults values), and how often each anomaly occurs?
- ▲Percent Populated per field: Note how often the field is populated
- ▲How many formats “templates” exist for the data?
- ▲The cardinality of the field: The number of distinct values
- ▲The frequency distribution: How often does each format occur?
- ▲How often does “data in the wrong domain” occur?

# Word investigation results

## Pattern Reports

^D?T	639 N MILLS AVE
^D?S	306 W MAIN ST
^D?T	3142 W CENTRAL AVE
^?T	843 HEARD AVE

## Word Frequency Reports

0000000869	ST
0000000791	RD
0000000622	STE
0000000566	AVE

## Word Classification Reports

ABBOTT	ABBOTT	?	;0000000001
ABERCON	ABERCON	?	;0000000001
ABERCORN	ABERCORN	?	;0000000007
ABERDEEN	ABERDEEN	?	;0000000001

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Parses free-form data into individual tokens

Tokens are classified to create patterns

Uses a set of rules for parsing and classifying the tokens

Discover tokens (key words) to be added to the classification table such as name prefixes, business terminology, street types, new abbreviations for cities

Create patterns of data tokens with the field context

Identify spelling, misspellings and representations of data

Identify parsing requirements for the conditioning process

**Patterns Reports:** Distinct patterns within the field

Pattern Reports

List of all patterns sorted by frequency (p.frq)

List of all patterns sorted alphabetical (p.srt)

List of each token and it's associated pattern (.pat)

**Word Frequency Reports:** The frequency distribution of distinct values

List of all alpha sorted by frequency (c.frq)

List of all alpha sorted alphabetically (a.frq)

May include numerics and mixed tokens

**Word Classification Reports:** The frequency distribution of "classified" and "unclassified" words

List of classified alpha (u.dlt)

List of not-classified alpha (n.dlt)

All alpha listed in the classification table are considered classified alpha

## Rule sets

- Rules for parsing, classifying, and organizing data
- Rule Set Domains
  - Country processing
  - Pre-processing
  - Domain Processing
    - Name: Business and Personal
    - Street Address
    - Area: Locality, City, State and Zip/Postal codes
  - Multinational Address Processing

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Parsing

- Parse free-form data with the SEPLIST and a STRIPLIST
  - SEPLIST - Any character in the SEPLIST will separate tokens, and become a token itself
  - STRIPLIST - Any character in the STRIPLIST will be ignored in the resulting pattern
- The SEPLIST is always applied first

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

There is a default seplist and striplist on the Word investigation: Advanced Options screen

Remember investigation is about discovery, feel free to changes the seplist and strip to experiment and identify the “best” parsing parameters”

The seplist and striplist only allow “simple” parsing (parsing by encountering the presence of a single character. More complex parsing can be done in the next phase, Conditioning.

If you really aren’t sure how to parse the data then be very conservative, that is separate by a space and only strip out a space. Add in more characters after analysis of the results

The rule is “Whenever in doubt don’t strip out”. If a character sometimes adds context and sometimes does not then DON’T strip out the character. Stripping the character loses context in all cases. Often we will choose the separate by the character but not strip it out:

Examples:

½ if we strip out the / we won’t know if this started as ½ or 12 or two independent digits 1 and 2.

C/O: again if we strip out the / we may not realize that this was an abbrev of “care of” and interpret the token as “company”

## Parsing example

Example: 120 Main St. N.W.

SEPLIST "↵."

STRIPLIST "↵"

<i>Token1</i>	<i>Token2</i>	<i>Token3</i>	<i>Token4</i>	<i>Token5</i>	<i>Token6</i>	<i>Token7</i>	<i>Token8</i>
<b>120</b>	<b>Main</b>	<b>St</b>	<b>.</b>	<b>N</b>	<b>.</b>	<b>W</b>	<b>.</b>

SEPLIST "↵"

STRIPLIST "↵."

<i>Token1</i>	<i>Token2</i>	<i>Token3</i>	<i>Token4</i>
<b>120</b>	<b>Main</b>	<b>St</b>	<b>NW</b>

SEPLIST "↵."

STRIPLIST "↵"

<i>Token1</i>	<i>Token2</i>	<i>Token3</i>	<i>Token4</i>	<i>Token5</i>
<b>120</b>	<b>Main</b>	<b>St</b>	<b>N</b>	<b>W</b>

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.

The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

↵ designates a space.

We talk about tokens back in the early parts of investigation. This slide helps make it more clear how we create "tokens"

Example

APT. is a period you can strip BUT \$10.00 is not a period you can strip without changing the meaning of the data.



## Data typing: classifying tokens

- Identify and type the token in terms of its business meaning and value

### MASK KEY:

N – Numeric token

A – Alpha token

AN – Mixed Token

120	Main	Street	Apt	6C
NNN	AAAA	AAAAAA	AAA	AN

### PATTERN KEY(USNAME rule set):

^ – Numeric token

? – Unclassified alpha token

@, <, > – Mixed Token

T – Street Type

U – Unit Type

120	Main	Street	Apt	6C
^	?	T	U	>

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Start with an example of simple classification. Is the data Alpha, Numeric or Mixed then introduce more sophisticated classification, the classification table.

## Example: word investigate

4	BCH	T
1	BLDG	L
71	BLVD	T
4	BND	T
212	BOX	B

Token report

<>T	N7283 110TH ST
<?T	N10030 FAIRGROUND AVE
<LT	\$9591 CHURCH ROAD
>&HH	6TH & NEW HAMPSHIRE

Pattern report

Produce Reports based on  
Patterns & Tokens

Classify known words  
and  
assign default tags

	^	?	T	U	^
Parse	10	MAPLE	STREET	APARTMENT	222

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

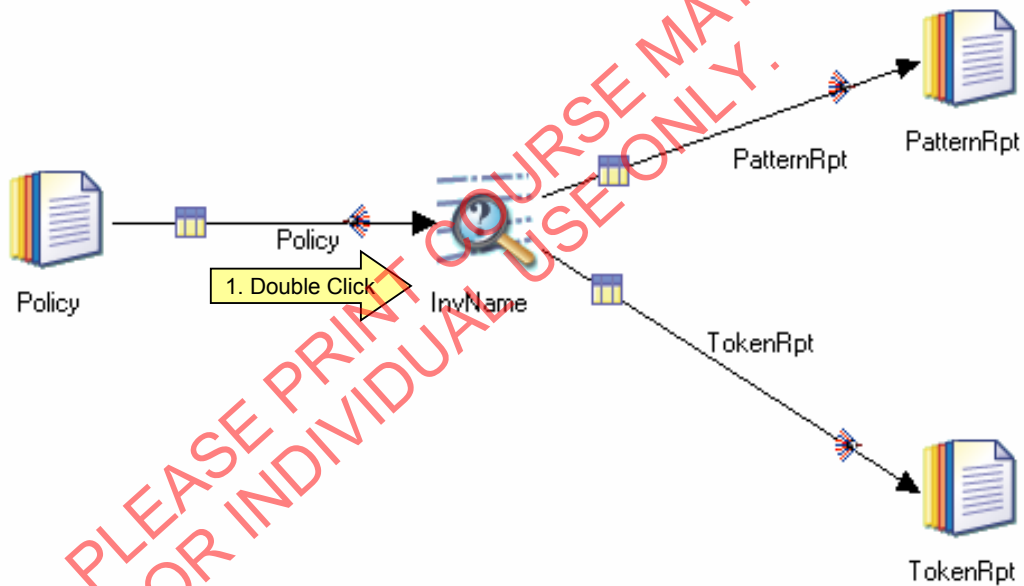
Illustrates the process for Parsing and classifying data tokens to create patterns

Example:

Rebuild	120	Apple	RD	APT	4B
DataType	HN	SN	ST	UT	UV
Pattern	^?TM>				
Classify2	^	?	T	M	>
Classify1	N	A	A	A	M
Parse	120	Apple	Road	Apart	4B

Start with the bottom line and build.

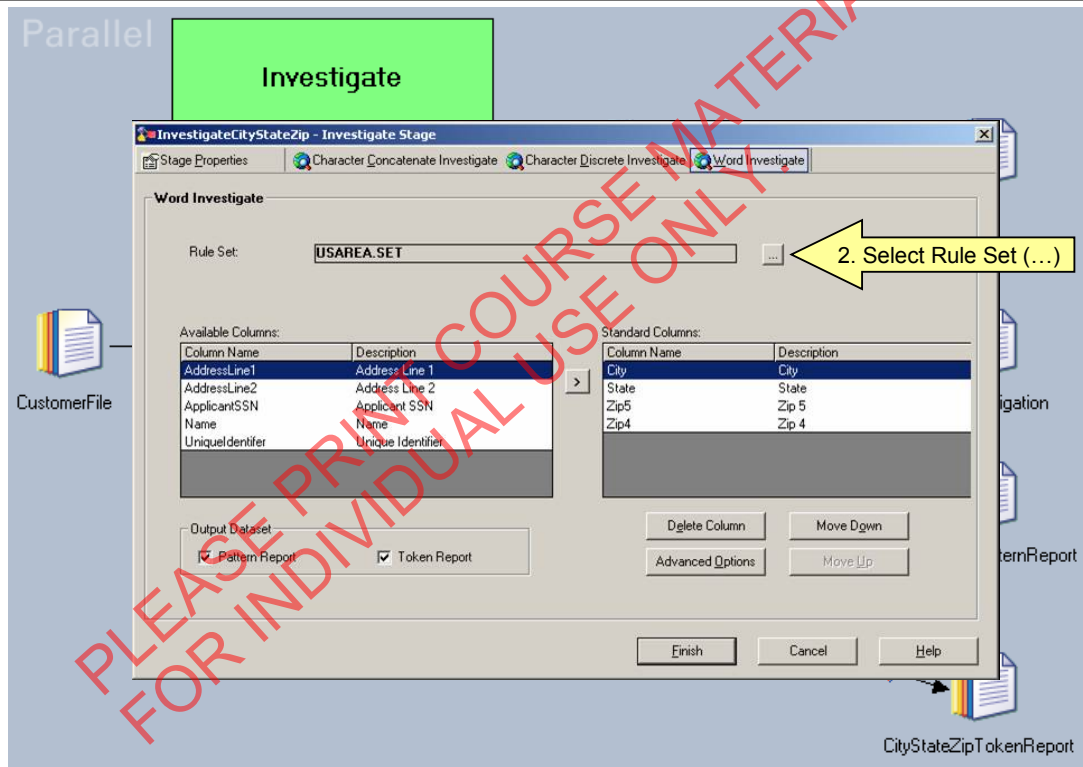
## Investigation - Word



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

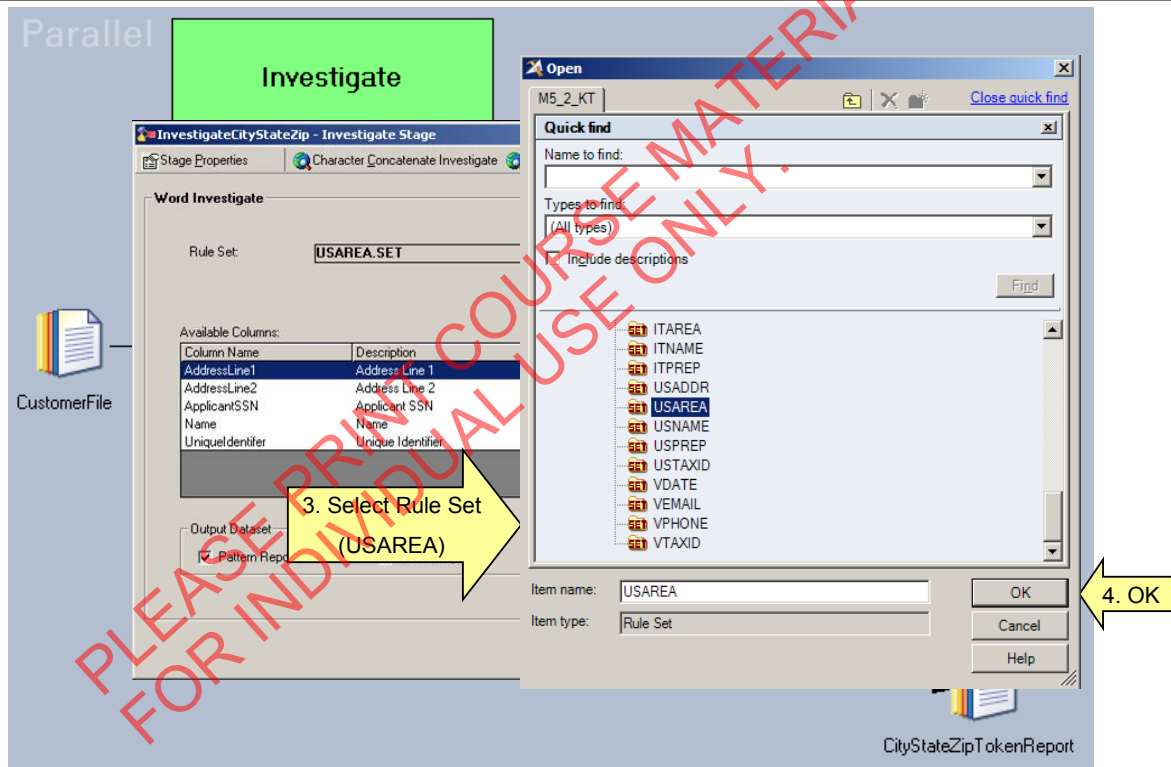
# Investigation - Word



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# Investigation - Word

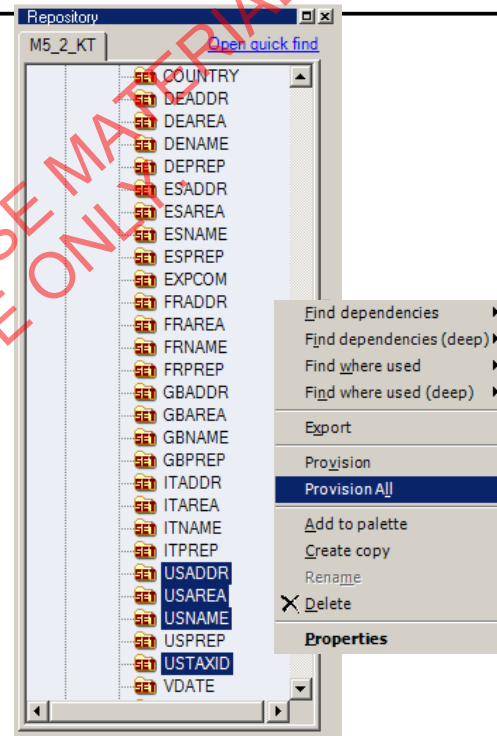


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Provision Rules to be used

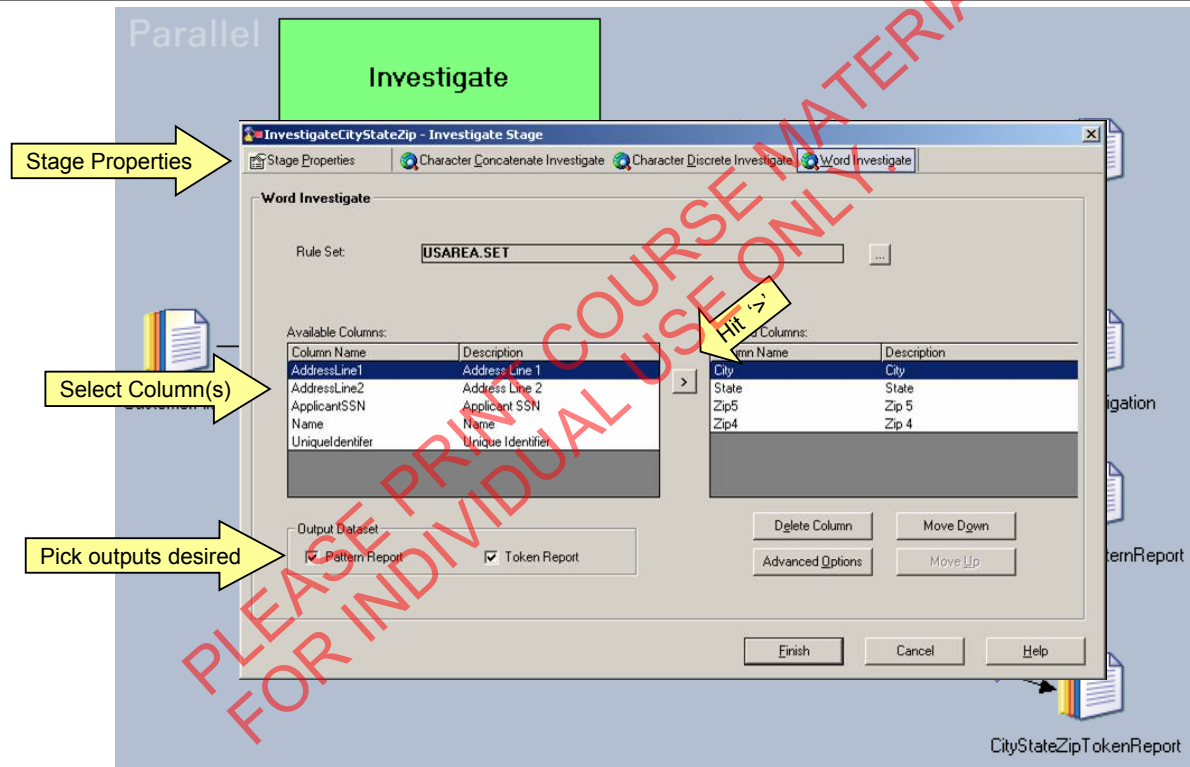
- Provisioning copies rules from repository to execution area
- Use 'Provision all'



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

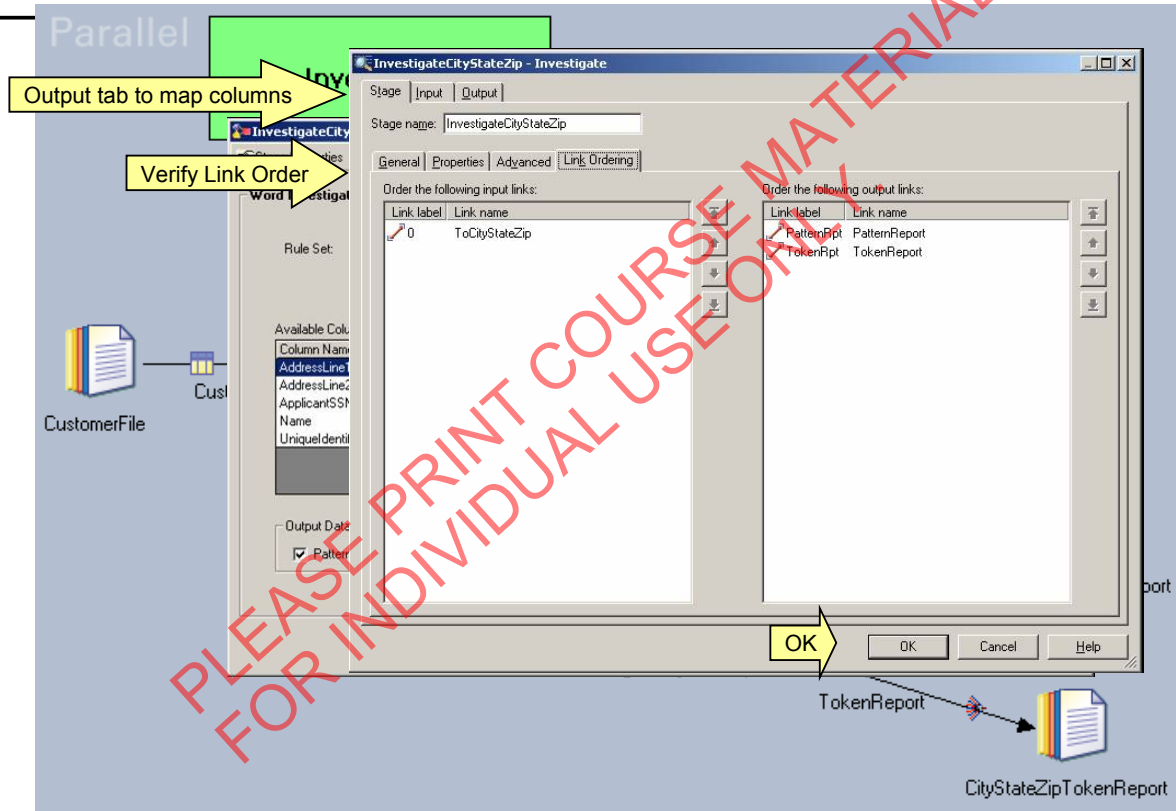
# Investigation - Word



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# Investigation - Word

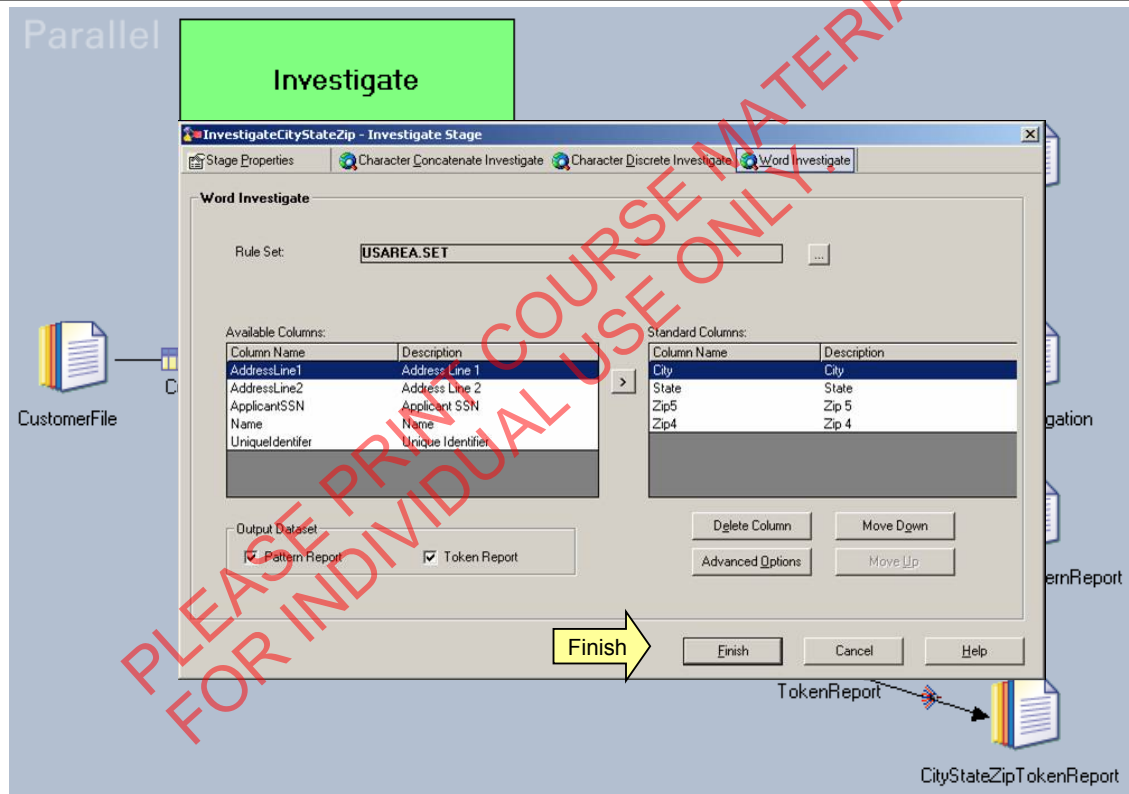


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



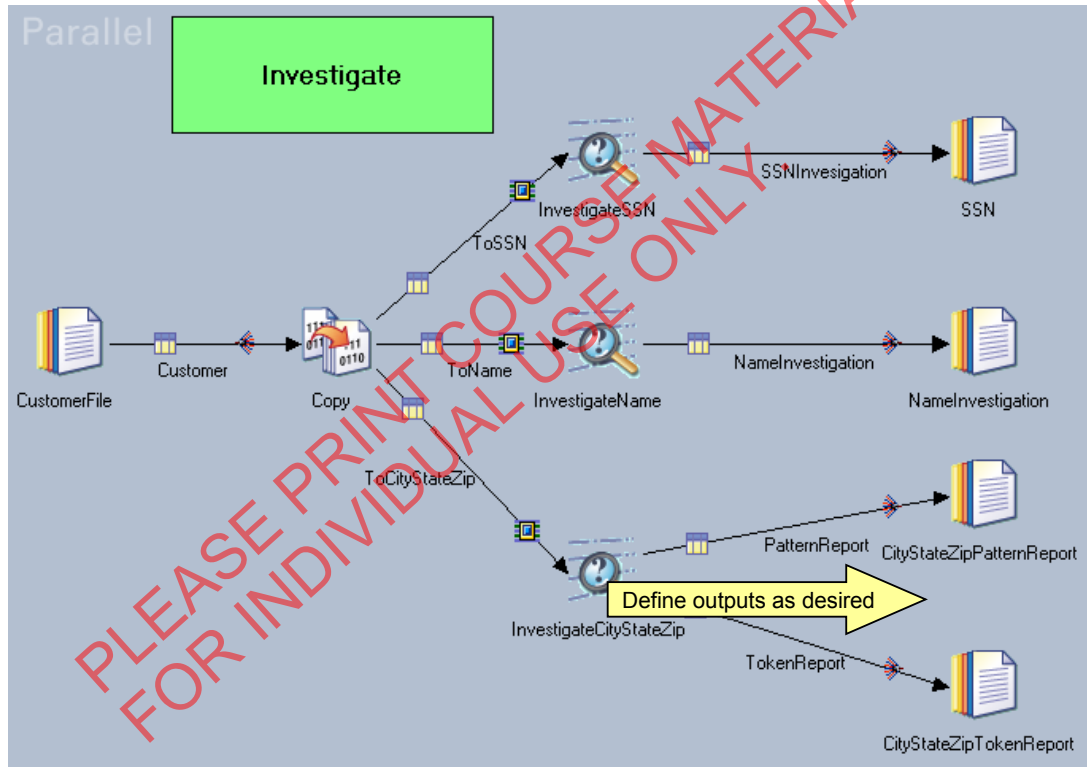
# Investigation - Word



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# Investigation - Word



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

# Investigation - Word

Parallel

Investigate..CityStateZipPatternReport.PatternReport - Data Browser

qsInvColumnName	qsInvPattern	qsInvSample	qsInvCount	qsInvPercent
City+State+Zip5+Zip4	??S^	PIKE ROAD AL 36064	12	2
City+State+Zip5+Zip4	??S^^	PIKE ROAD AL 36064 2634	24	4
City+State+Zip5+Zip4	?S^	MONTGOMERY AL 36111	188	31.3333
City+State+Zip5+Zip4	?S^^	TROY AL 36081 0000	376	62.6667

Investigate..CityStateZipTokenReport.TokenReport - Data Browser

qsInvCount	qsInvWord	qsInvClassCode
573	AL	S
3	ALBERTVILLE	?
1	ALEXANDER	?
26	BANKS	?
26	BIRMINGHAM	?
3	BRUNDIDGE	?
2	CECIL	?
8	CITY	?
4	DOTHAN	?
1	ECLECTIC	?
3	FAIRFAX	?
2	FITZPATRICK	?
6	FL	S
4	GA	S
3	GLENWOOD	?
13	GOSHEN	?
4	GRADY	?
2	GREENVILLE	?
3	GUNTERSVILLE	?
6	HOPE	?
3	HOUSTON	?
6	HULL	?
3	LA	S
4	LAWTON	?
2	LEVEL	?
1	MACON	?
1	MARTINEZ	?
3	MATEO	?
6	MATHEWS	?
3	MELTON	?
1	MILLBROOK	?

CustomerFile

Customer

SSN

Investigation

NameInvestigation

PatternReport

CityStateZipPatternReport

TokenReport

CityStateZipTokenReport

PLEASE PRINT FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Data quality assessment

- Review and analyze each field for the following information:
  - How often is the field populated?
  - What are the anomalies and out-of-range values? How often does each one occur?
  - How many unique values were found?
  - What is the distribution of the data or patterns?
- Use Investigate results to:
  - Update business requirements
  - Define development plan and application design

© Copyright IBM Corporation 2007  
The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint

---

1. (T/F) Character discrete investigation examines a single domain.
2. (T/F) Word investigation examines a single domain.
3. Name the three character masks.

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Checkpoint solutions

1. (T/F) Character discrete investigation examines a single domain.

*Answer: True*

2. (T/F) Word investigation examines a single domain.

*Answer: False*

3. Name the three character masks.

*Answer: C, T, and X*

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Unit summary

Having completed this unit, you should be able to:

- Define data investigation
- Build Investigate stages
- Use character discrete, concatenate, and word investigations to analyze data fields
- Locate and review results

PLEASE PRINT COURSE MATERIALS.  
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 6: Build investigate job

- Character with C mask

**InvestigateCmask - Investigate Stage**

Stage Properties | Character Concatenate Investigate | **Character Discrete Investigate** | Word Investigate

**Character Discrete Investigate**

Available Data Columns:

Column Name	Length	Description
AddressLine1	35	
AddressLine2	35	
City	35	
FullName	46	
RecKey	5	
State	5	
Zip	10	

Add To Selected Columns

Scheduled Process:

Column Name	Mask
SourceSystem	C
PolicyNumber	CCCCCCCCCCCC
FEDID	CCCCCCCCCCCC
DOB	CCCCCCCC
DOD	CCCCCCCC

Delete

Change Mask

Advanced Options

OK Cancel Help

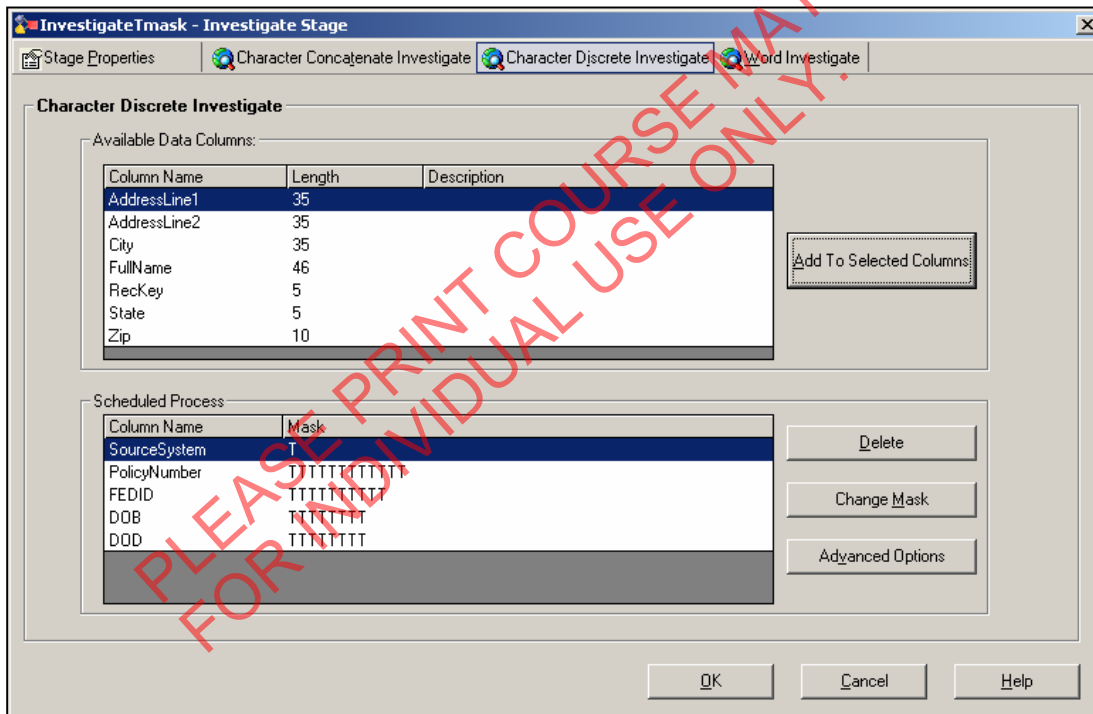
© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.



## Exercise 7: Build investigate job

- Character with T mask

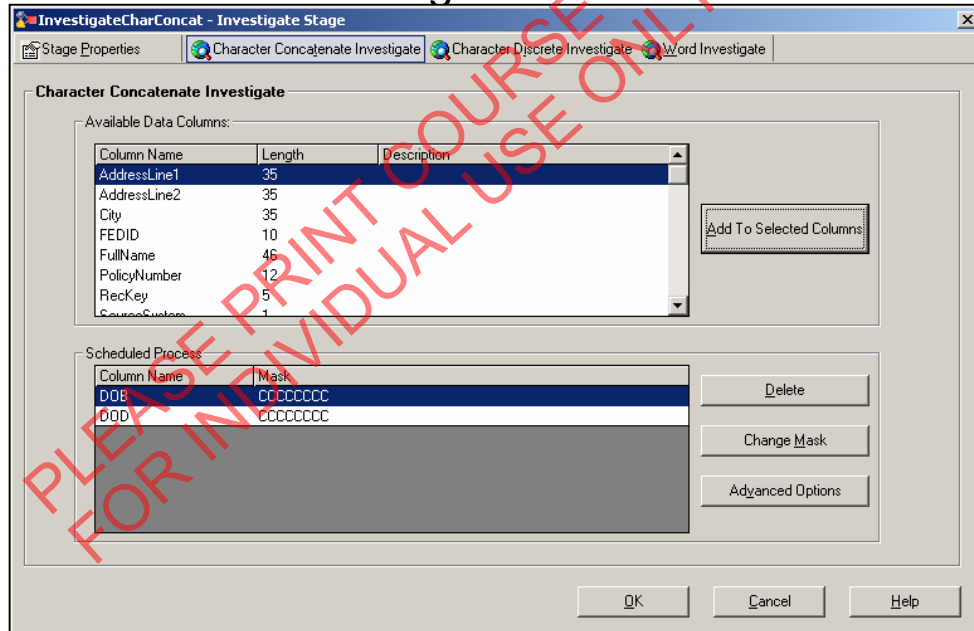


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 8: Build investigate job

- Character concatenate
- Useful for auditing results of other processes, such as standardization and matching

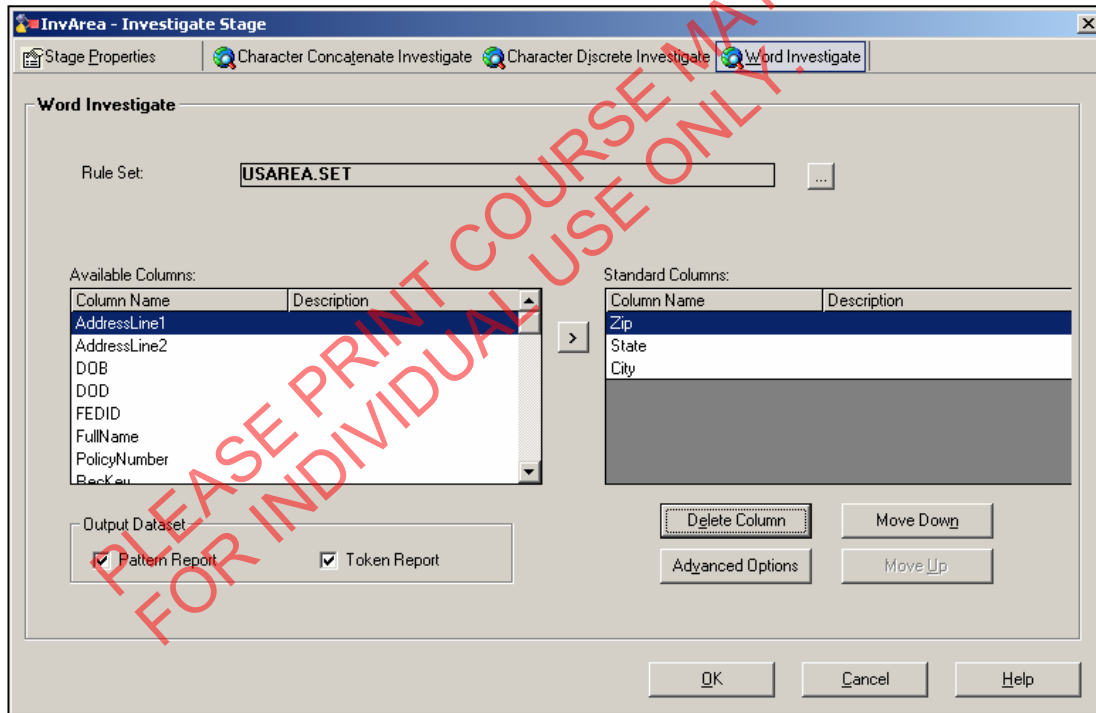


© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

## Exercise 9: Build investigate job

- Word investigation



© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual receiving the materials.  
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.