



PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

Standardize



© 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 Standardize stage in the Data Re-engineering Methodology
 - Identify rule sets
 - Apply the Standardize stage
 - Interpret standardization results
 - Investigate unhandled data and patterns

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.

Standardize

- Transformation
 - Parsing free form fields
 - Comparison threshold for classifying like words
 - Bucketing data tokens
- Standardization
 - Applying standard values and standard formats
- Phonetic Coding for use in Matching
 - NYSIIS
 - Soundex

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.

- decompose free-form fields into single component fields
- assign data to its appropriate metadata fields

Standardize example

Input File:

Address Line 1

1721 W ELFINDALE ST
 1721 W ELFINDALE ST # 20
 16200 VENTURA BOULEVARD
 C/O JOSEPH C REIFF
 1705 W St
 1655 PONCE DE LEON AVENUE

Address Line 2

UNIT 20
 SUITE 201
 12 WESTERN AVE
 PHILADELPHIA
 15TH FLOOR

Result File:

House #	Dir	Str. Name	Type	Unit Type	Unit Value	Floor Type	Floor Value	NYSIIS	Soundex	City
1721	W	ELFINDALE	AVE	UNIT	20					
1721	W	ELFINDALE	ST		20					
16200		VENTURA	BLVD							
12		WESTERN	AVE							
1705	W	ST								PHILADELPHIA
1655		PONCE DE LEON	AVE			FLOOR	15			

© 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.

Notice:

- ▲Unit 20 & # 20, both bucketed as unit type and value
- ▲C/O Joseph C Reiff, moved to another field (can't see on screen)
- ▲12 Western Ave, recognized as an address

Standardize process

Output File →

House Number	Street Name	Street Type	Unit Type	Unit
10	MAPLE	ST	APT	222

Key:

^ = Single numeric

? = One or more unknown alphas

T = Street type

U = Unit type

Process Patterns and Bucket Data

Classify & 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 them. The materials may not be modified, copied, distributed or transferred without the express written consent of IBM Corporation.

Notice the “standardization” of St & APT

1. First we parse the data
2. Classify known words (classification table)
3. Apply general (default tags) to unclassified tokens
4. Create new output fields (dictionary file)
5. Process the patterns (Pattern action File) to move data into correct field and apply standard values and formats.

Standardize stage

- Standardize Stage
 - Uses Rule sets for:
 - Country processing
 - Pre-domain processing
 - USPREP
 - Domain processing
 - USADDR
 - USAREA
 - USNAME
 - Multi-national Address
 - WAVES

PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

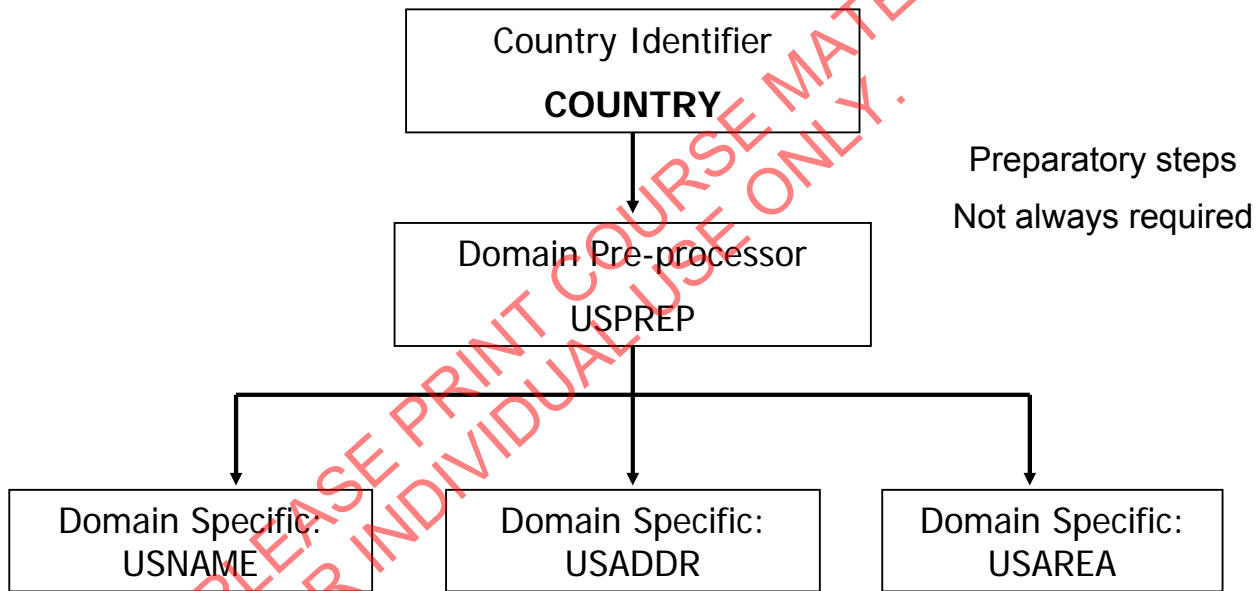
The course materials are provided for internal use only by the individual receiving the materials.

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

- Identify new fields based on underlying data
- Examples:
 - set a gender flag
 - Name type flag - identify individual address from an organization address.
I=Individual, O=Organization
 - Address type flag – S=Street address, B= box address, R=rural route address, O=other type of address
 - Transformation rules are created both for matching and creating the load file (sometimes these rules are different)

WAVES (World-wide address verification) performs address validation against a database

Types of 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.

Three levels of rule sets

- Country Identifier - Identify the country and append ISO country code based on address format
- Domain Pre-processing – separates mixed domain fields into specific domain fields
- Domain Specific

Example: country identifier

Input Record

100 SUMMER STREET 15TH FLOOR BOSTON, MA 02111
 SITE 6 COMP 10 RR 8 STN MAIN MILLARVILLE AB TOL 1K0
 28 GROSVENOR STREET LONDON W1X 9FE
 123 MAIN STREET

Output Record

US Y 100 SUMMER STREET 15TH FLOOR BOSTON, MA 02111
 CA Y SITE 6 COMP 10 RR 8 STN MAIN MILLARVILLE AB TOL 1K0
 GB Y 28 GROSVENOR STREET LONDON W1X 9FE
 US N 123 MAIN STREET

© Copyright IBM Corporation 2007

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

Note: US default is assumed using ZQUSZQ country code delimiter

ISO Country Code

Indicator Flag (Y or N, where Y = country code verified)

The format of the country code delimiter is:

- ZQ<ISO Country Code>ZQ

For example, the country code delimiter for the United States would be:

- ZQUSZQ

Allows easy processing of multi-national data

Assigns an ISO country code to each record

When the Country Identifier Rule Set can not determine the country code, the default value will be taken from the country code delimiter

Example: domain preprocessor

Input Record

Field 1	JIM HARRIS (781) 322-2426
Field 2	92 DEVIR STREET MALDEN MA 02148

Mixed domain

Output Record

Name Domain	JIM HARRIS
Address Domain	92 DEVIR STREET
Area Domain	MALDEN MA 02148
Other Domain	(781) 322-2426

© Copyright IBM Corporation 2007

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

Categorize input data into one of the following domain-specific column sets:

- Name - individual and organization names, attention instructions, and secondary names
- Address - low level geography including street, rural, box, unit, and building addresses
- Area - high level geography including city name, postal code, and country code
- Other - non-name and non-geography data

The metadata delimiters indicate what kind of information you are expecting to find within the fields of your input data

If the pre-processor can not determine the domain of a token, it will be defaulted based on its metadata delimiter

The format of the metadata delimiter is:

- ZQ<Domain>ZQ

There are four accepted metadata delimiters:

- ZQNAMEZQ - Name delimiter
- ZQADDRZQ - Address delimiter
- ZQAREAZQ - Area delimiter
- ZQOTHRZQ - Other delimiter

Example: domain specific

Input Record

100 SUMMER STREET 15TH FLOOR

Output Record

House Number	100
Street Name	SUMMER
Street Suffix Type	ST
Floor Type	FL
Floor Value	15
Address Type	S
NYSIS of Street Name	SANAR
Reverse Soundex of Street Name	R520
Input Pattern	^+T>U

© Copyright IBM Corporation 2007

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

Evaluate domain-specific input

Generate business intelligence fields:

- Create all subordinate domain elements needed for data storage and presentation
- Apply consistent representations to data
- Incorporate applicable standards such as postal standards for addresses

Generate matching fields:

- Blocking keys
- Primary match keys

Rule sets

- Rule sets contain logic for:
 - Parsing
 - Classifying
 - Processing data by pattern and bucketing data
- Three required files
 - Classification Table
 - Dictionary File
 - Pattern Action File
- Optional files
 - Lookup tables
 - Override tables

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.

Rule set files

Classification Table

Contains standard abbreviations that identify and classify key words.

Dictionary File

Define the output file fields to store the parsed and conditioned data

Pattern Action File

Contains a series of patterns and programming commands to condition the data

Reference Tables

Optional conversion and lookup tables for converting and returning standardized values

Override Tables

Tables for storing overrides entered into the Designer GUI

© 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 first three files are required for each rule set.

If the rule set name is USNAME, then the classification table is named USNAME.CLS, USNAME.PAT, USNAME.DCT

Rule sets can be copied.

Dictionary File: Defines two-character field abbreviations that are used in the output file for a particular rule set.. Each field is referenced in the PAT file to determine where individual tokens will be bucketed.

First Name Lookup table: Applies an enhanced first name, e.g. Barbara to Barb or Barbie, Kenneth to Ken, Kathleen to Kathy.

Example:

Classification table ST. is classified as a street type (T) with a standard value of (ST)

Classification table

- Contains the words for classification, standardized versions of words, and data class
- Data class (data tag) is assigned to each data token
- Default classes are the same across all rule sets
- User-defined classes are assigned in the classification table
 - Users may modify, add or delete these classes
 - User-defined classes are a single letter

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.

Default classes

Class	Description
^	A single numeric
+	A single unclassified alpha (word)
?	One or more consecutive unclassified alphas
@	Complex mixed token, e.g., C3PO
>	Leading numeric, e.g., 6A
<	Trailing numeric, e.g. A6
Zero (0)	Null class

© 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 classes are the same across all rule sets.

If any of the characters in the class column actually occur in the data they are represented by a ~. A ~ also represents itself.

Null does not mean database null. It means the PAT file will delete the string.

User-defined classes

Class	Description
USNAME	
G	Generational, e.g., Senior, I, II
P	Prefix, e.g. Dr., Mr., Miss
USADDR	
T	Street Type
D	Directional
B	Box Type
USAREA	
S	State Abbreviation

© 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.

User defined classes are specific to the rule set.

Classification table example

```
; USADDR Classification Table
```

```
;
```

```
; Classification Legend
```

```
;
```

```
; B - Box Types
```

```
; D - Directionals
```

```
; F - Floor Types
```

```
; H - Highway Modifiers
```

```
; R - Rural Route, Highway Contract, Star Route
```

```
; T - Street Types
```

```
; U - Unit Types
```

```
;
```

```
; Table Order: 51-51 Ascending, 26-50 Ascending, 1-25 Ascending
```

```
;
```

Token	Standard form	Classification
DRAW	"PO BOX"	B
DRAWER	"PO BOX"	B
PO	"PO BOX"	B
POB	"PO BOX"	B
POBOX	"PO BOX"	B
POBX	"PO BOX"	B
PODRAWER	"PO BOX"	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.

Comparison threshold

- May be used in the Classification table
- Used to efficiently make entries into the classification table
- Helps overcome spelling and data entry errors
- Not required
- Threshold uses a logical string comparator

Threshold level	
900	Exact match
850	Almost certainly the same
800	Most likely equivalent
750	Most likely not the same
700	Almost certainly not the same

© 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.

Two passes are made through the classification table:

First pass looks for an exact match

Second pass looks for a fuzzy match based on the threshold level

Classification table example with comparison threshold

```

; USADDR Classification Table
;-----
; Classification Legend
;-----
; B - Box Types
; D - Directionals
; F - Floor Types
; H - Highway Modifiers
; R - Rural Route, Highway Contract, Star Route
; T - Street Types
; U - Unit Types
;-----
; Table Sort Order: 51-51 Ascending, 26-50 Ascending, 1-25 Ascending
;-----
DRAW          "PO BOX"          B
DRAWER        "PO BOX"          B
.....
NORTHEAST     NE                 D 850
NORTHWEST     NW                 D 850
NW            NW                 D
S             S                 D
SO            S                 D
SOUTH         S                 D

```

© 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.

850 – almost certainly the same

Dictionary file

- Defines the field definitions for the output file
- When data is moved to these output fields it is called “bucketing” the data
- The order that the fields are listed in the dictionary file defines the order the fields appear in the output file
- Dictionary file entries are similar to field definitions

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.

Dictionary file example

```

;;QualityStage v8.0
\FORMAT\ SORT=N
;-----
; USADDR Dictionary File
;-----
; Total Dictionary Length = 411
;-----
; Business Intelligence Fields
;-----
HouseNumber C 10 S HouseNumber ;0001-0010
HouseNumberSuffix C 10 S HouseNumberSuffix ;0011-0020
StreetPrefixDirectional C 3 S StreetPrefixDirectional ;0021-0023
StreetPrefixType C 20 S StreetPrefixType ;0024-0043
StreetName C 25 S StreetName ;0044-0068
StreetSuffixType C 5 S StreetSuffixType ;0069-0073
StreetSuffixQualifier C 5 S StreetSuffixQualifier ;0074-0078

```

© 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.

Pattern-Action file

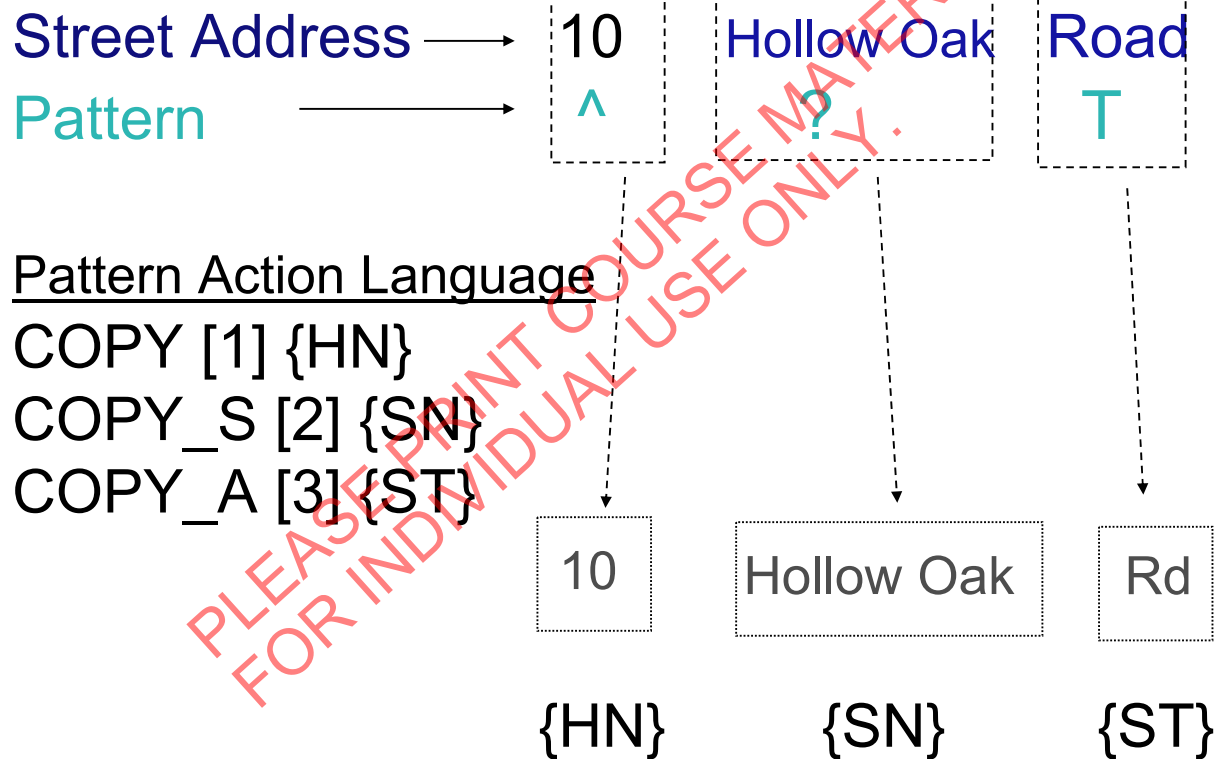
- Contains the rules for standardization; that is, the actions to execute with a given pattern of tokens
- Records are processed from the top down
- Written in Pattern Action Language (PAL)
- Complex parsing can be coded in this file

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.

Pattern Action file process



© Copyright IBM Corporation 2007

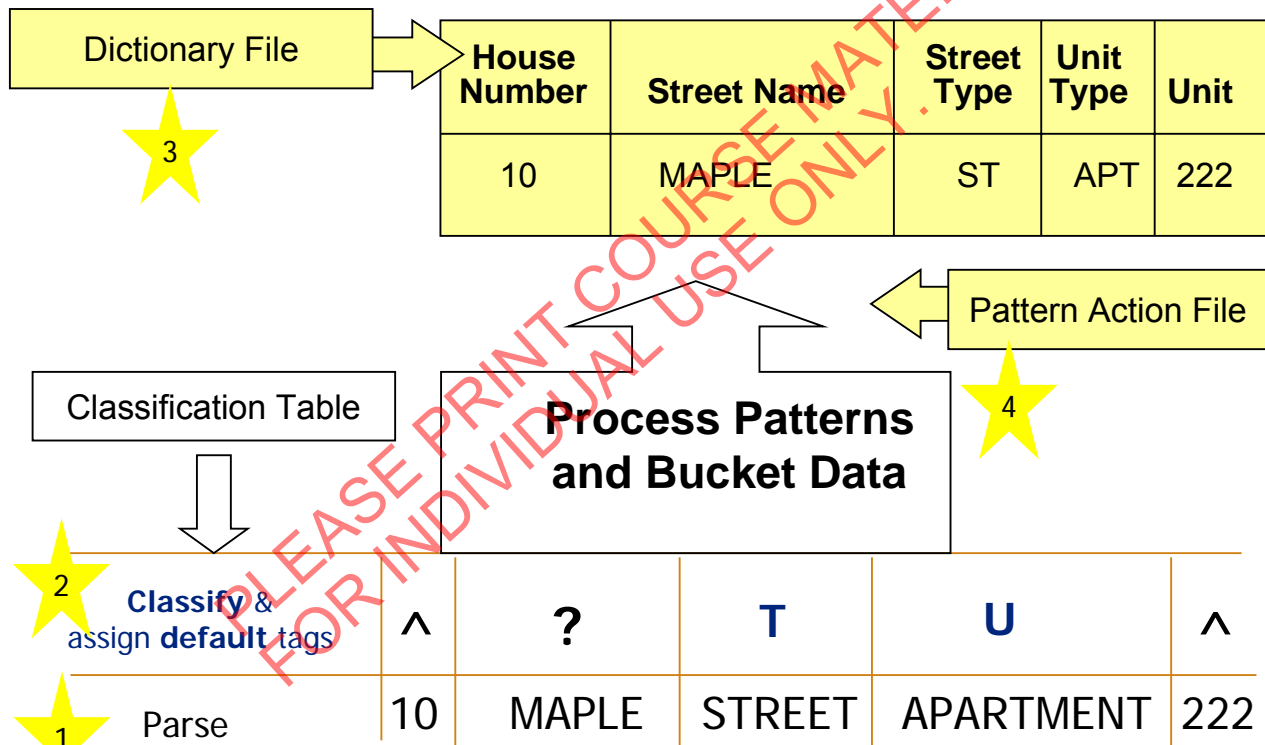
The course materials are provided for internal use only by the individual receiving the materials.
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Optional lookup tables

- Called from the Pattern Action File
- Rule sets may contain lookup tables such as:
 - Common First Names and Enhanced First Names
 - Barb & Barbara
 - Ted & Edward
 - Gender based on name
 - State abbreviations
 - Common city abbreviations
 - NYC = New York City
 - LA = Los Angeles

© 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.

Standardize process



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

1. (Pattern Action file) parses the data
2. (classification table) Classifies known words
3. (dictionary file) Defines new output fields
4. (Pattern action File) Processes the patterns to move data into correct field and applies standard values and formats.

Standardizing international data

- Two methods
 - Method 1: Use country pre-processor, domain pre-processor, and domain-specific rules
 - Uses out-of-the-box, included functionality/rules
 - Method 2: Use Multinational Standardize and WAVES
 - Requires purchase of WAVES database

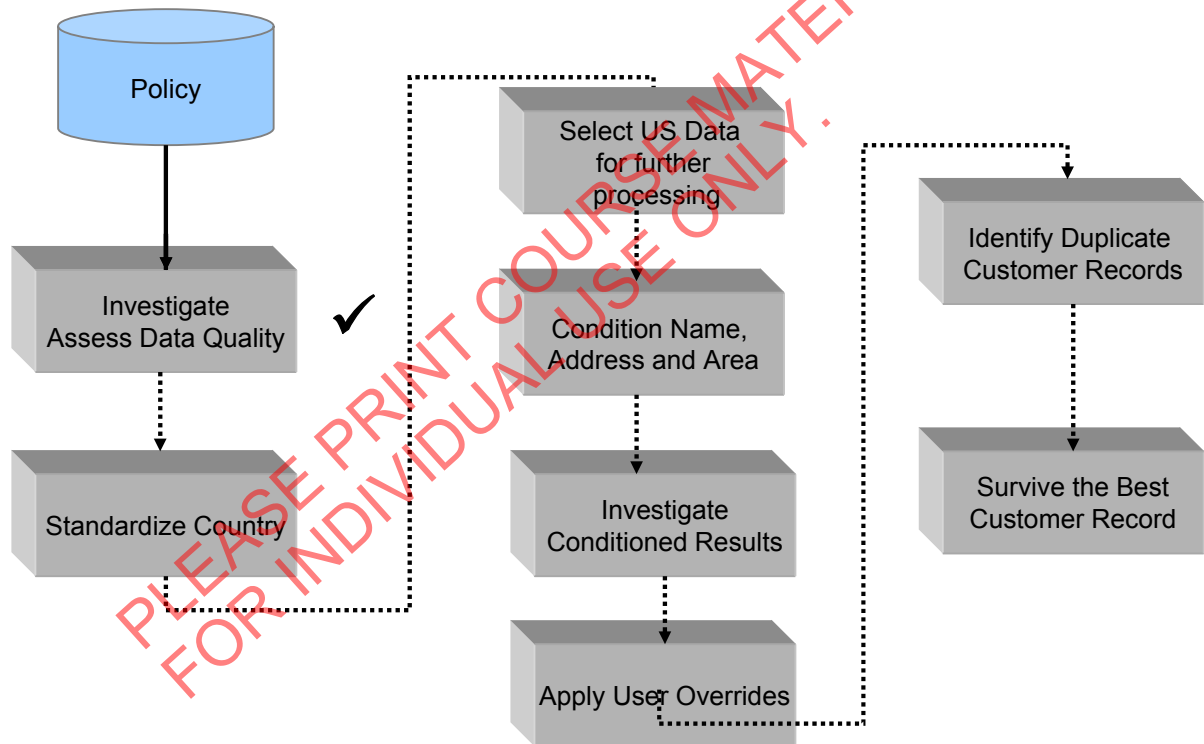
PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

© 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.

WAVES first uses the Multinational standardization, then validates the result against a database.

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.

Country rule set

- Country Rule set appends the two byte ISO country code
- Input to the country rule set includes:
 - Street Address
 - City or locality
 - State
 - Zip or Postal code
 - Country field (if it exists)
- Output:
 - Two byte ISO country code
 - Flag identifying explicit or default decision

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.

Standardization implementation

PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

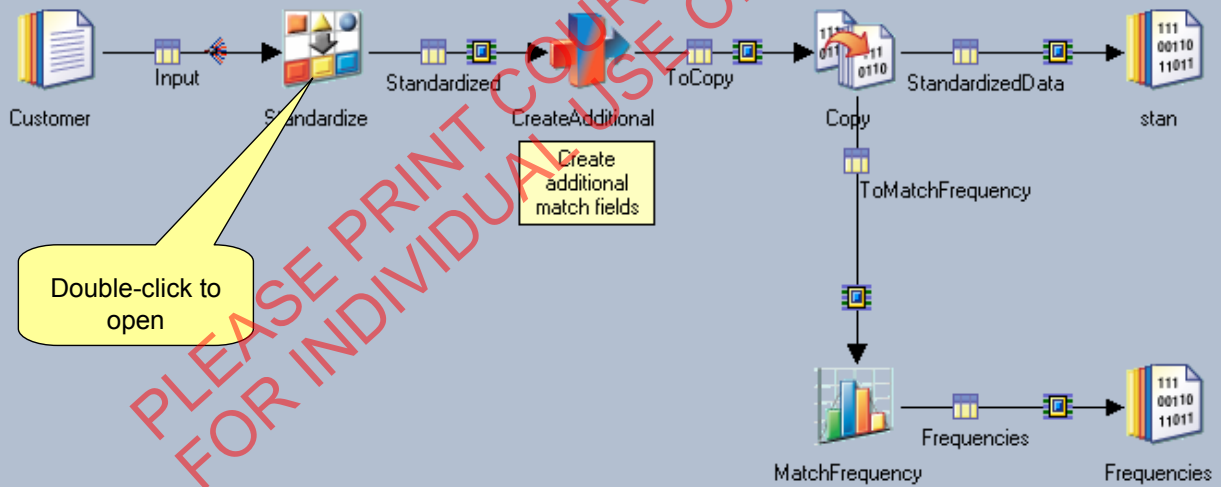
© Copyright IBM Corporation 2007

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.

Standardization

Parallel

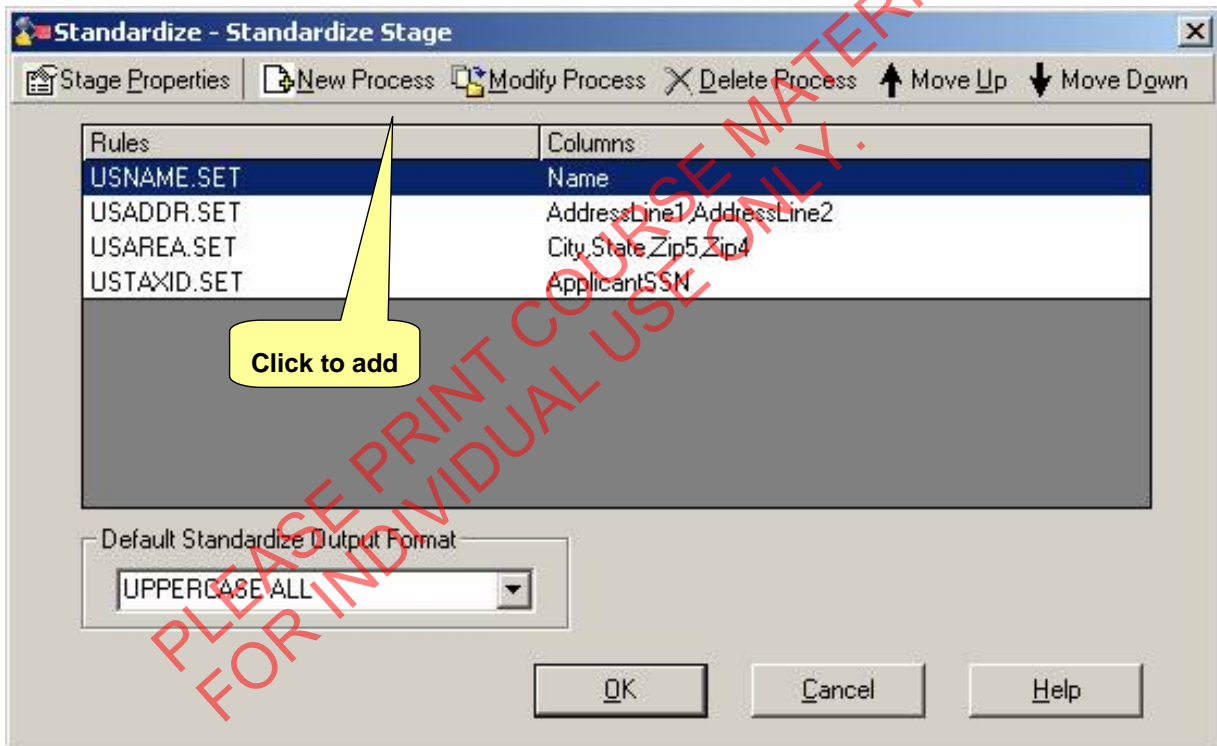
Standardize and generate match frequencies



© 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.

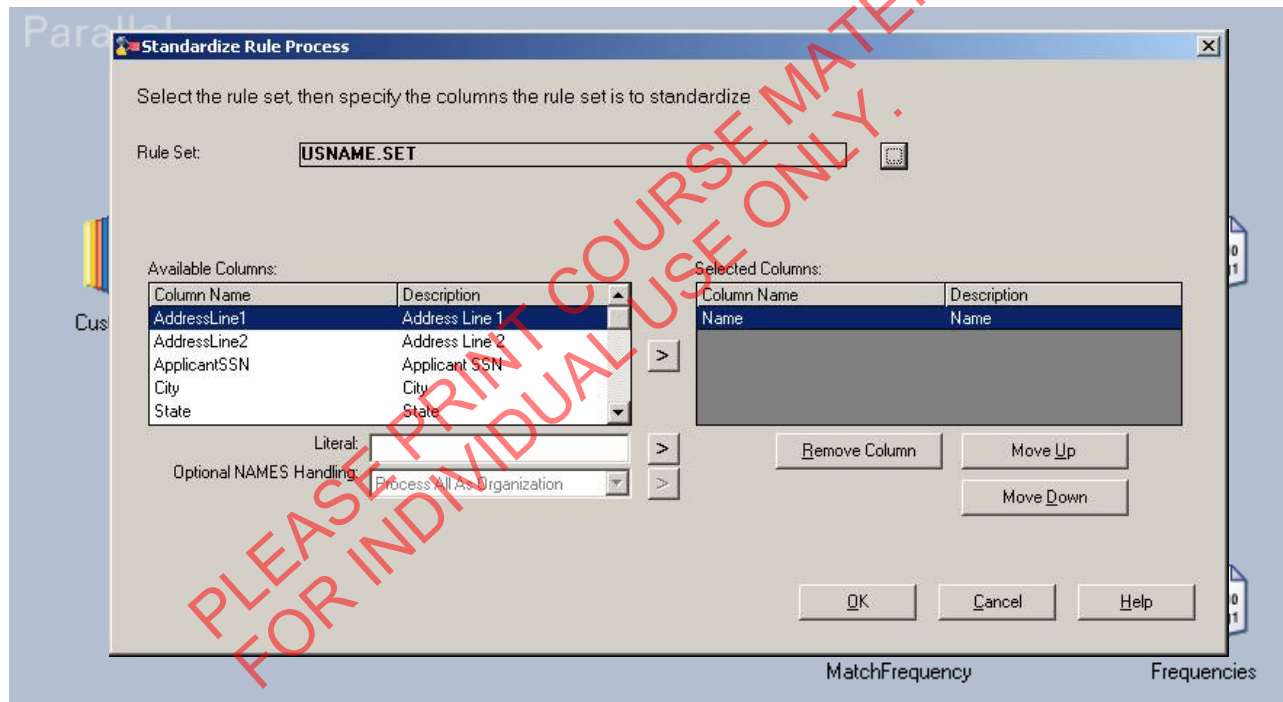
Standardization



© 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.

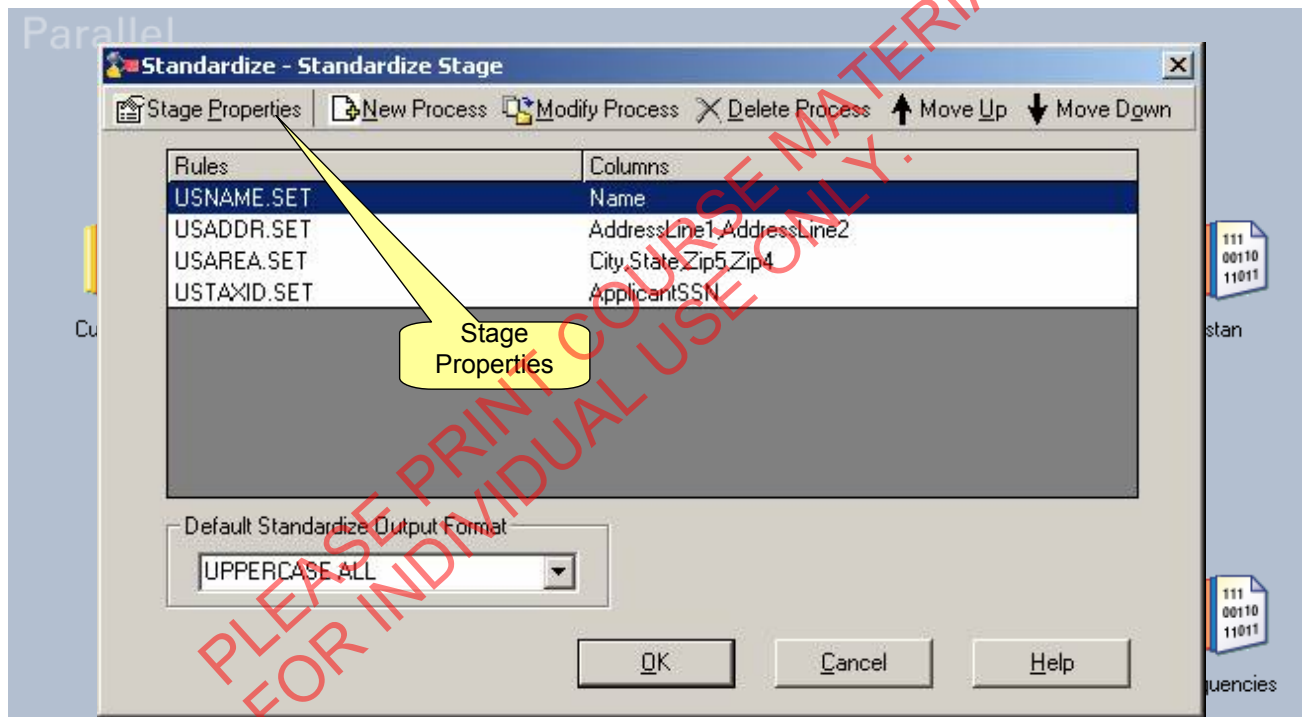
Standardization - Word



© 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.

Standardization



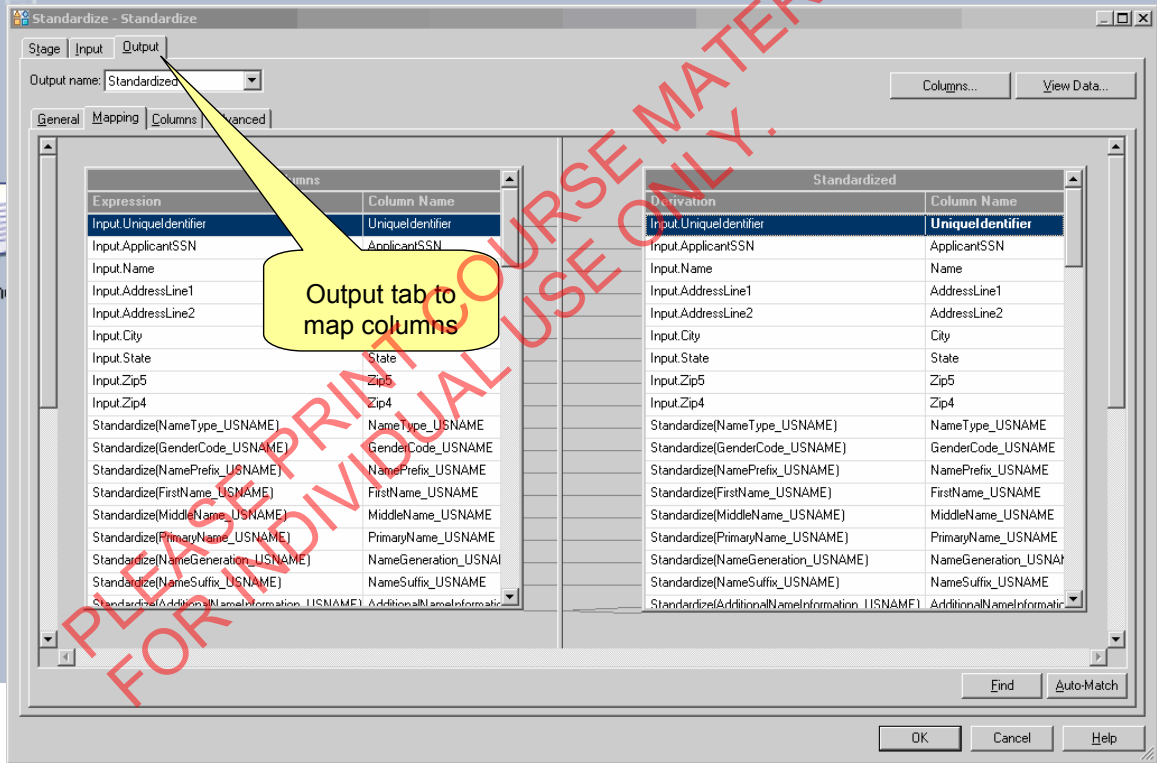
© Copyright IBM Corporation 2007

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

Parallel

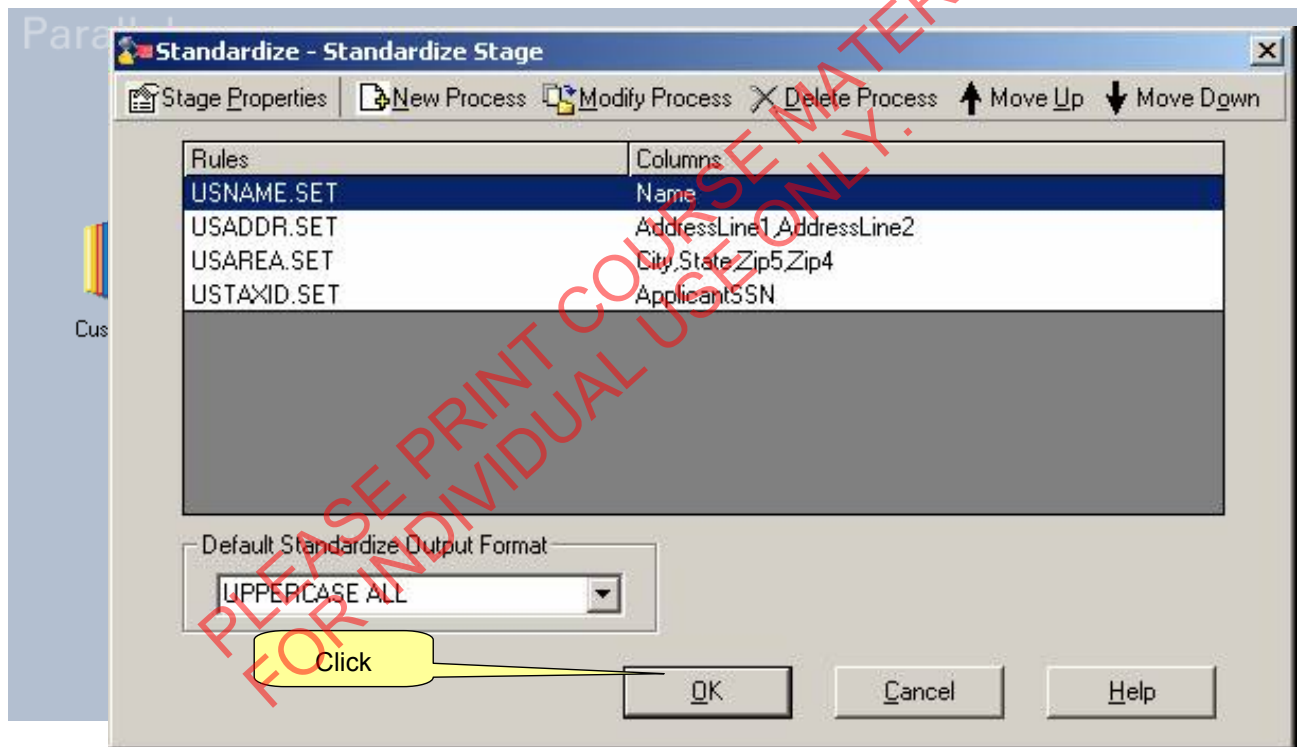
Custom



© 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.

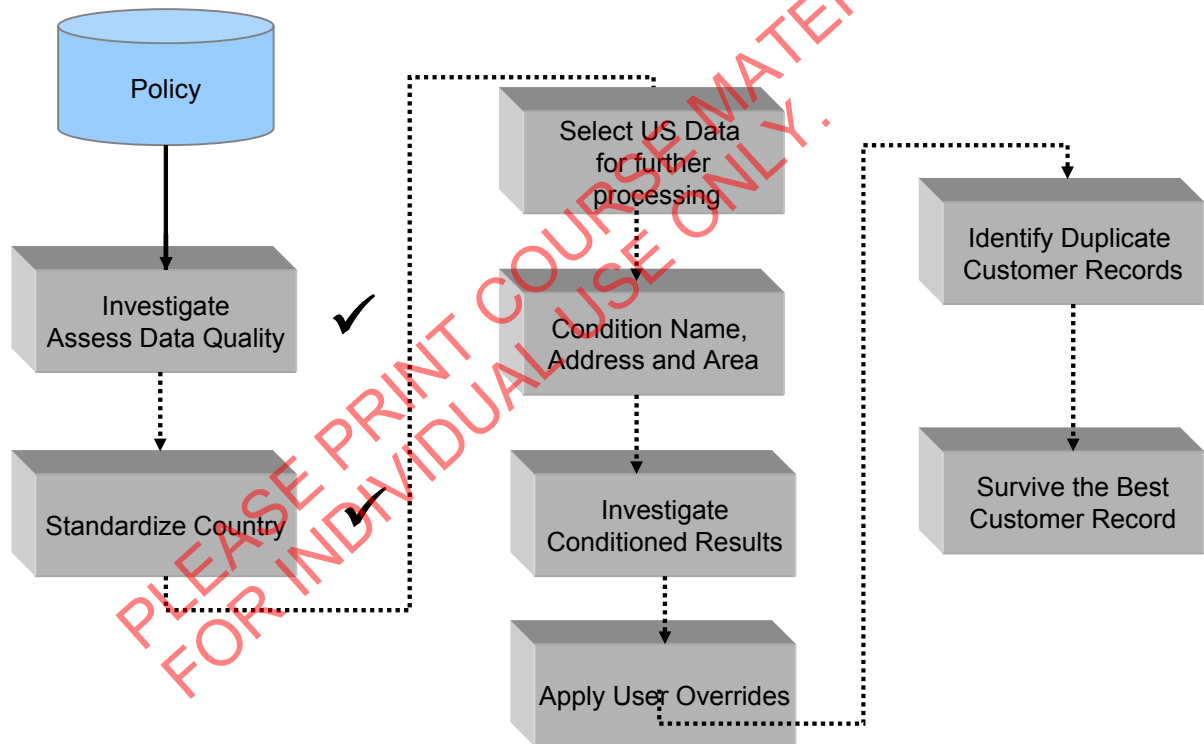
Standardization



© 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.

Course exercise project design

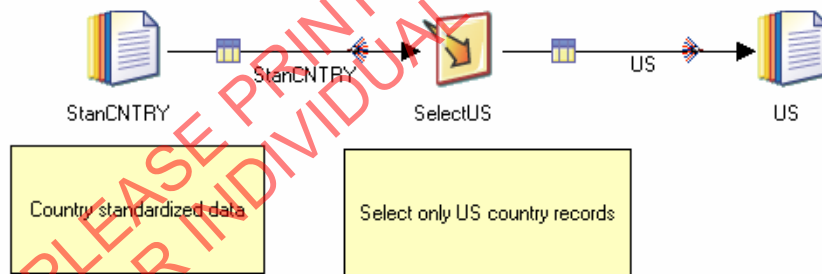


© Copyright IBM Corporation 2007

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.

Selecting US data

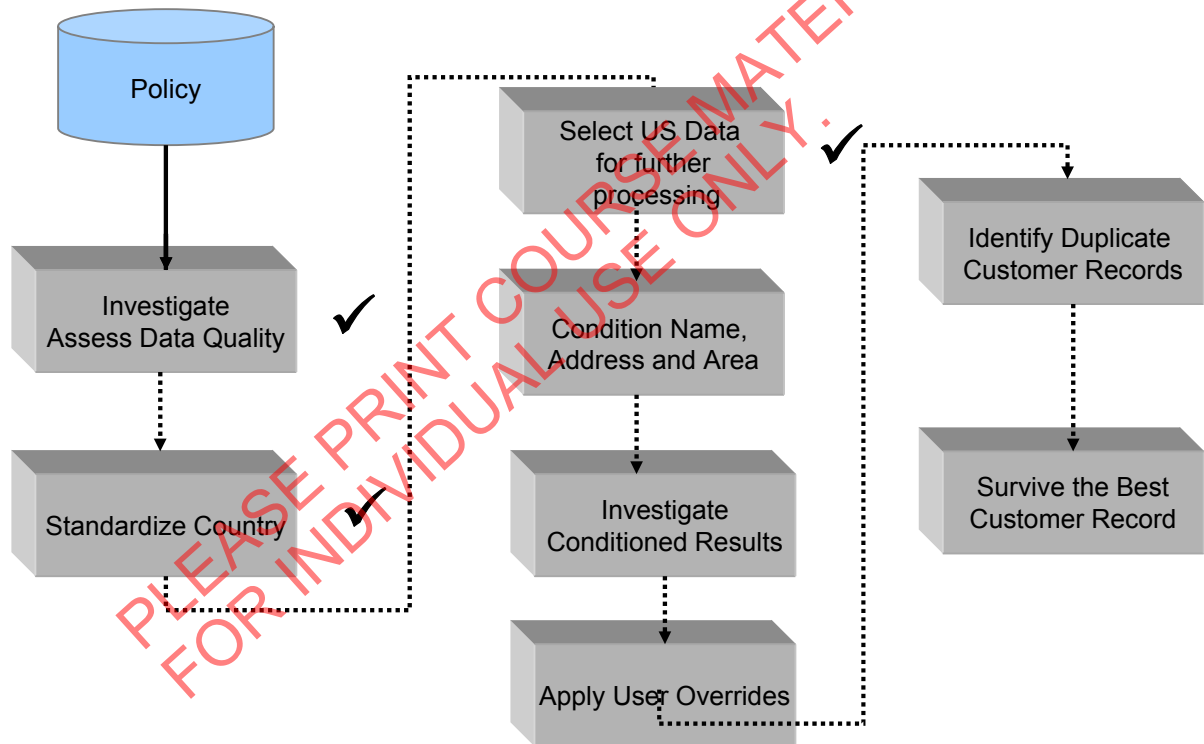
- The DataStage Filter Stage provides the capability of selecting and/or rejecting records based on a set of values for a field
- Selecting or splitting data requiring compound or complex logic may require Transformer stage



© 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.

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.

Domain pre-processor rule sets

- Pre-processor rule sets are designed to filter name, street address and area (city, state, zip) data
 - For example, if the city, state and zip is found in ADDRESS LINE 2, the pre-processor rule set will attempt to recognize this data and move it into the area domain
- The pre-processor rule set prepares the data for processing by domain specific 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.

Domain rule sets

- Domain rule sets expect only data for that domain as the input
- Domain rule sets that come with QualityStage are:
 - Name
 - Street address
 - Area (city, state and zip)

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.

USNAME rule set

- The USNAME rule set works on both personal names and organization names for US data
- Data is parsed into name components
- Phonetic coding of the First Name and Primary Name are created for matching

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.

See USNAME.DCT file for all parsed names (PUT COPY OF USNAME.DCT, CLS KEY, PAT PARSING REQ. IN SG)

A "Name Type" Flag is applied to help identify individual names vs. organization names which may require different match strategies

Data is parsed into individual name fields

USADDR rule set

- This rule set is applied to street address fields
- The “Address Type” flag identifies different types of addresses
 - ‘S’ Street address
 - ‘B’ Box address
 - ‘R’ Rural route address
- Phonetic coding of the Street Name is created for matching

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.

Note: (PUT COPY OF USADDR.DCT & CLS KEY & PARSING REQ FROM PAT FILE IN SG)

Address Type Flag: S=street address (street name is populated), B=Box address (no street name but a box type is populated), R=Rural Route (no street name, no box type, but the rr type is populated).

Different types of addresses may require different match strategies.

Data is parsed into individual street address fields

USAREA rule set

- This rule set is applied to city, state and postal code fields
- Data is parsed into city name, state abbreviation, zip code and zip plus four
- Phonetic coding of the city name is created for matching

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.

Standardize results

- Business Intelligence fields
 - Parsed from the original data, they may be used in matching and generally they are moved to the target system
- Matching Fields
 - Generally these fields are created to help during the match process and are dropped after successful matching
- Reporting fields
 - Specifically created to help review results of Standardize and recognized handled and unhandled 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.

Field types are documented in the dictionary file of the rule set.

Business intelligence fields

- Intelligent data parsed and bucketed from the input free-form field

USNAME Examples

- Title
- First Name
- Middle Name
- Primary Name
- Generational

USADDR Examples

- HouseNumber
- Directional
- Street Name
- Unit Types
- Box Types
- Unit Values
- Building Names

USAREA Examples

- City
- State
- Zip5
- Zip4

© 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.

Fields are parsed from the source data.

Standardize matching fields

- Phonetic coding
 - NYSIIS
 - Reverse NYSIIS
 - Soundex
 - Reverse Soundex
- Hash keys
 - First 2 characters of the first five words
- Packed Keys
 - Data concatenated, or packed

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.

Fields used to facilitate matching.

Soundex field somewhat more stable for reverse phonetics.

Phonetics usually used for name fields.

Hash keys and packed keys are name-specific.

Standardize reporting fields

Unhandled Pattern

The pattern generated for tokens not processed by the selected rule set.

Unhandled Data

The remaining tokens not processed by the selected rule set.

Input Pattern

The pattern generated for the stream of input tokens based on the parsing rules and token classifications.

Exception Data

The tokens not processed by the rule set because they represent a data exception.

User override flag

© 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 components are present in every domain-specific rule set.

Investigate NAME unhandled patterns and data

- Identify the unhandled patterns for the NAME field. In the report include the unhandled data, input pattern, original data and the record key.
 1. Build a Character Concatenate Investigation using the following fields
 2. Increase the number of samples to 5

Field Name	Field Description	Type
UPUSNAM	Unhandled Pattern	C
UDUSNAM	Unhandled Data	X
IPUSNAM	Input Pattern	X
NAUSPRE	Name domain data	X

© 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 w

Standard practice: investigate handled and unhandled data

- Review the business intelligence fields to ensure accurate bucketing of the data
- Build a Character Discrete Investigation for each field and review the contents and the format
- Build Investigation to review:
 - Unhandled Patterns
 - Unhandled Data
 - Input Pattern
 - Input Fields

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.

Example: Ensure that the House Number field (HNUSADD) contains only numeric data

Is the House Number field always blank?

Directions fields contain: N, S, E, W, or NW, NE, SW, SE

You may complete some quick visual inspections, the Investigation reports allow you to quantify the changes and improvements

Customizing rule sets

- A rule set may require modification if some input data is:
 - Not processed
 - Incorrectly processed
- QualityStage provides functions to:
 - Modify classification table
 - Apply user Overrides
 - Test strings for classifications using the Rules Analyzer

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.

Modify classification table

- In repository
 - Copy the rule set
 - Modify the copy

To modify:
USNAME.CLS

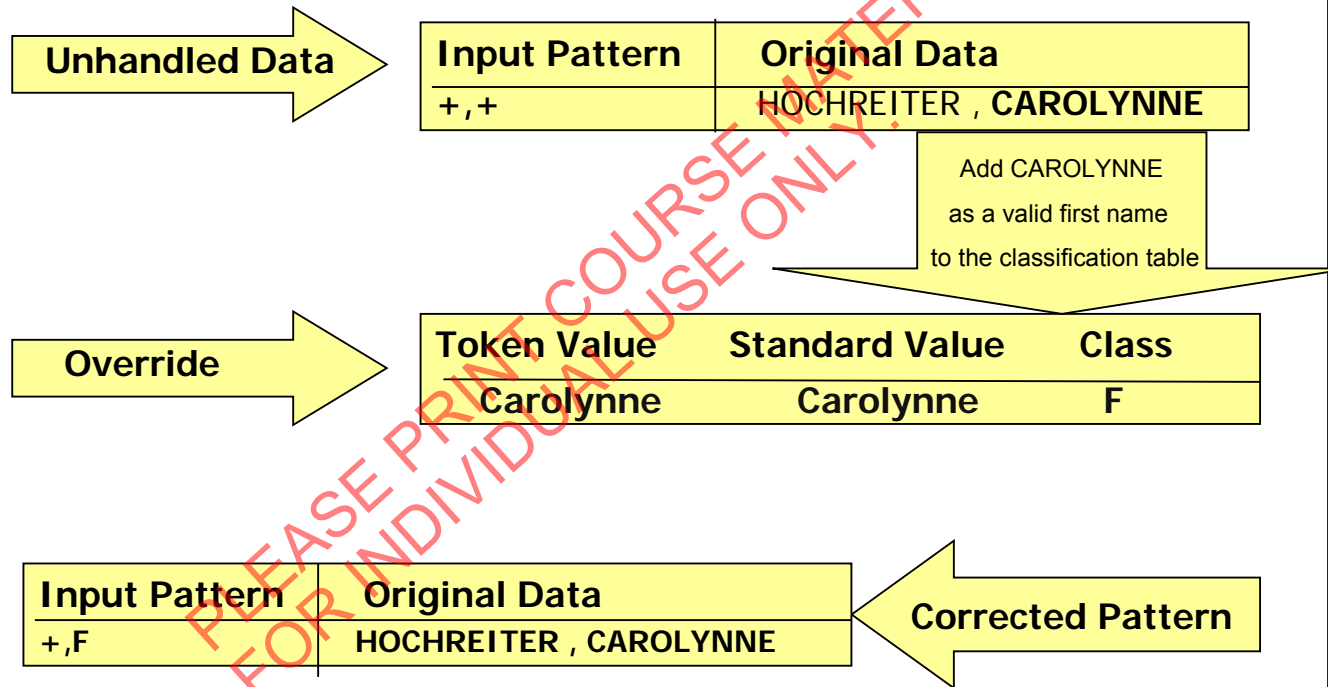
```

USNAME.CLS - Notepad
File Edit Format View Help
;QualityStage v7.0
\FORMAT\ SORT=Y
-----
USNAME Classification Table
-----
Classification Legend
-----
A - Abbreviations (Misspellings)
C - Common words
F - First Names
G - Individual Name Generations
I - Initials
L - Last Name Prefixes
O - Organization Name Suffixes
P - Individual Name Prefixes
Q - Additional Name Qualifiers
S - Individual Name Suffixes
W - Organization Name words
Z - Delimiters
-----
Table Sort Order: 51-51 Ascending, 26-50 Ascending, 1-25 Ascending
-----
END          ENDOWMENT          W
AN           AN                C
AND          AND               C
AS           AS                C
AT           AT                C
BY           BY                C
FOR          FOR               C
FROM         FROM             C
IN           IN                C
OF           OF                C
  
```

© 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.

Apply classification override



© 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 word (alpha) Carolynne is not recognized as a First Name. Review the Name Word INV frequency report and note that Carolynne appears 5 times in the data. This frequency influences the decision to add Carolynne as a Classification Override so that it is recognized as a first name.

The input pattern before the override is: +, + → unknown alpha comma unknown alpha

After the override the pattern is +, F → unknown word, comma first name

The second pattern is recognized and processed by the pattern action file.

User overrides

- Provides the user with the ability to modify rule sets
- The following types of rule sets can be modified using User Overrides
 - Domain Pre-processor rule sets
 - Domain rule sets
- There are five types of user overrides relating to: classifications, patterns, and text strings
- User overrides are
 - GUI Driven
 - Stored in separate lookup tables
- Rule set should be provisioned after modifications applied

© 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.

Provides the user the ability to specify their own Standardize rules

User overrides are GUI-driven

The user does not need to know pattern-action language syntax

The user does not need to edit the classification table or the pattern-action file

Overrides require the following information:

- Dictionary field name to move the token to
- Original value or standard value of token
- Leading space or no leading space for multiple tokens moved to the same dictionary field

User classification override

- Recognized as a keyword and classified
 - Additional words
 - New abbreviation, variation
 - Misspelling of a word
- User Classifications may override or add:
 - Original values (Token values)
 - Standard value
 - Class

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.

Example:

Carolynne is not recognized as a first name. Since the name carolynne occurs 5 times in the data (review word inv report on names, word frequency count). We might want to add this name to the classification table so that it is recognized.

Use the Word Investigation Word Frequency reports to check the frequency that a word, abbreviation or misspelling occurs.

Classification overrides take precedence over the classification table

Classification overrides are available in both domain pre-processor rule sets and domain-specific rule sets

Text overrides

- Allow the user to specify overrides based on an entire text string
- Use this override for special cases and specific handling of a string of text
- Input Text Overrides
 - Applied to the original text string
- Unhandled Text Overrides
 - Applied to the Unhandled Data field

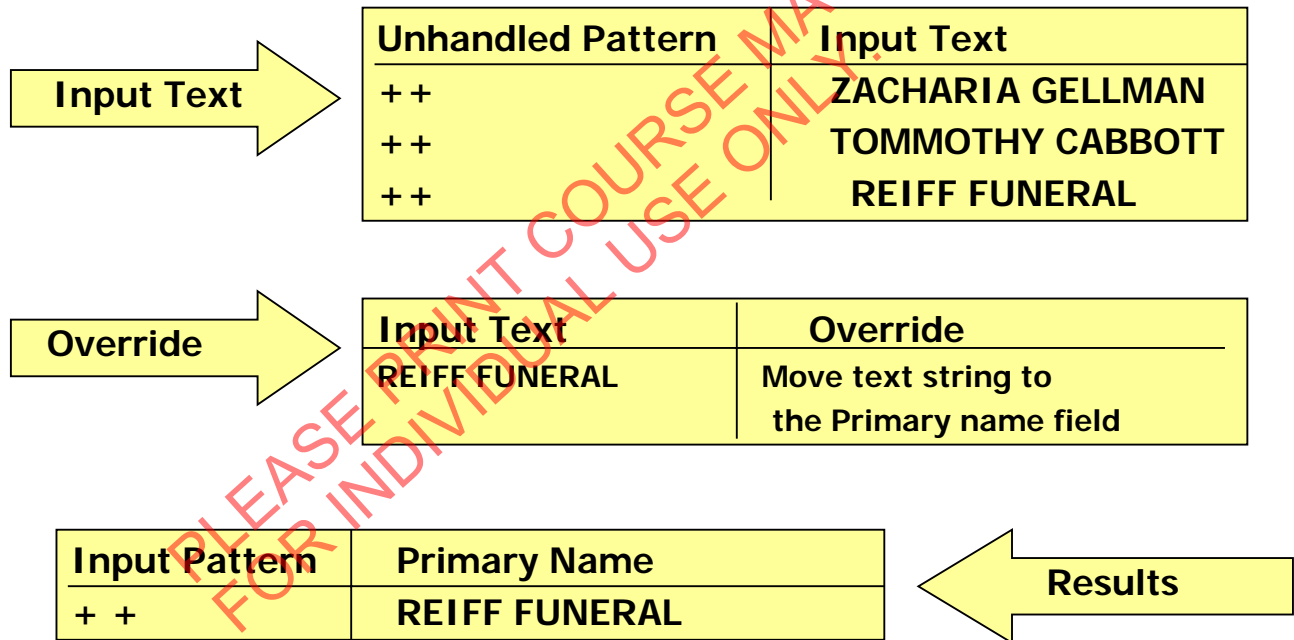
PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY

© 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.

No partial string matching, only complete string matching

Input text overrides



© Copyright IBM Corporation 2007

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

The example REIFF Funeral is a "Special case" as it needs to be handled different than the rest of the data with the Unhandled pattern of ++

The remaining Unhandled Patterns of ++ may be handled the same way. The best type of override to

Pattern overrides

- Allow the user to specify overrides based on an entire pattern
- Use this override when most or all records should be processed with identical logic
- Input Pattern Overrides
 - Applied to the original text string
- Unhandled Pattern Overrides
 - Applied to the Unhandled Data field

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.

Again no partial matching, only complete pattern matching

Pattern Overrides are the most general.

Whenever possible use a pattern override as it is more general and will be applied to many records one override improves the data quality on many records vs. a text override which is very specific to a string of text

Unhandled pattern overrides

Unhandled Pattern	Unhandled Pattern	Input Text
	+, +	HAYWARD, WINSLOW
	+, +	ESHAGHIAN , JOUBI
	+, +	BOULDER, CORONA

Override	Unhandled Pattern	Override
	+, +	Move + to Primary Name Comma provides context Move + to First Name

Unhandled Pattern	First Name	Primary Name
+, +	WINSLOW	HAYWARD
+, +	JOUBI	ESHAGHIAN
+, +	CORONA	BOULDER

Results

© Copyright IBM Corporation 2007

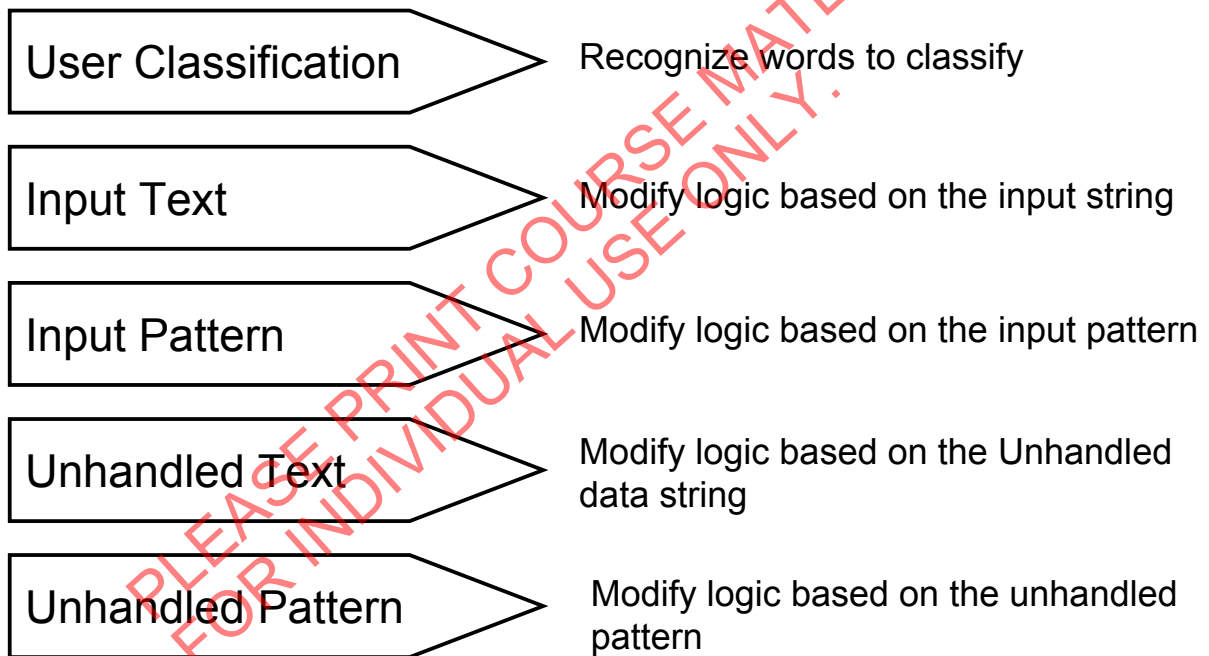
The course materials are provided for internal use only by the individual receiving the materials.
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Why an Unhandled Pattern and not Input Pattern?

Below is an example of a record that has a different input pattern than other records in this category, however it has the same unhandled pattern. All records with this unhandled pattern are to be "handled" (processed) the same way, it is more efficient to use one unhandled pattern override rather than having to apply multiple input pattern overrides

+,+ SANCHEZ-CIFUENTES , RYLMA +-+,+

User override precedence



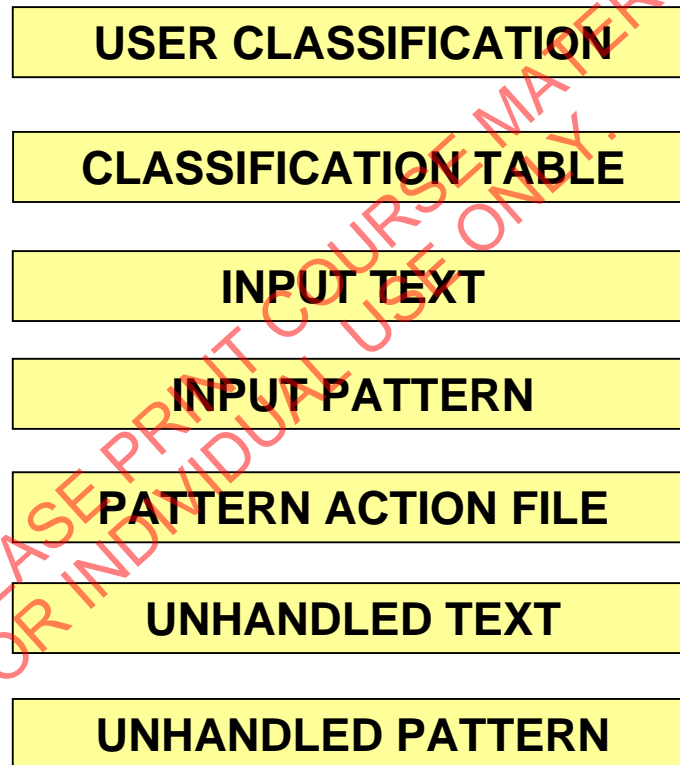
© Copyright IBM Corporation 2007

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.

Text overrides take precedence over pattern overrides because they are more specific

Input overrides take precedence over all other patterns in the pattern-action file

Rule set customization precedence



© 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.

Text overrides take precedence over pattern overrides because they are more specific

Input overrides take precedence over all other patterns in the pattern-action file

Order of what I to look for:

1. Words to classify
2. Input Pattern Overrides
3. Unhandled Pattern Overrides
4. Input Text Overrides
5. Unhandled Text Overrides

Investigate address and area unhandled patterns

- Identify the unhandled patterns for the Address and AREA fields. In the report include the unhandled data, input pattern, original data and the record key.
 1. Build a Character Concatenate Investigation using the following fields
 2. Increase the number of samples to 5

Field Name	Field Description	Type
UPUSADD	Unhandled Pattern	C
UDUSADD	Unhandled Data	X
IPUSADD	Input Pattern	X
ADUSPRE	Address Domain	X

© 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 w

Overrides

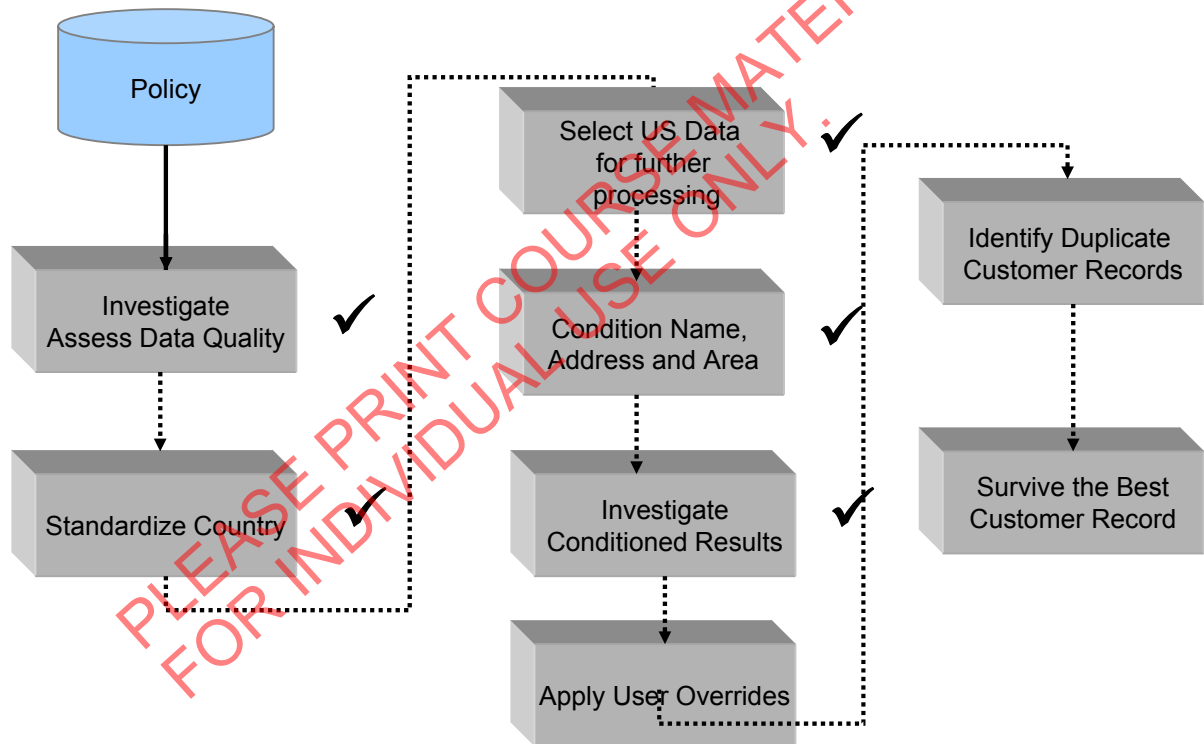
- Purpose
 - Correct problems found during standardization
- Rule set may require overrides because you have data
 - Not processed
 - Incorrectly processed
- Override types
 - Classification
 - Input pattern
 - Input text
 - Unhandled pattern
 - Unhandled text
- Can be tested with rules analyzer

© 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.

Overrides are used to customized rule sets without applying changes to the Pattern Action File.

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.

Overrides screen

Classification - USADDR

Classification | Input Pattern | Input Text | Unhandled Pattern | Unhandled Text

Input Token:

Standard Form:

Classification: A - Abbreviations

Comparison Threshold:

Override Summary

Add Copy Edit Delete

Input Token:	Standard Form:	Classification:	Tolerance:

© 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) WAVES can standardize name fields.
2. (T/F) Rule sets are used in standardization processing.
3. Name the components of rule sets.

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) WAVES can standardize name fields

Answer: False

2. (T/F) Rule sets are used in standardization processing.

Answer: True

3. Name the components of rule sets.

Answer: Classification table, dictionary, pattern action file, lookup tables

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:

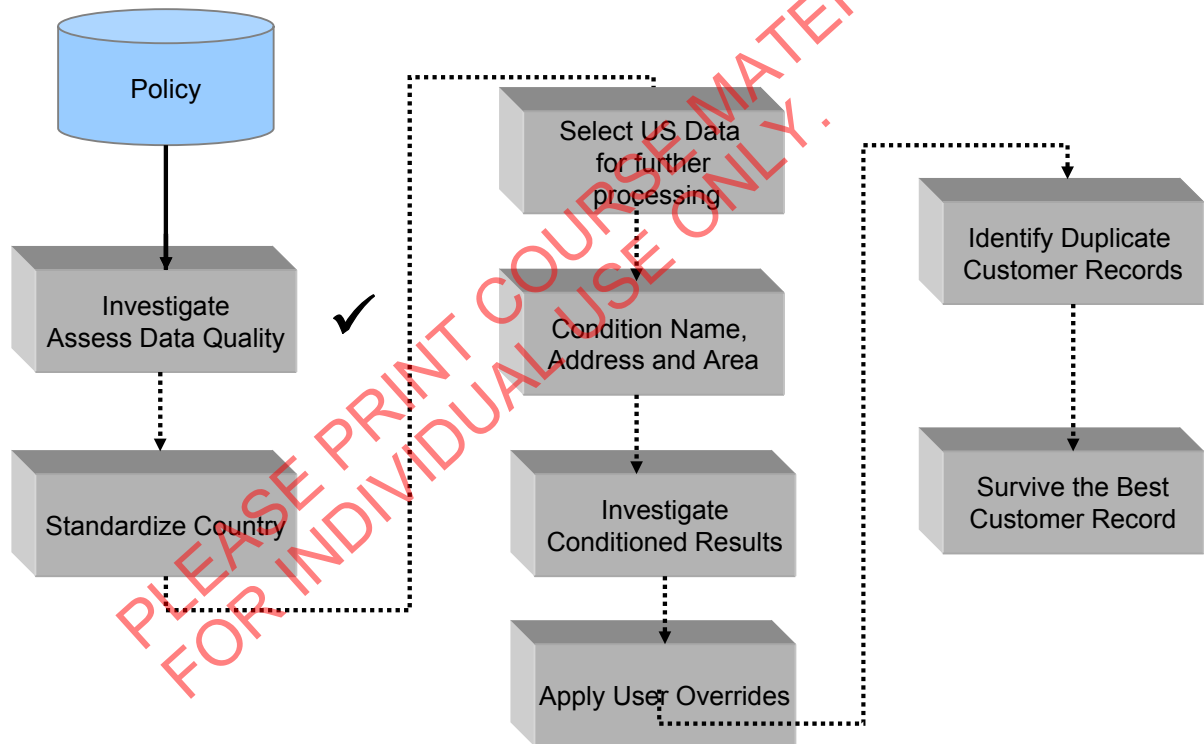
- Describe the Standardize stage in the Data Re-engineering Methodology
- Identify rule sets
- Apply the Standardize stage
- Interpret standardization results
- Investigate unhandled data and patterns

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.

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.

Exercise 10: Standardize country

- Word investigation
 - Uses COUNTRY rule set
- Rule set found in Other folder
- Adds ISO country code to records

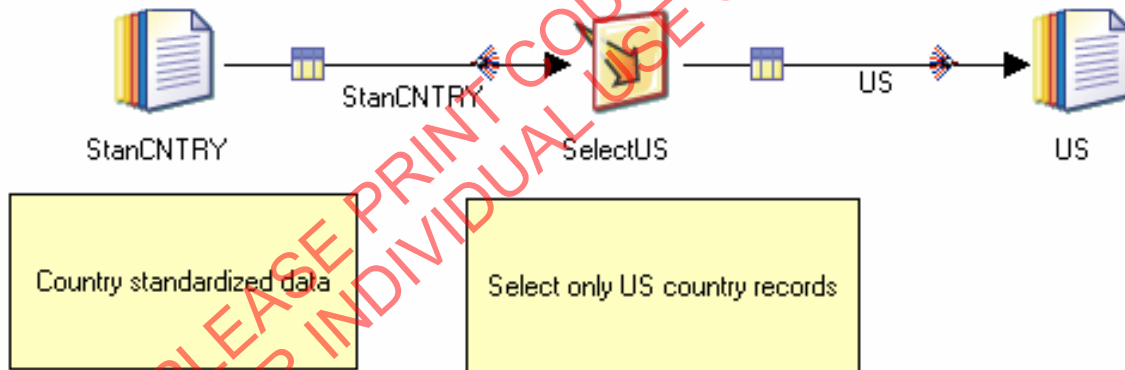


© Copyright IBM Corporation 2007

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 11: Select US records

- Uses Select stage to separate records with US ISO code
- Could also use Transformer stage

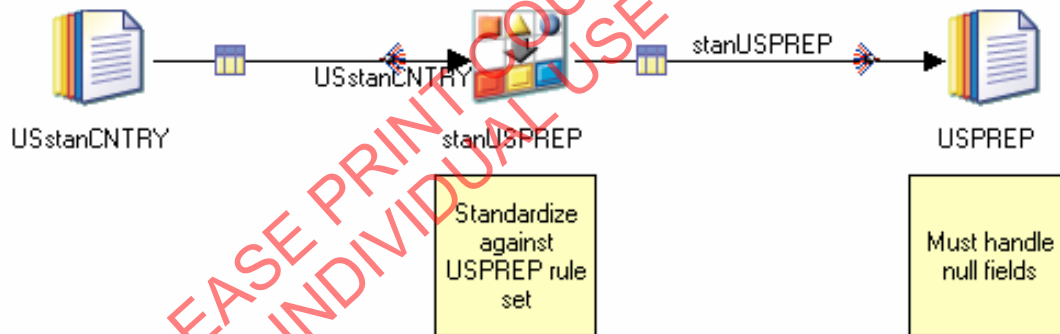


© Copyright IBM Corporation 2007

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 12: Standardize USPSREP

- Word investigation
 - Uses rule sets, must be provisioned
- Rule set found in US folder



© 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 13: Standardize USNAME, USADDR, USAREA

- Word investigation
 - Uses rule sets, must be provisioned
- Rule set found in US folder



© 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 14: Investigate standardization results

- Character concatenate investigation
- C mask used to produce histogram
- X mask used to display other fields of interest

Character Concatenate Investigate

Available Data Columns:

Column Name	Length	Description
AdditionalAddress_USADI	50	
AdditionalName_USNAME	50	
AddressDomain_USDPREF	100	
AddressLine1	35	
AddressLine2	35	
AddressType_USADDR	1	
AreaDomain_USDPREF	100	
AreaType_USADDR	7	

Add To Selected Columns

Scheduled Process

Column Name	Mask
UnhandledPattern_USNAME	CCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
UnhandledData_USNAME	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
InputPattern_USNAME	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
NameDomain_USDPREF	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
FullName	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
RecKey	XXXXX

Delete

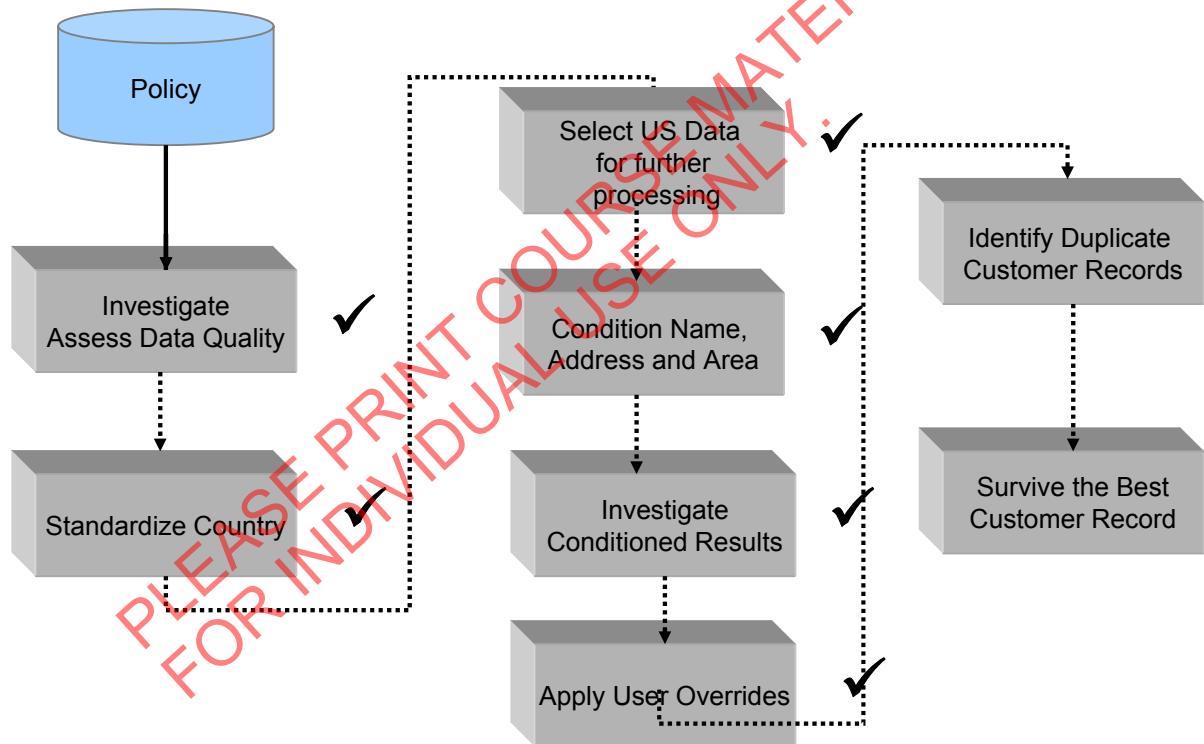
Change Mask

Advanced Options

© 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.

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.

Exercise 15: Apply user overrides

- Classification
- Input pattern
- Unhandled pattern

PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

Classification	Input Pattern	Input Text	Unhandled Pattern	Unhandled Text
Input Token:				
<input type="text"/>				
Standard Form:				
<input type="text"/>				
Classification:				
A - Abbreviations (Misspellings) ▼				
Comparison Threshold: <input type="text"/>				
Override Summary				
<input type="button" value="Add"/> <input type="button" value="Copy"/> <input type="button" value="Edit"/> <input type="button" value="Delete"/>				
Input Token:	Standard Form:	Classification:	Tolerance:	
WINSLOW	WINSLOW	F		
JOUBIN	JOUBIN	F		
MARVYN	MARVYN	F		
HARUKO	HARUKO	F		
CAROLYNNE	CAROLYNNE	F		

© 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.



PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

Match



© 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 QualityStage Match concepts
 - Define the type of matching algorithms
 - Describe the importance of blocking
 - Apply multiple match passes to increase efficiency/efficacy
 - Interpret and improve match results

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

Match stage

- Statistically-based method for determining matches
- Over 24 match comparison algorithms providing a full spectrum of fuzzy matching functions
- Ability to measure informational content of data
- Identify duplicate entities within one or more files
- Match Designer
- Critical field settings

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.

Probabilistic record linkage theory is a subset of comparison of data.

What constitutes a good match?

Which of the following record pairs is a match? And how do you know?

W HOLDEN 12 MAIN ST

W HOLDEN 12 MAINE ST

W HOLDEN 128 MAIN PL 02111 12/8/62

W HOLDEN 128 MAINE PL 02110 12/8/62

WM HOLDEN 128A MAIN SQ 02111 12/8/62 338-0824

WILL HOLDEN 128A MAINE SQ 02110 12/8/62 338-0824

- Do you compare all the shared or common fields?
- Do you give partial credit?
- Are some fields (or some values) more important to you than others? Why?
- Do more fields increase your confidence?
- By how much? What is enough?

© 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.

Before we discuss the technology of matching, think about the human process that you would apply in making a decision about these record pairs.

Do you feel 'comfortable' about the first pair? Is there really enough information to suggest that these records should be linked? The two locations could be anywhere in the world and the first name initial doesn't offer much supporting information either.

What about the second pair? Now we know that we're dealing with the same geographical area and the same birthdate. Has this additional information given you greater confidence? Do you find yourself assigning more-or-less importance to some of the fields or some of the values? For instance, does the abbreviation of PLACE (PL) carry a little more weight in your mind than the abbreviation of STREET (ST) even though they are both just 2 characters? Does the 3 digit building number and the matching PLACE words give you sufficient confidence that these two versions of MAIN are likely the same even though there is a one-digit conflict in the Zip Code? Is the date-of-birth sufficient to say that these are the same person, or is there still some risk of them being twins?

By the time you get to the 3rd pair your confidence should be very high. We now have phone number to further support the location data, and enough first name information to eliminate the risk of twins.

These are the issues that automated matching must consider as well. Being accurate, consistent and justifiable are essential; being able to navigate the "gray-areas" of missing and conflicting values is what separates the simplistic from the industrial strength methods. Now let's look at the methods.

The value of information content

- Information content measures the significance of one field over another (Discriminating Value)
 - A Gender Code contributes less information than a Tax-Id Number
- Information content also measures the significance of one value in a field over another (Frequency)
 - In a First-Name Field, JOHN contributes less information than DWEZEL
- Significance is determined by a value's reliability and its ability to discriminate, both can be calculated from your 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.

So what does INFORMATION CONTENT mean, and why should I care?

It's the phrase that describes the scientific process of measuring the amount of emphasis, meaning, significance, usefulness, or decide-ability that a piece of data contributes to a process -- in this case the process of determining a match.

Its actually a rigorous and mathematically defined concept based on INFORMATION THEORY. And QualityStage is the premier commercial implementation of that theory. QualityStage investigates your actual data, as a step in the matching process, and dynamically adjusts field and value-level scoring based on the characteristics of the data.

You care because it automates, with far greater precision, the human intuitions that cause you to give more or less emphasis to certain values even within the same field.

It results in greater accuracy and because it gives your matching process a legitimacy and justification not possible through other techniques. And that's often essential to enterprise and mission critical projects whose success is measured by the confidence and trustworthiness of the resulting information.

Now lets take a closer look at the Probabilistic process of measuring information content....

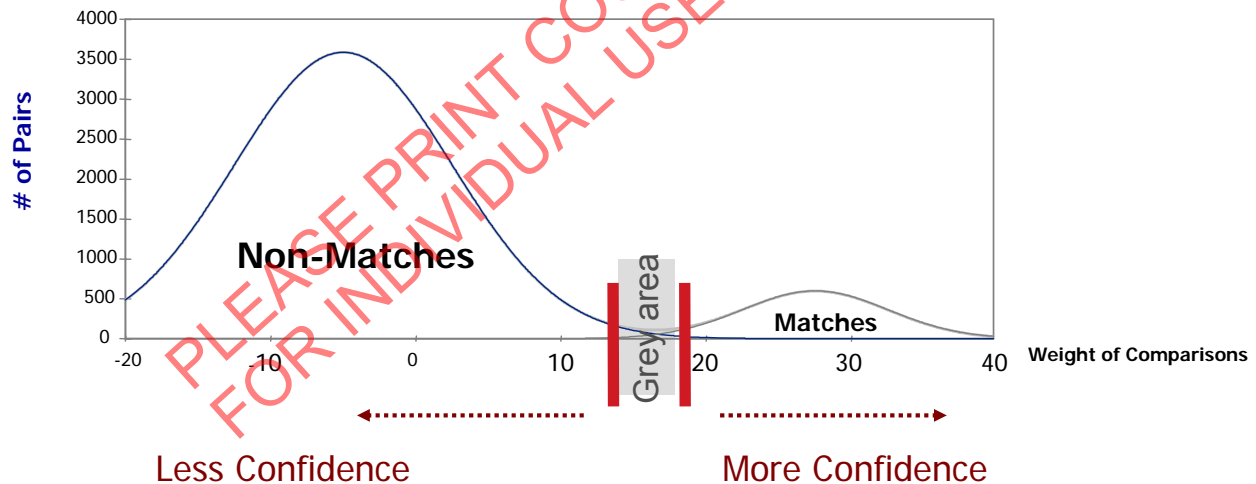
Distribution of weights

WILLIAM J HOLDEN 128 MAIN ST 02111 12/8/62

WILLIAM JOHN HOLDEN 128 MAINE AVE 02110 12/8/62

+1 +1 +17 +2 +4 -1 +7 +9 = 40

The weighted score is a relative measure of the probability of a match



© Copyright IBM Corporation 2007

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

The scores (composites weights are relative to all the other scores). Plot the scores to “see” the distribution of the scores.

▲ This is the distribution of weights for matched and unmatched records. The more variables added to the match, the further apart these two “humps” will be. It’s the point where the two groups intersect which can cause problems.

▲ In our previous example the score of “31.64” based on the distribution this is a fairly high score indicating a high confidence in the match.

▲ This is the distribution of weights for matched and unmatched records. The more variables added to the match, the further apart these two “humps” will be. It’s the point where the two groups intersect which can cause problems.

Including more fields is better as long as each field supports your matching goals

Consider how you would either use or omit fields depending on what your match goals are

Weights

- Measures the information content of a data value
- Each field contributes to the confidence (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.

- ▲ We measure the contribution with a weight.
- ▲ The more contribution the higher the weight, the less the lower the weight.
- ▲ Weight can also be defined as the “discriminating power” of a field

Types of weights

- If a field matches, the agreement weight is used
 - Agreement weight is a positive value
- If a field doesn't match, the disagreement weight is used
 - Disagreement is a negative value
- Partial weight is assigned for non-exact or “fuzzy” matches
- Missing values have a default weight of zero
- Weights for all field comparisons are summed to form a composite weight

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.

▲It is important to combine the theory with “Business Knowledge” to obtain the desired goals.

Example: Statistic program have become so “easy” to use that anyone can build a regression and run it the achieve a result. But it takes an expert, “a knowledgeable business person” to understand the results and the relationship of the input to the result!

Now let's look at that match again and this time we will apply weights.....

Matching terminology

Informational Content	Measures the significance of one field value over another
Weight	Measures the informational content of a data value
Composite Weight	Measures the confidence of a match
Match Cutoffs	Distinguish matches from non-matches
False Positives	Records with a score above the High cutoff that really aren't a match
False Negatives	Records below the low cutoff that really are a match

© Copyright IBM Corporation 2007

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

Measuring the conditions of uncertainty

- Reliability of the data in a given field
 - Estimated as the probability that the field agrees given the record pair is a match
- Probability of a random agreement of values
 - Estimated as the probability the field agrees given the record pair is not 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.

Reliability: How “correct” are the data values. How often are the filled-in (non-missing) and when they are filled-in how often are they correct.

▲**Chance of Random Agreement:** Measures the “rareness” or “uniqueness” of a value. The more frequent (or common) a value occurs in the data the less weight, “confidence” it contributes to the match.

Example: If another Barbara walked into the room would you think AH HAH they must be the same person? ...Well maybe but I’m not convinced. Now, if there were a Vladimir in the room and another Vladimir entered the room, instinctively I would have more confidence that they two Vladimir’s are a match than I would the Barbara’s.

Reliability (m-probability)

- Approximated as, 1 - error rate for the given field
- The higher the m-probability, the higher the disagreement weight will be for the field not matching since the data is considered reliable

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.

- The disagreement weight is proportional to the reliability score
- The m-prob value is entered by the user. It does not need to be an exact measurement as QualityStage will use the user-entered m-prob and improve the measurement based on a sample of the data.
- The more reliable the data in the field the more records are penalized for not agreeing, since errors are relatively rare.
- Estimating the m-prob, if you really don't know then assume a 10% error rate, (m-prob =90% or .9)
- The m-prob is between .001 and .999. It can never be 1 (100%), there is always a chance of error. And it can never be 0 (completely reliable).
- Data is reliable = errors are rare
- Data is not reliable = errors are common
- Example: of M-PROB: If the variable street type has a 12% error rate, then the m-probability is 0.88

Chance agreement (u-probability)

- The u-probability can be approximated as the probability that a field agrees at random (by chance)
- QualityStage uses a frequency analysis to determine the probability of chance agreement for all values
- Rare values bring more weight to 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.

▲ Rare values have less chance of accidental agreement and contribute more to a match

▲ Frequency analysis determines the probability of chance agreement for any values (INTEGRITY calculates)

Example: If two records in a matched pair have a name of John Smith specified, you would be less sure that the record pair represented a true match than if Vladimir Horowitz were matched on both records

Rare events have more discriminating power than common events

Frequencies should not be calculated for fields such as individual identification numbers since all values are rare

Calculating weights

- Agreement weight is estimated as:
 - $\log_2(m/u)$
- Disagreement weight is estimated as:
 - $\log_2((1-m)/(1-u))$

$$M \text{ (m-prob)} = .9$$

$$U \text{ (u-prob)} = .01$$

$$\text{Agreement weight} \quad \log_2 (.9/.01) \quad = \quad 6.49$$

$$\text{Disagreement weight} \quad \log_2 (1-.9)/(1-.01) \quad = \quad -3.31$$

© 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.

Note: Use a sample size of 10,000 with a value frequency of 100 for a u-prob of .01 (1 in a hundred)

Blocking

- Grouping together like records that have a high-probability of producing matches
- Only “like” records are compared to each other making the match more efficient and computationally feasible
- Records in a “block” match exactly on one to several blocking fields

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.

Within each group of “blocked” records, each record is compared to every other record according to the matching variables.

Blocking example: sample data

- Block on NYSIIS of Last Name

NYSIIS LNAME	NAME	ADDRESS	ZIP
YANG	YUNG , WAYNE D	9000 SHEPARD DRIVE	78753
GARAS	GEROSA, FRAN X	29 AARONS CT	06877
YANG	YOUNG , JONATHAN A	1767 TOBEY ROAD	30341
GARAS	GERISA, FRANCIS	29 AARONS CT	06877
GARAS	GEROSA, FRANCIS XAVIER	29 AARONS COURT	06877
MATAC	MARCUS MATIC	100 SUMMER STREET	02111
GARAS	GEROSA, MARY	29 AARONS CT	06877
JANCAN	RENEE JENKINS	100 SUMMER STREET	02111
YANG	YOUNG THERESA C	1767 TOBEY ROAD	30341

© Copyright IBM Corporation 2007

The course materials are provided for internal use only by the individual rece
The materials may not be modified, copied, distributed or transferred without the expres

1. We could compare every record to every other record or we could block them by: Last Name.
2. Enter to show block color-coding
3. Notice the Jerosa record did not make the same block group as the "Gerosa" records. It did not match exactly on Last Name. This is one reason we create the Phonetic coding of some fields in the Conditioning phase. The phonetic coding fields are very useful for blocking as they introduce "fuzziness" to a rigid set of criteria (blocking).

Blocking example – NYSIIS of Last Name

NYSIIS	NAME	ADDRESS	ZIP
YANG	YOUNG , WAYNE D	9000 SHEPARD DRIVE	78753
YANG	YOUNG , JONATHAN A	4220 BELLE PARK DR.	77072
YANG	YOUNG THERESA C	1767 TOBEY ROAD	30341

GARAS	GEROSA, FRAN X	29 AARONS CT	06877
GARAS	GEROSA, FRANCIS XAVIER	29 AARONS COURT	06877
GARAS	GEROSA, MARY	29 AARONS CT	06877
GARAS	GARISA, FRANCIS	29 AARONS CT	06877

MATAC	MARCUS MATIC	100 SUMMER STREET	02111
-------	--------------	-------------------	-------

JANCAN	RENEE JENKINS	100 SUMMER STREET	02111
--------	---------------	-------------------	-------

Blocks with only one record are considered residuals

© 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.

Blocks with one record are considered “residuals”. There are not any other records in the group to compare to.

“Due to an error (potentially) in the last name Jerosa it did not make the same block group as the Gerosa records.”

Balance scope and accuracy

Balance the scope and accuracy to compare a reasonable amount of “like” records

Accuracy

The quality of the candidate records

Scope

The number of records



© 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.

▲ If the accuracy is loose the scope is very large

▲ If the accuracy is too tight then the scope is too small

Example: If you are matching bank records to customers, it is NEVER OK to match the wrong record to a customer. The tolerance for error is very low. The accuracy must be high which causes a narrower scope of records.

Marketing Campaign: In order to market a reasonable number of customers you might be willing to tolerate more error (less accuracy) to get a sufficiently wide scope, or “marketing list”.

Blocking strategy

- Choose fields with reliable data
- Choose fields with a good distribution of values
- Combinations of fields may be used

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.

▲The goal of blocking is to group together like records that have a high probability of producing matches. The character discrete INV reports will help with these decisions, they tell you how often a field is populated.

▲If you choose fields with reliable data then you are “truly” grouping together like records since the data values are reliable (usually correct).

▲Choose fields that make business sense to meet your objective. If you are trying to identify unique customers then blocking by house number isn't the best choice.

▲Gender usually doesn't have enough values to break records into groups of 100-200 (guideline). If all the data is from a few states the state may not be the best field.

▲Again Inv reports tell you how often a field is populated, the distribution of the data.

Examples of blocking strategies

- Zip code for matching addresses
- NYSIIS of last name for matching individuals
- Brand name for matching products
- Combination of zip code and NYSIIS of street name for matching addresses
- Combination of NYSIIS of last name and first letter of first name for matching individuals

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.

Blocking summary

- Blocking groups together “like” records
- Matching is more efficient for small block sizes
 - Blocks should have less than 1000 records
- Blocking fields must match exactly for a candidate set to be created/evaluated

PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

© Copyright IBM Corporation 2007

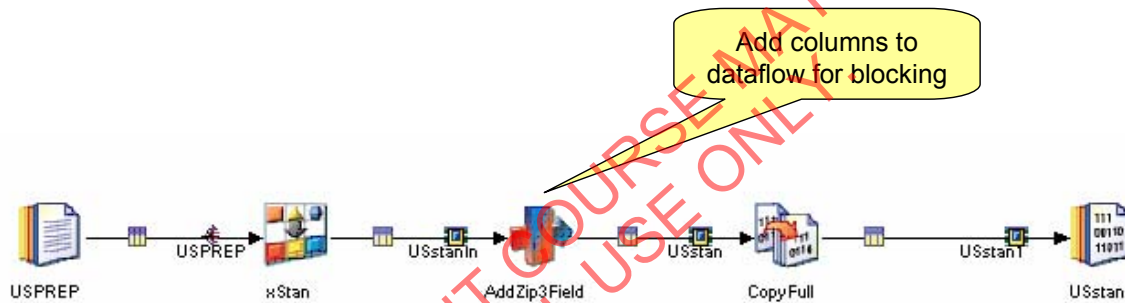
The course materials are provided for internal use only by the individual receiving the materials.

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

Example:

100*100 records (10,000 comparisons) is much faster than 200 (200*200 or 40,000 records)

DataStage job with Transformer



Column definition:

```
If IsNull(USStanIn.ZipCode_USAREA) Then
SetNull() Else USStanIn.ZipCode_USAREA[1,3]
```

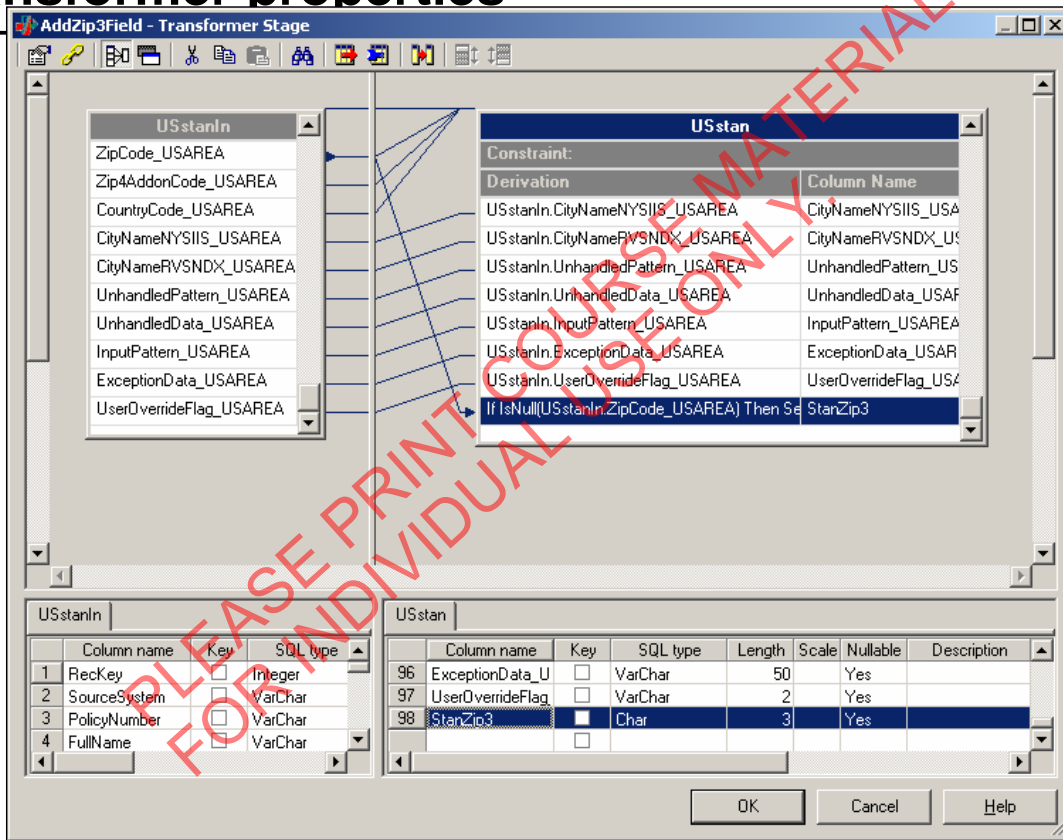
© Copyright IBM Corporation 2007

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.

Transformer properties

IN

OUT



© 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 16: Add fields using Transformer stage

- Create new field to be used in blocking



© 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 types

- Unduplication
 - Identifies duplicates candidates in one file
- Reference Match (Two File)
 - One-to-one correspondence
 - For every record on stream link we expect to find a match to one record on reference link
 - Many-to-one correspondence
 - More than one record on stream link can match to the same record on reference link

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.

Examples:

One-to-one means Customer records from the billing system should have one-to-one correspondence with customers from the marketing database.

One-to-Many means Many Visa transactions will match to the same credit card number. Many addresses match to one postal code.

Comparing data values

- Different comparisons for different data
- Over 24 comparison methods
- Most common
 - CHAR - (character comparison) character by character, left to right.
 - UNCERT - (character uncertainty) tolerates phonetic errors, transpositions, random insertion, deletion, and replacement of characters
 - CNT_DIFF – Counts keying errors in numeric fields. You set a tolerance threshold
 - NAME_UNCERT – Can be used to compare and character values, if the strings are different lengths then the shorter of the two lengths is used

© 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.

Over 24 ways to compare data values

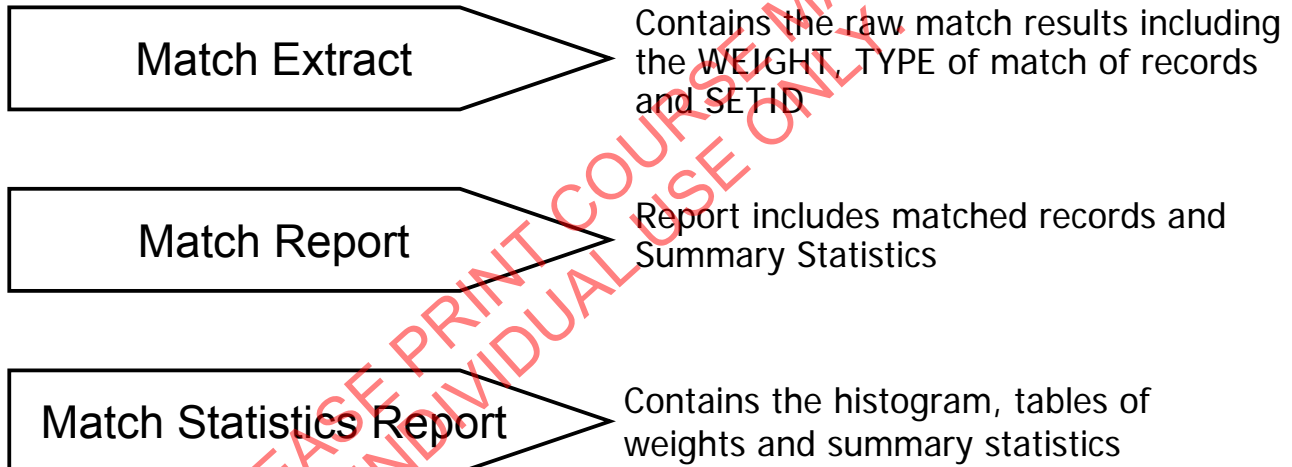
Char = Exact = Total agreement weight

Uncert = Fuzzy = Agreement weight is prorated based on how close to exact.

If records are not “close enough” then the disagreement weight is assigned.

These two are the most popular ways to compare data in fields.

Match output files



Must examine both sequential files and 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.

Review Match extract for match results

Review the Match Report with clients and Business Analysts

Match Debug file, check for block overflow and review the Histogram

Match extract

SETID	TYPE	PASS	WEIGHT	ALL_OF_THE_DATA
393	XA	1	55.32	MICHAEL F DOHERTY
393	DA	1	41.36	MICHAEL F DOUGHERTY
468	XA	1	50.40	EUGENE B BOROWITZ
468	DA	1	24.01	BOROWITZ FAMILY TRUST
468	DA	1	47.26	GENE BOROWITZ
520	XA	1	52.75	FRAN X GEROSA
520	DA	1	40.95	FRANCIS XAVIER GEROSA
520	DA	1	52.75	FRANCIS X GEROSA
520	DA	1	41.22	FRANK X GEROSA
1035	RA	1		DARRYL F LINDBERG

© 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 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.

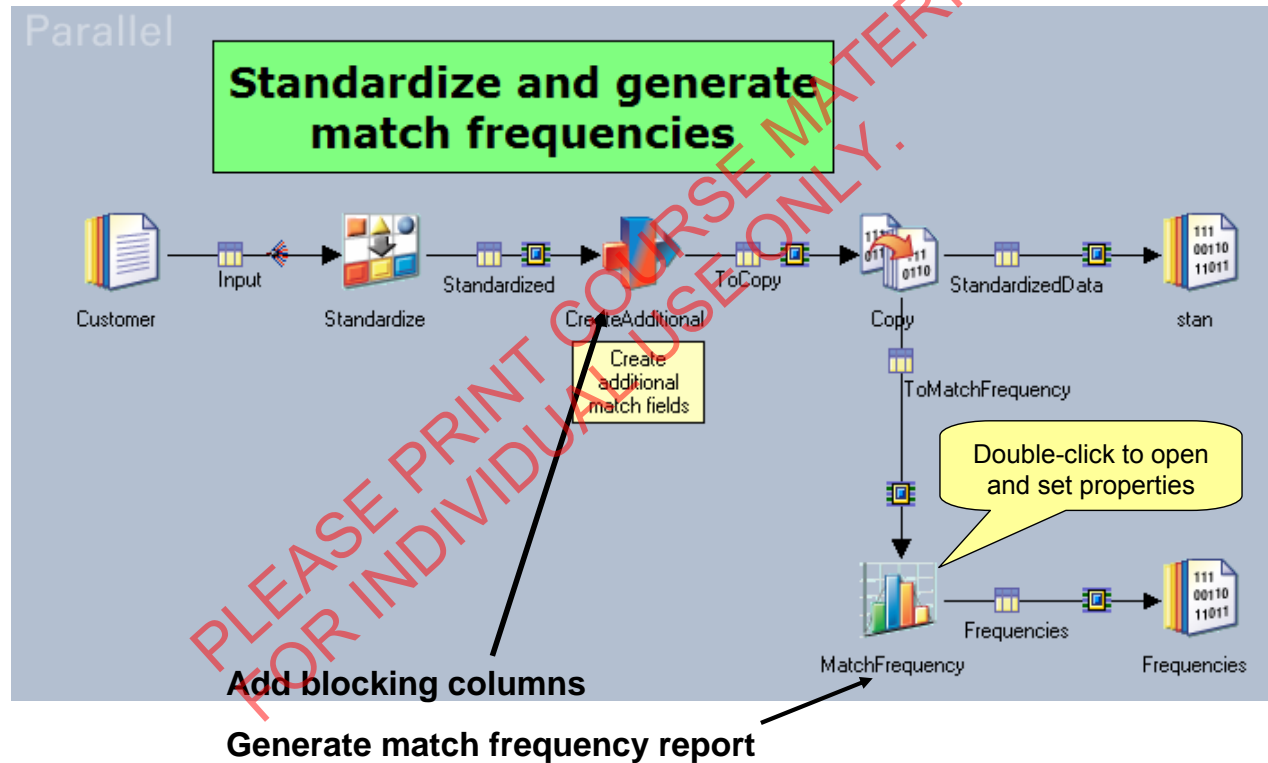
Tasks required in match process

- Standardize the data
- Add data columns needed for blocking
- Generate match frequency report
- Build match specification in Match Designer
 - Add pass
 - Blocking columns
 - Match commands
 - Configure match test results environment
- Run pass
- Review results
- Tune the match
 - Add cutoffs
 - Set overrides
 - Add more passes
- Repeat steps until match results are acceptable

© 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.

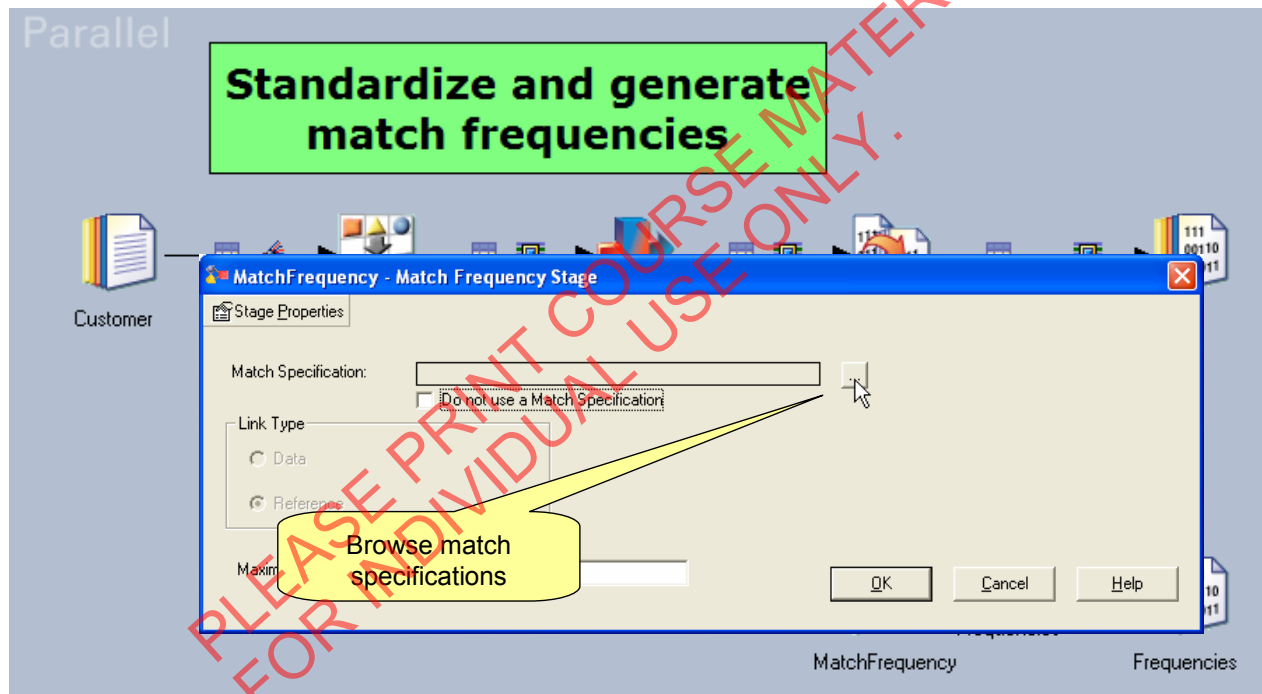
Add columns and generate match frequency



© 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 data is used by the match process.

Match frequency generation

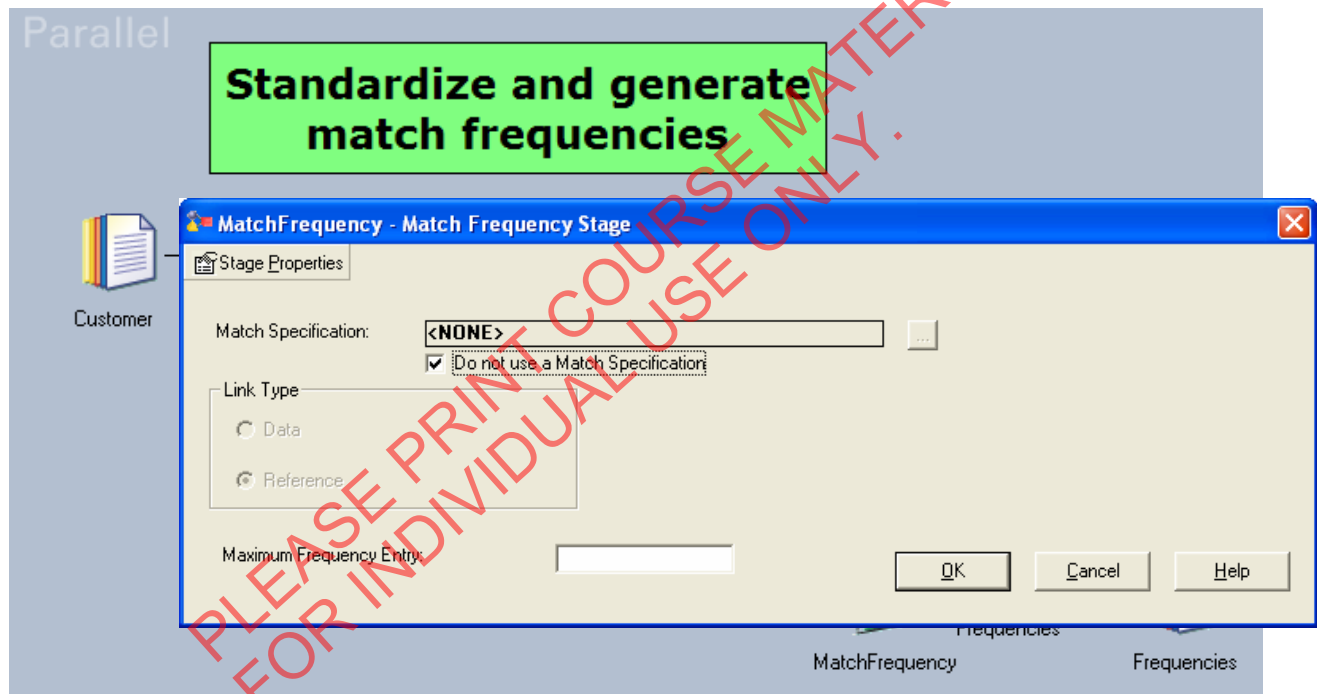


Option 1: Use an existing Match Specification

© 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 frequency generation



Option 2: Build the frequencies and develop the match

© 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 frequency generation

Parallel

Customer

MatchFrequency - Match Frequency

Stage | Input | Output

Output name: Frequencies

Columns... View Data...

General Mapping Columns Advanced

Columns		Frequencies	
Expression	Column Name	Derivation	Column Name
MatchFrequency(qsFreqValue)	qsFreqValue	MatchFrequency(qsFreqValue)	qsFreqValue
MatchFrequency(qsFreqCounts)	qsFreqCounts	MatchFrequency(qsFreqCounts)	qsFreqCounts
MatchFrequency(qsFreqColumnID)	qsFreqColumnID	MatchFrequency(qsFreqColumnID)	qsFreqColumnID
MatchFrequency(qsFreqHeaderFlag)	qsFreqHeaderFlag	MatchFrequency(qsFreqHeaderFlag)	qsFreqHeaderFlag

Map output columns

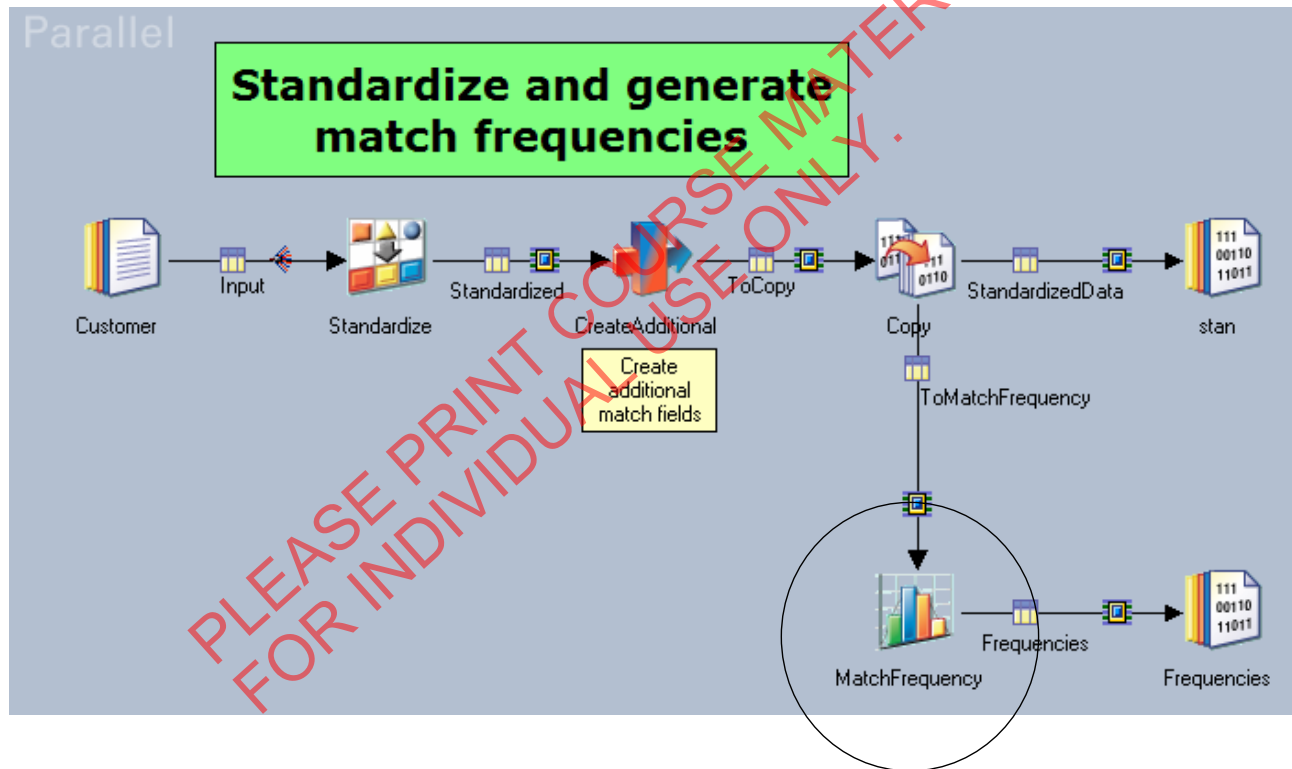
Find Auto-Match

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.

Match frequency generation



© 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 frequency generation

qsFreqValue	qsFreqCounts	qsFreqColumnID	qsFreqHeaderFlag
ADAMSON	00000003 00000000 00000000	97	1
ARMSTRONG	00000002 00000000 00000000	97	1
AVINGER	00000003 00000000 00000000	97	1
BAKER	00000002 00000000 00000000	97	1
BARNETTE	00000003 00000000 00000000	97	1
BARR	00000002 00000000 00000000	97	1
BELL	00000003 00000000 00000000	97	1
BERNEY	00000009 00000000 00000000	97	1
BERNEY FNCL SYSTEMS	00000003 00000000 00000000	97	1
BRUNSON	00000003 00000000 00000000	97	1
BULLINGTON	00000003 00000000 00000000	97	1
CANNON	00000003 00000000 00000000	97	1
CARLISLE	00000002 00000000 00000000	97	1
CARNELL WILKES	00000003 00000000 00000000	97	1
CARROLL	00000004 00000000 00000000	97	1
CHANCE	00000002 00000000 00000000	97	1
CHENEY	00000004 00000000 00000000	97	1
CHRISTENBERG	00000006 00000000 00000000	97	1
COGBORN	00000005 00000000 00000000	97	1



stan



Frequencies

MatchFrequency

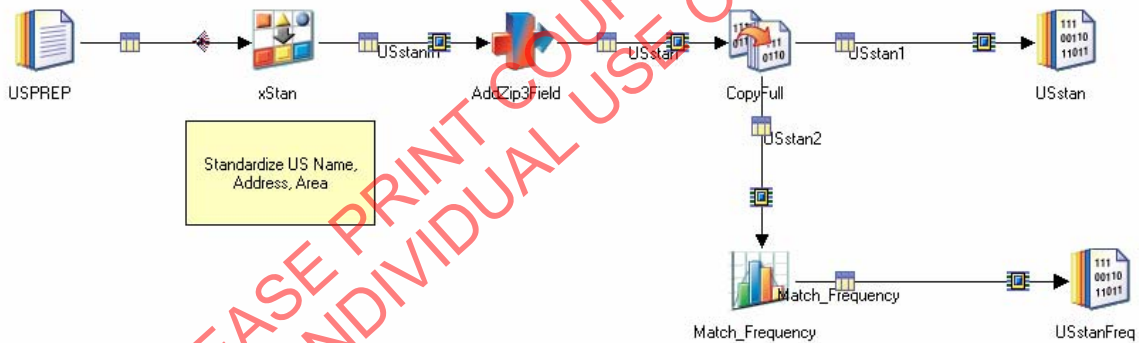
© Copyright IBM Corporation 2007

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 17: Match frequency

- Use Match Frequency stage in a match job



© 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 Designer

- Used to build a match specification that will be addressed in a match job

Features

- Design control center
- A data-centric, graphical, self-contained environment
- Graphical representation of statistics
- Match design is independent of job design
- Reusable components
- Separate design from the physical data representation
- Iterative development

© 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 Design - Unduplicate

PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

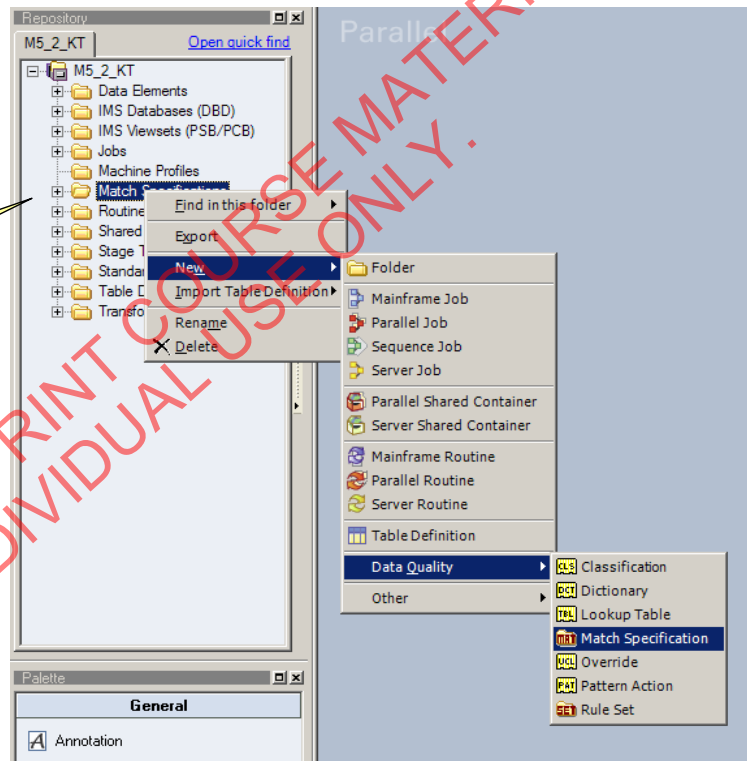
© Copyright IBM Corporation 2007

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.

Match design - unduplicate

How to create a new match specification

Right-click in non-root area of repository

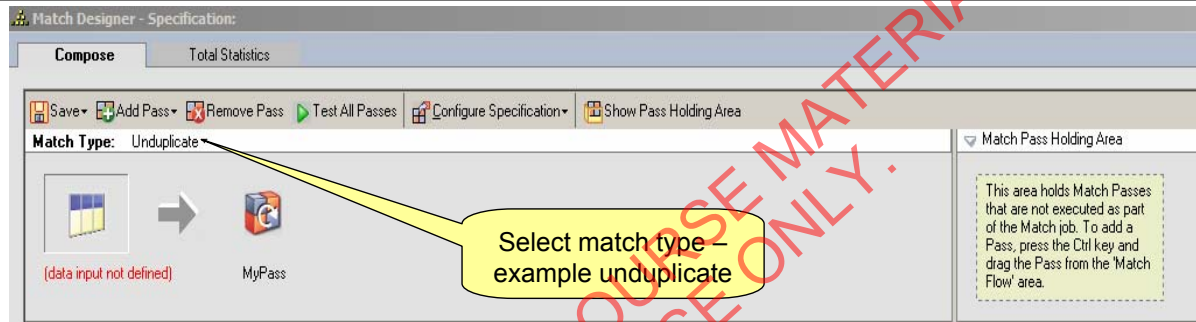


© Copyright IBM Corporation 2007

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.

To start the Match Designer, right-click in the repository view.

Match design - unduplicate



Will initially get one pass called
MyPass

© 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 design - unduplicate

Click table definition icon

Use load button to access table definition of standardized data set

Data Table Definition
Name: StandardizedData

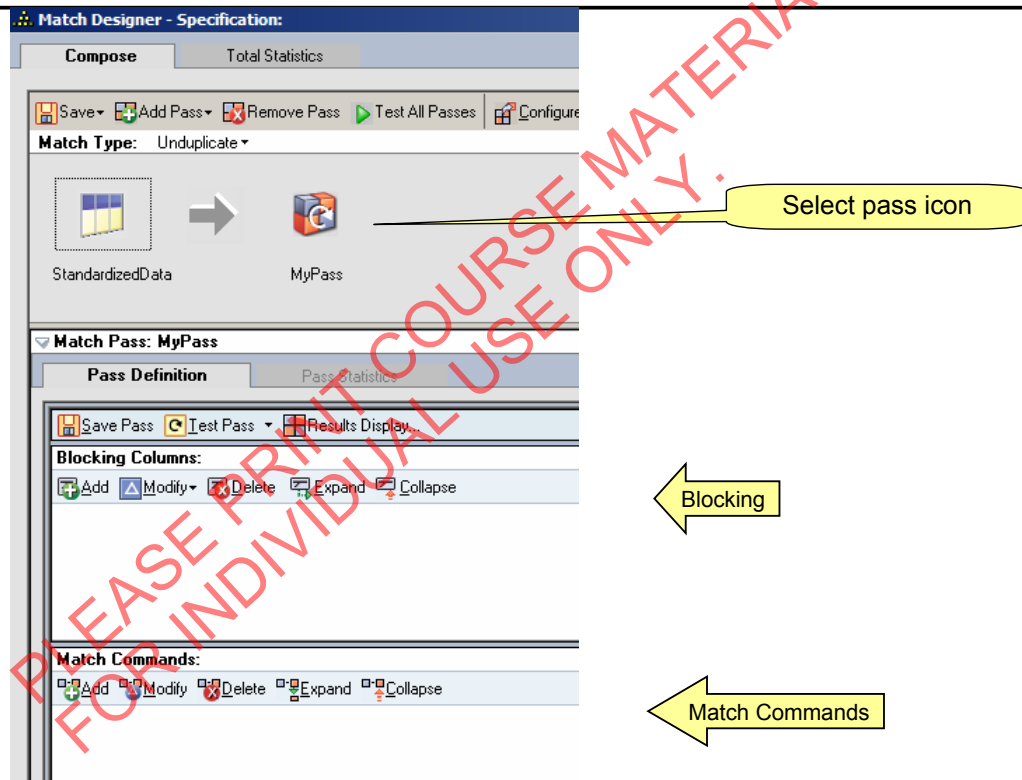
Column Name	Sql Type	Length	Description
MatchFirst1	VarChar	1	Match First Name First Character
HouseNumberFirstChar	VarChar	1	First Character of House Number
ZipCode3	VarChar	3	First Three Characters of Zip Code
UniqueIdentifier	VarChar	10	Unique Identifier
ApplicantSSN	VarChar	9	Applicant SSN
Name	VarChar	70	Name
AddressLine1	VarChar	50	Address Line 1
AddressLine2	VarChar	50	Address Line 2
City	VarChar	40	City
State	VarChar	2	State
Zip5	VarChar	5	Zip Code
Zip4	VarChar	4	Zip + 4
NameType	VarChar	1	
GenderCode	VarChar	1	
NamePrefix	VarChar	20	
FirstName	VarChar	25	
MiddleName	VarChar	25	
PrimaryName	VarChar	50	
NameGeneration	VarChar	10	
NameSuffix	VarChar	20	
AdditionalName	VarChar	50	

Load Edit... 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.

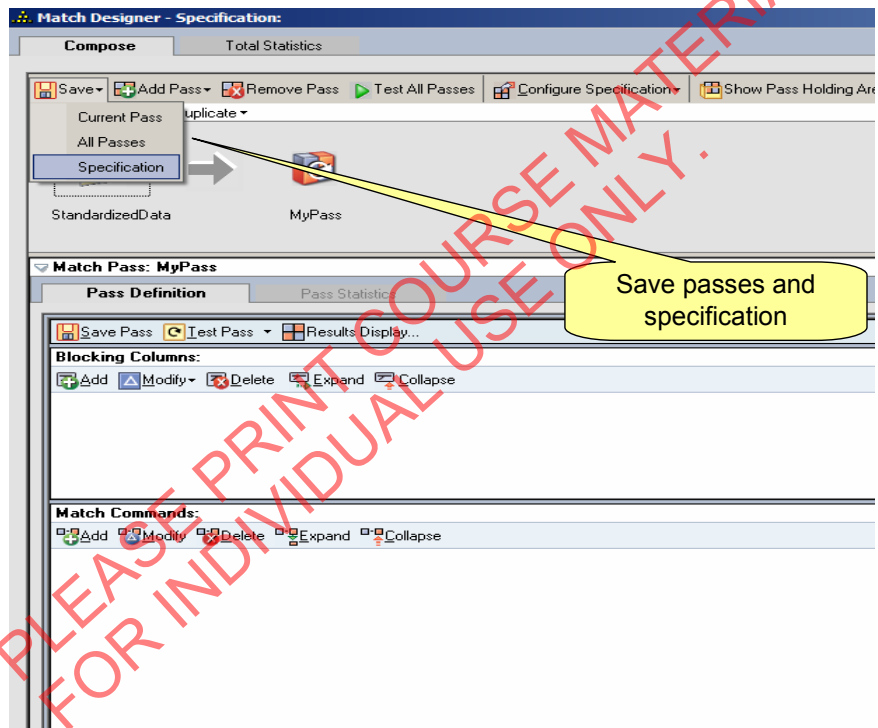
Match design - unduplicate



© 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.

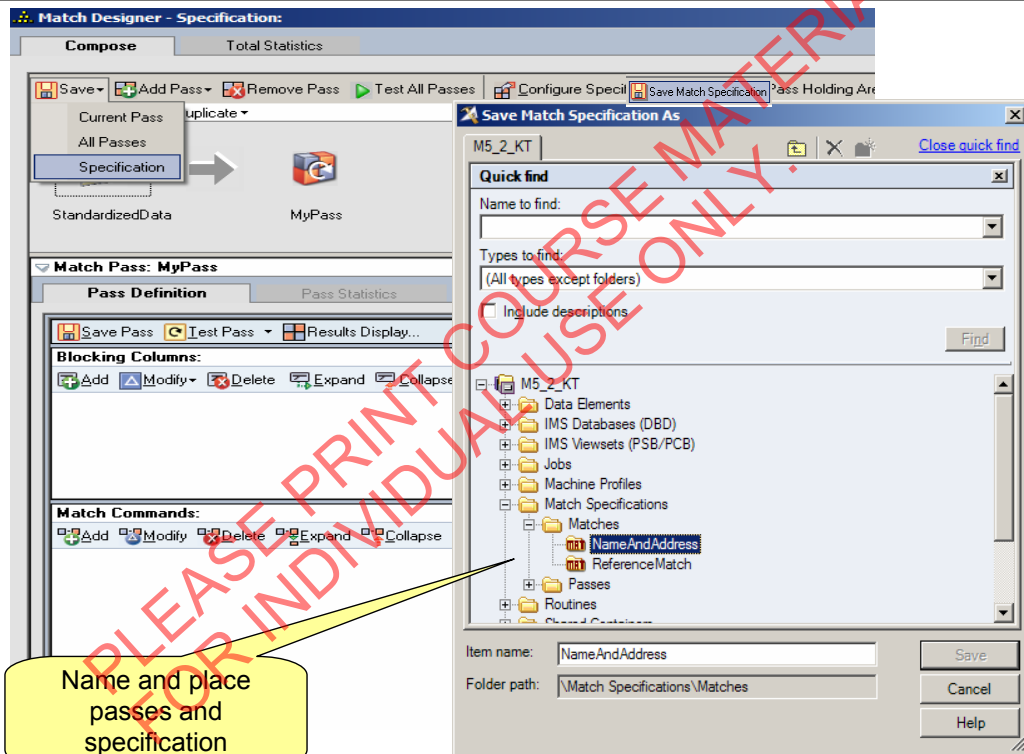
Match design - unduplicate



© 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.

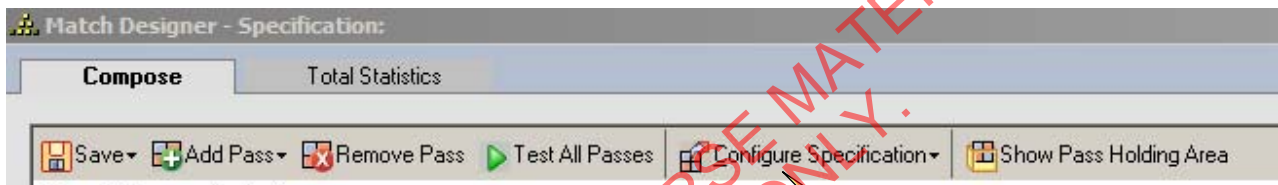
Match design - unduplicate



© 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.

Match design - unduplicate



Questions:

Where is the standardized data?

Where is the frequency report?

What ODBC-accessed database will store test results?

Set up test results 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.

Match design - unduplicate

Note: these must be data sets

Standardized sample data

Frequencies data set

Data Source Name
User Name
Password

Test Environment

To maximize update efficiency uncheck items whose name, format, and content have not changed...

Sample Information

☒ Data Sample Data Set: [] ...

☐ Reference Sample Data Set: [] ...

Frequency Information

☒ Data Frequency Data Set: [] ...

☐ Reference Frequency Data Set: [] ...

☒ Maximum Frequency Value: [100]

Test Results Database

☒ ODBC Data Source Name: [(Please select a DSN)]

Username: []

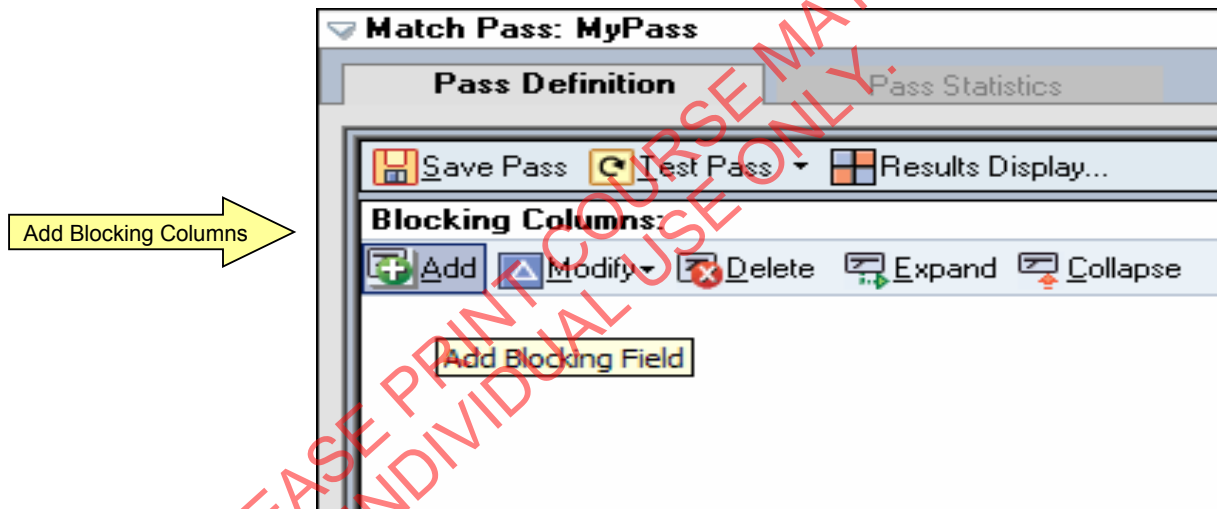
Password: [] [Test Connection]

[Update] [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.

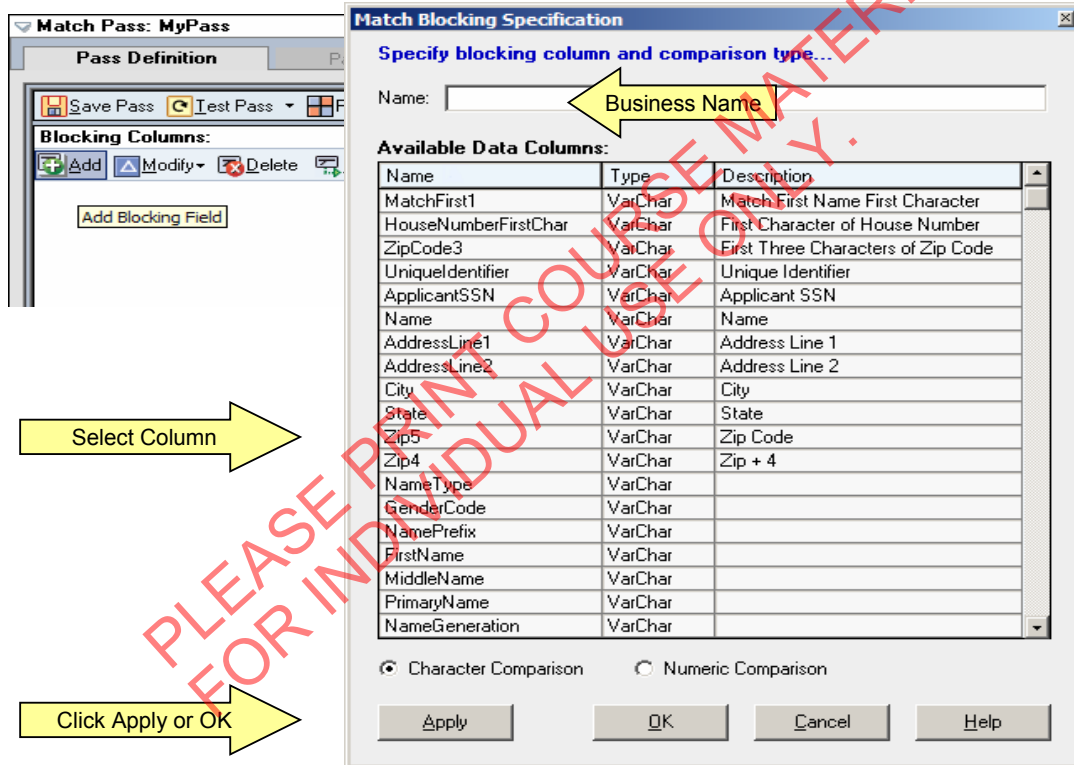
Match design - unduplicate



© Copyright IBM Corporation 2007

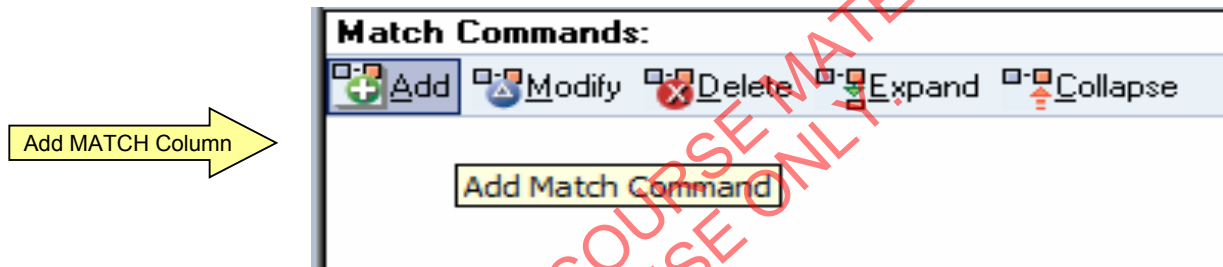
The course materials are provided for internal use only by the individual receiving the materials.
The materials may not be modified, copied, distributed or transferred without the express prior written consent of IBM.

Match design - 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.

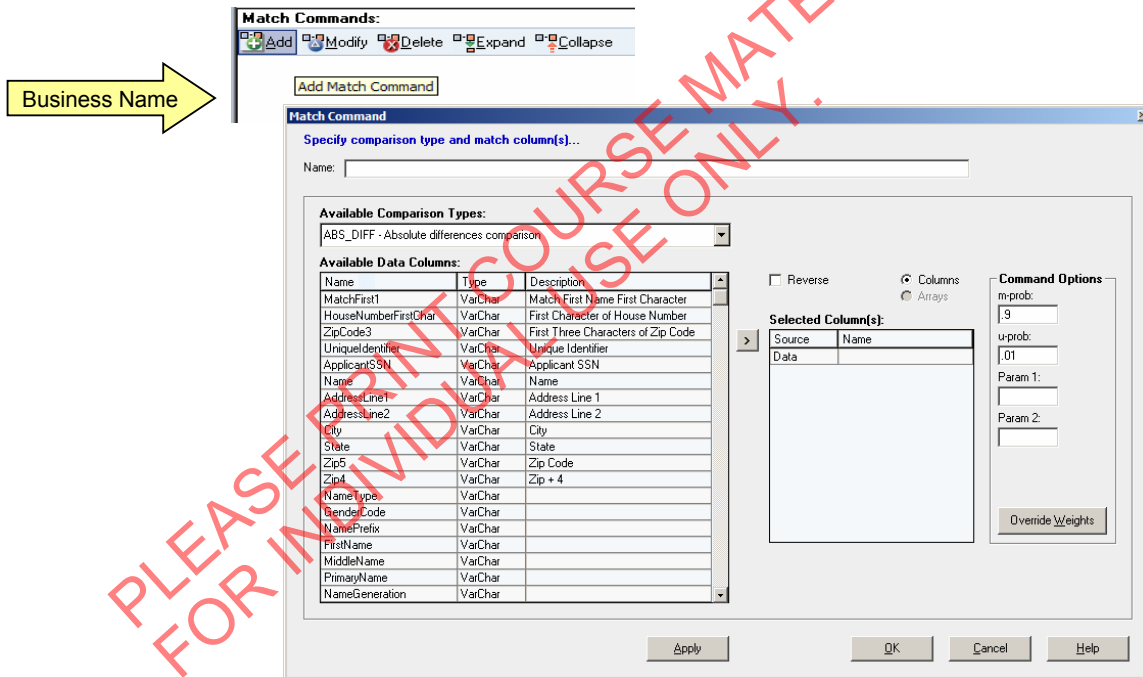
Match design - unduplicate



© 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.

Match design - 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.

Match design - unduplicate

Match Commands:

Match Command

Specify comparison type and match column(s)...

Name:

Available Comparison Types:

Comparison	Description
ABS_DIFF - Absolute differences comparison	Absolute differences comparison
CHAR	Character comparisons
CNT_DIFF	Counting errors in fields
DATE8	Date comparisons in the YYYYMMDD form
DELTA_PERCENT	Delta percentage comparisons
DISTANCE	Geometric distance comparisons
INT_TO_INT	Interval interval comparisons
MULT_EXACT	Multi-word exact
MULT_UNCERT	Multi-word uncertainty
NAME_UNCERT	First name uncertainty string comparison
NUMERIC	Numeric comparisons
PREFIX	Prefix comparisons for truncated fields
PRORATED	Prorated comparisons
TIME	Time comparisons
UNCERT	Character uncertainty comparisons
USPS_INT	United States Postal Service intervals

☐ Reverse ☒ Columns ☐ Arrays

Selected Column(s):

Source	Name
Data	

Command Options

m-prob:

u-prob:

Param 1:

Param 2:

© 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 design - unduplicate

Match Commands:

☐ Add
 ☐ Modify
 ☐ Delete
 ☐ Expand
 ☐ Collapse

Add Match Command

Match Command

Specify comparison type and match column(s)...

Name: Last Name

Available Comparison Types:

UNCERT - Character uncertainty comparisons

Available Data Columns:

Name	Type	Description
NamePrefix	VarChar	
FirstName	VarChar	
MiddleName	VarChar	
PrimaryName	VarChar	
NameGeneration	VarChar	
NameSuffix	VarChar	
AdditionalName	VarChar	
MatchFirstName	VarChar	
MatchFirstNameNYSIIS	VarChar	
MatchFirstNameRVSNDX	VarChar	
MatchPrimaryName	VarChar	
MatchPrimaryNameHashK	VarChar	
MatchPrimaryNamePackK	VarChar	
NumofMatchPrimaryWords	VarChar	
MatchPrimaryWord1	VarChar	
MatchPrimaryWord2	VarChar	
MatchPrimaryWord3	VarChar	
MatchPrimaryWord4	VarChar	
MatchPrimaryWord5	VarChar	

Properties

Selected Column(s):

Source	Name
Data	

☐ Reverse

☒ Columns

☐ Arrays

Command Options

m-prob: .9

u-prob: .01

Param 1:

Param 2:

Override Weights

Apply OK Cancel Help

Data Column

Right-Click to view data frequencies

© 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 design - unduplicate

Match Commands

☐ Add
 ☐ Modify
 ☐ Delete
 ☐ Expand
 ☐ Collapse

Descriptions

Add Match Command

Match Command

Specify comparison type and match column(s)...

Name: Last Name

Available Comparison Types:

UNCERT - Character uncertainty comparisons

Available Data Columns:

Name	Type	Description
NamePrefix	VarChar	
FirstName	VarChar	
MiddleName	VarChar	
PrimaryName	VarChar	
NameGeneration	VarChar	
NameSuffix	VarChar	
AdditionalName	VarChar	
MatchFirstName	VarChar	
MatchFirstNameNYSIIS	VarChar	
MatchFirstNameRVSNDX	VarChar	
MatchPrimaryName	VarChar	
MatchPrimaryNameHashKi	VarChar	
MatchPrimaryNamePackKi	VarChar	
NumOfMatchPrimaryWords	VarChar	
MatchPrimaryWord1	VarChar	
MatchPrimaryWord2	VarChar	
MatchPrimaryWord3	VarChar	
MatchPrimaryWord4	VarChar	
MatchPrimaryWord5	VarChar	

Properties

Apply

Column Details

Find:

String: inc:

Name: MatchPrimaryName

Description:

Type: VarChar

Frequency:

(Click on column name to sort or reverse order...)

VALUE	QUANTITY	PERCENTAGE
JONES	18	58
WINDHAM	10	32
BERNEY	9	29
SIPPIAL ENTERPRISES	8	26
HARRIS	7	23
HUGHEY	7	23
JOHNSON	7	23
CHRISTENBERRY	6	19

Cancel Help

Frequencies

© 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 design - unduplicate

Match Commands

☐ Add
 ☐ Modify
 ☐ Delete
 ☐ Expand
 ☐ Collapse

Descriptions
 Add Match C

Match Command

Specify comparison type and match column(s)...

Name:

Available Comparison Types:

UNCERT - Character uncertainty comparisons

Available Data Columns:

Name	Type	Description
NamePrefix	VarChar	
FirstName	VarChar	
MiddleName	VarChar	
PrimaryName	VarChar	
NameGeneration	VarChar	
NameSuffix	VarChar	
AdditionalName	VarChar	
MatchFirstName	VarChar	
MatchFirstNameNYSIIS	VarChar	
MatchFirstNamePVSNDX	VarChar	
MatchPrimaryName	VarChar	
MatchPrimaryNameHashKi	VarChar	
MatchPrimaryNamePackKi	VarChar	
NumofMatchPrimaryWords	VarChar	
MatchPrimaryWord1	VarChar	
MatchPrimaryWord2	VarChar	
MatchPrimaryWord3	VarChar	
MatchPrimaryWord4	VarChar	
MatchPrimaryWord5	VarChar	

Selected Column(s):

Source	Name
Data	MatchPrimaryName

Command Options

☐ Reverse
 ☒ Columns
 ☒ Arrays

m-prob:
 u-prob:
 Param 1:
 Param 2:

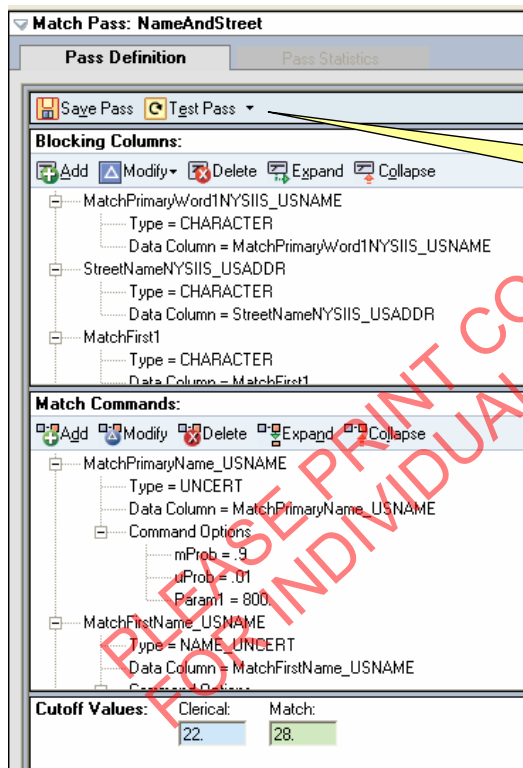
Override Weights

Apply 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.

Fully configured pass



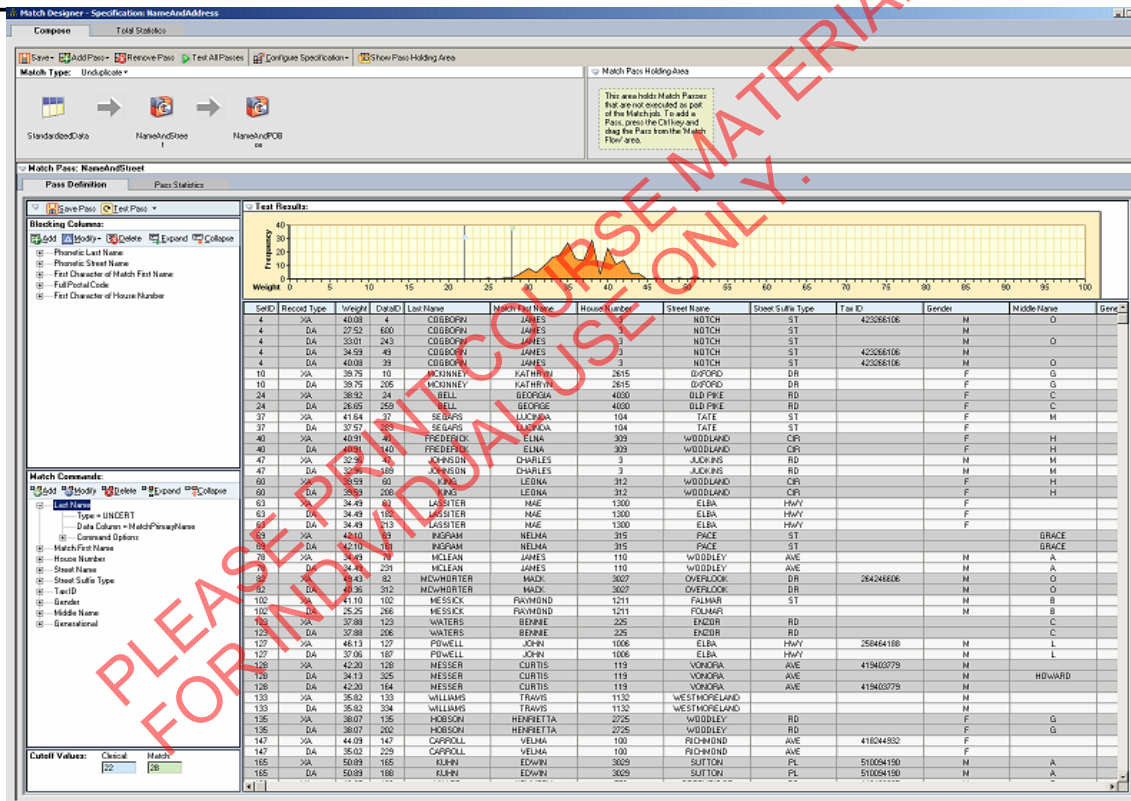
Click test pass to run the pass against the data

Expanded view will show details of blocking and match commands

© 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.

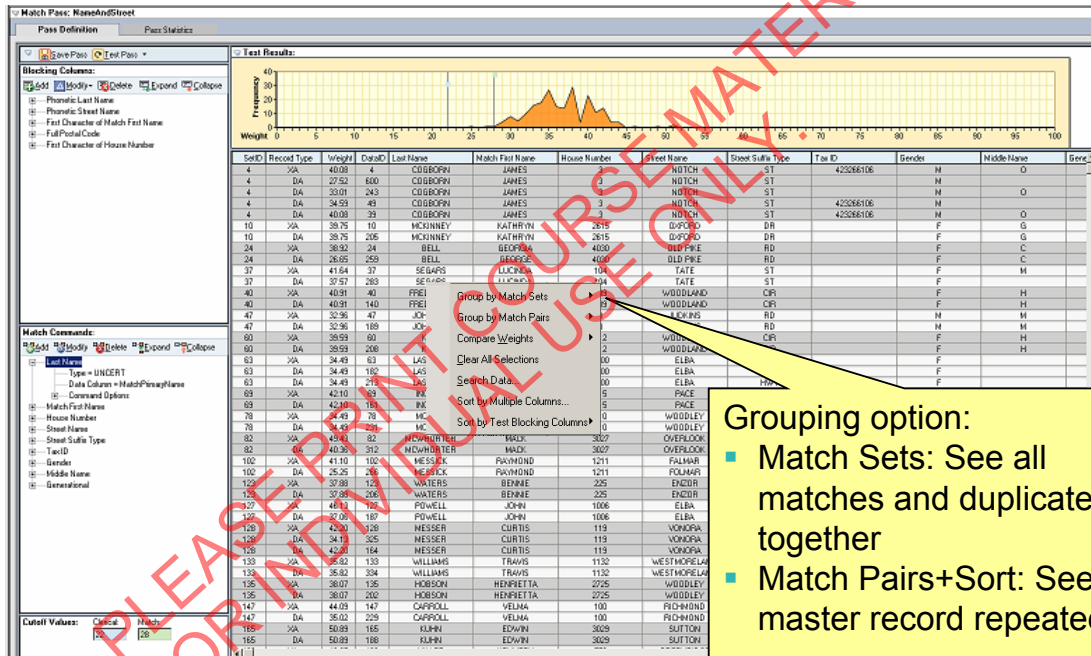
Match design – after test pass run



© 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 design - 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.

Match design - unduplicate

Default Display (Grouped by Match Sets)

SetID	Record Type	Weight	DataID	MatchPrimaryName_US	MatchFirstName_USNA	HouseNumber_USADD	StreetName_USADDR	StreetSuffixType_L
314	XA	40.08	314	COGBORN	JAMES	3	NOTCH	ST
314	DA	33.01	59	COGBORN	JAMES	3	NOTCH	ST
314	DA	40.08	349	COGBORN	JAMES	3	NOTCH	ST

Grouped by Match Pairs and then sorted Ascending by Weight

Group by Match Sets	1	CARROLL
Group by Match Pairs	Exclude Residuals	
Compare Weights	Include Residuals	
Clear All Selections		
Search Data	357	JOHNSON
	18	HOBSON
	445	HOBSON

SetID	Record Type	Weight	DataID	MatchPrimaryName_US	MatchFirstName_USNA
2	XA	42.30			WILLIAM
2	DA	33.23			WILLIAM
4	XA	50.89			EDWARD
4	DA	50.89			
5	XA	32.96	357	JOHNSON	
5	DA	32.96	357	JOHNSON	

Add to Blocking Fields	
Add to Match Commands	
Sort	Ascending
	Descending

SetID	Record Type	Weight	DataID	MatchPrimaryName_US	MatchFirstName_USNA	HouseNumber_USADD	StreetName_USADDR	StreetSuffixType_L
5	XA	32.96	357	JOHNSON	CHARLES	3	JUDKINS	RD
5	DA	32.96	357	JOHNSON	CHARLES	3	JUDKINS	RD
314	XA	40.08	59	COGBORN	JAMES	3	NOTCH	ST
314	DA	33.01	59	COGBORN	JAMES	3	NOTCH	ST
582	XA	33.15	583	WINDHAM	AMY	2219	GEORGIA	RD
582	DA	33.15	583	WINDHAM	AMY	2219	GEORGIA	RD

SetID	Record Type	Weight	DataID	MatchPrimaryName_US	MatchFirstName_USNA	HouseNumber_USADD	StreetName_USADDR	StreetSuffixType_L
314	XA	40.08	349	COGBORN	JAMES	3	NOTCH	ST
314	DA	40.08	349	COGBORN	JAMES	3	NOTCH	ST

© 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 design - unduplicate

PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

■ Compare Weights:
 ■ See how any two records score

SetID	Record Type	Weight	DataID	MatchPrimaryName_US	MatchFirstName_USNA
314	XA	40.08	314	COGBORN	JAMES
314	DA	33.01	59	COGBORN	JAMES
314	DA			COGBORN	JAMES
334	XA			BELL	GEORGIA
334	CP			BELL	GEORGE
350	XA				
350	DA				
379	XA				
379	DA			NGRAM	NELMA

Group by Match Sets ▶
 Group by Match Pairs ▶
 Compare Weights ▶
 Clear All Selections
 Search Data

Based on Last Match Run
 Based on Current Match Settings

Match Weights

	GenderCode_USNAME	MiddleName_USNAME	NameGeneration_USNA	MatchFirstName_USNA	MatchPrimaryName_US	HouseNumber_USADD	StreetName_USADD
Record 1:	M	0		JAMES	COGBORN	3	NOTCH
Record 2:	M	0		JAMES	COGBORN	3	NOTCH
Weight	0.80	5.49	0.00	4.32	6.75	5.90	6.26

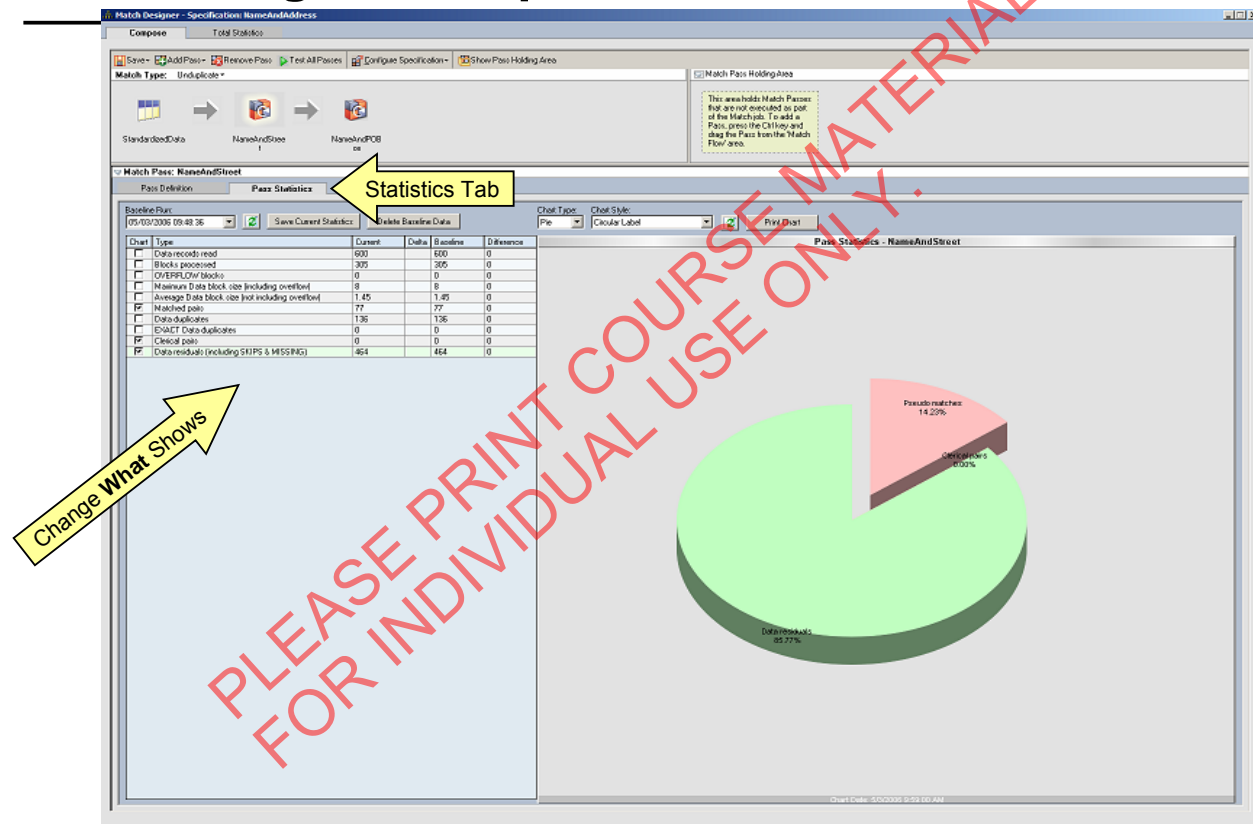
Composite Weight: 33.01

OK 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.

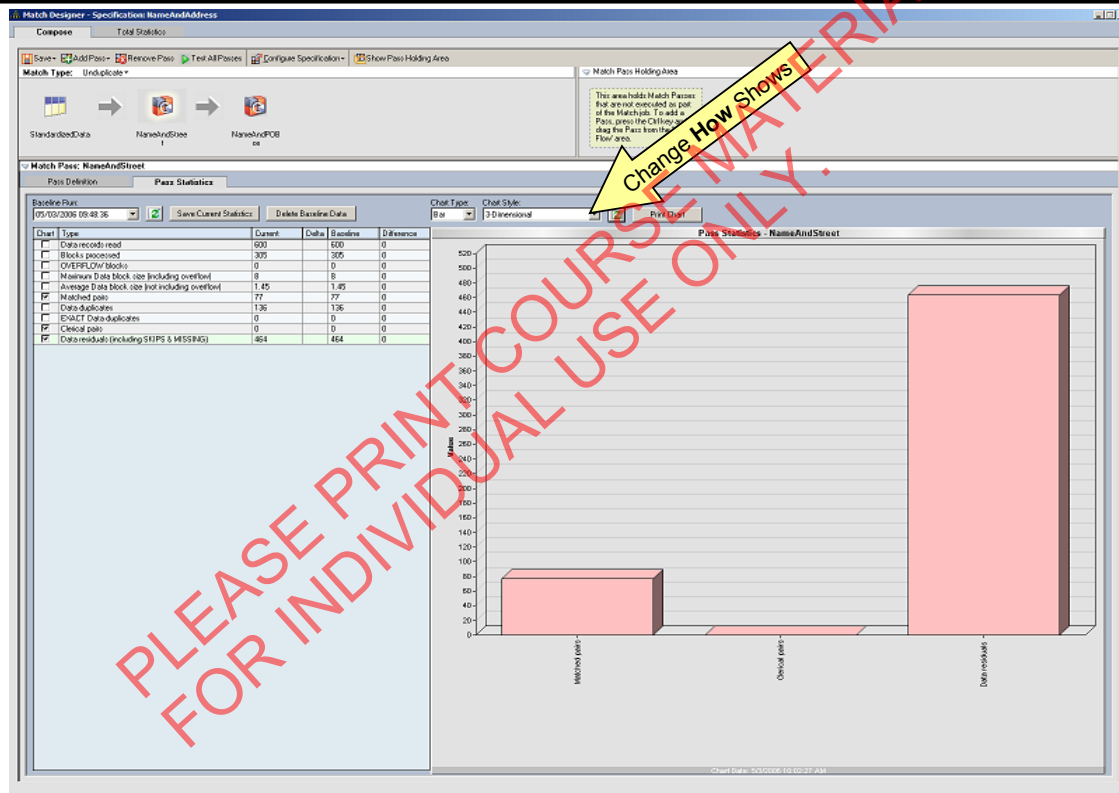
Match design - unduplicate



© 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.

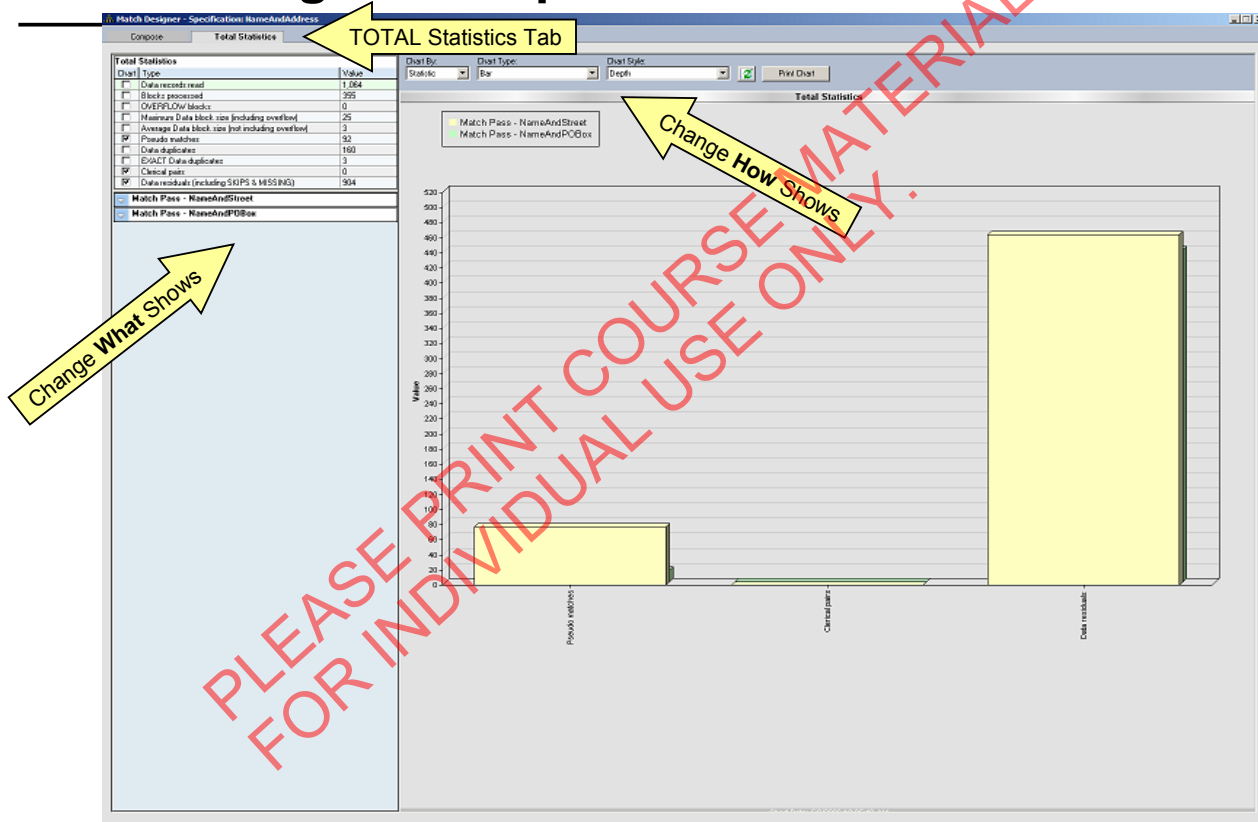
Match design - unduplicate



© 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.

Match design - 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.

Exercise 18: Configure test results database

- Build a DB2 database to contain match test results
- Build an ODBC source to connect the database to QualityStage

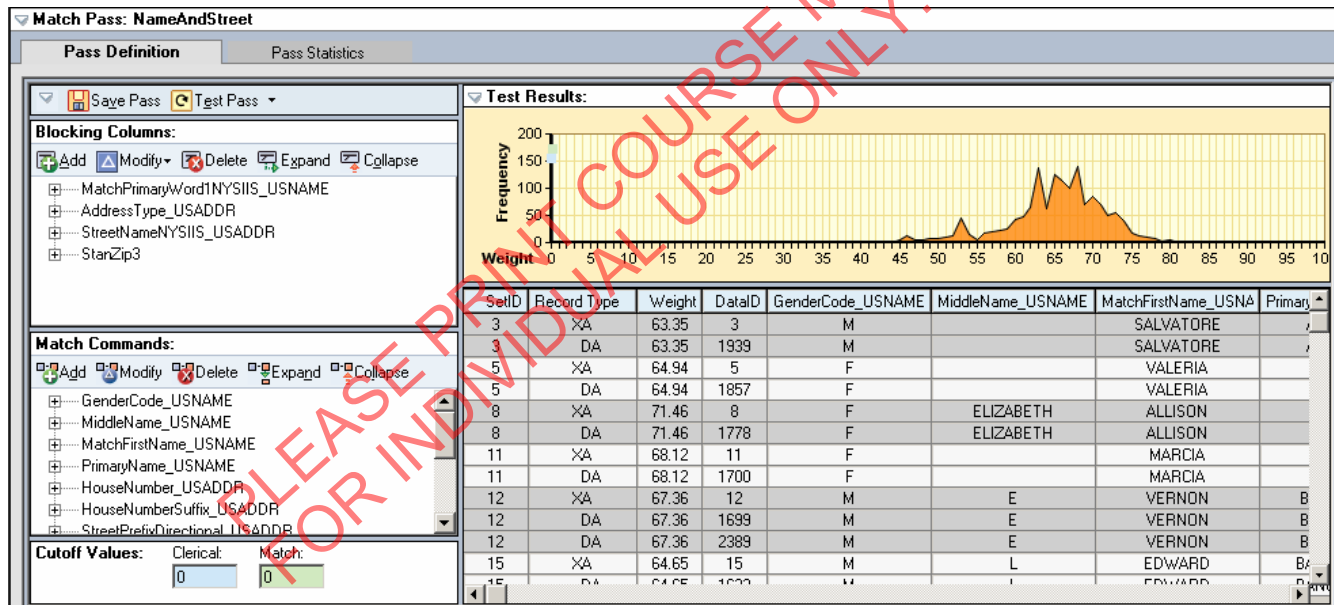
PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

© 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.

Exercise 19: Match specification

- Use Match Designer to build specification for unduplicate job
- Configure test results 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.

Match improvement strategy

1. Set critical values for important fields
2. Review calculated weights
 - Adjust weights using weight overrides
3. Set cutoffs
4. Add additional passes

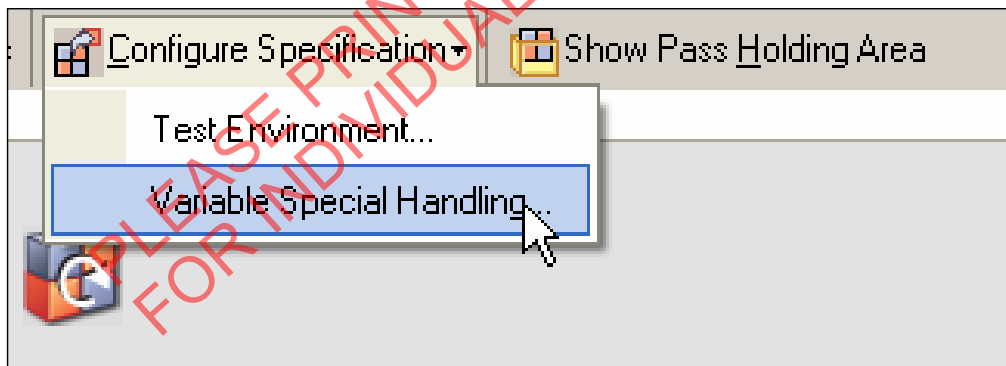
PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

© 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.

Critical fields

- Used to identify fields that must agree in order for records to be linked
 - Critical – Fields values must agree exactly or the records cannot be linked (considered a match)
 - Critical Missing OK – Field values must agree exactly on values not considered “missing values”
- QualityStage feature: VARTYPE



© 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.

Variable special handling

Variable Special Handling

Select the desired Actions and corresponding Columns that require special handling...

Define Special Handling

Actions: (click arrow to view) • Data Columns • Reference Columns Table Definition Name: stan

Available: (click arrow to view)

Name	Description
Applicant5	CRITICAL MISSINGOK
MatchFirst	CLERICAL
HouseNum	CLERICAL MISSINGOK
ZipCode3	NOFREQ
UniqueIdentifier	CONCAT
Name	
AddressLine1	
AddressLine2	
City	
State	
Zip5	

Column1:

Summary of Special Handling

Action	Source	Column Name(s)
CRITICAL	Data	ApplicantSSN

Add Remove

Cancel OK 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.

Weight overrides

- Allows you to adjust both the agreement and/or disagreement weights for specific situations
 - Add to calculated weight
 - Replace weight

On Match Commands
screen

The screenshot shows a software interface for configuring match commands. It includes a section for selecting columns, a 'Reverse' checkbox, and a 'Command Options' panel with input fields for 'm-prob', 'u-prob', 'Param 1', and 'Param 2'. A 'Weight Overrides' button is at the bottom right.

Selected Column(s):

Source	Name
Data	MatchPrimaryName_USN/

Command Options

m-prob:

u-prob:

Param 1:

Param 2:

© 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.

Weight override screen

Enter one or more values for each weight override...

Compose Weight Override

☒ Replace Agreement Weight [AW]: Data Source Missing Weight [AM]:
☐ Add Disagreement Weight [D'W]: Reference Source Missing Weight [BM]:
Conditional Data Source Value [AV]: Both Missing Weight [XM]:
Conditional Reference Source Value [BV]:

Summary of Weight Overrides

A/R	AV	BV	AW	D'W	AM	BM	XM

© 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.

Cutoffs

- There are two cutoffs
 - Match cutoff (high cutoff)
 - Clerical cutoff (low cutoff)
- Records with a weight equal to or above the Match cutoff are considered matches
- Records with a weight below the low cutoff are not matches
- Records with a weight greater than or equal to the low cutoff and less than the high cutoff are considered clerical records for manual review
- Cutoffs can be set at the same value eliminating clerical records

© 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.

Setting the match cut-off

	Weights	Data fields					
Definite Match	27.82	PO	BOX	930202			
	27.82	PO	BOX	930202			
	27.82	PO	BOX	930202			
Definite Match	38.65	35	COLLIER	RD	NW	STE	610
	38.65	35	COLLIER	RD	NW	STE	610
Questionable Match	25.81	928	S	1ST	ST		
	14.45		S	1ST	ST		

© 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.

Multiple match passes

- Additional passes are helpful in overcoming data errors and missing values in block fields
- You should always create at least two match passes
- Change blocking strategies for each pass

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.

Note: You may create up to seven match passes. Usually 2-3 are sufficient.

Example: multiple match passes

Pass	Weights					
Data fields						
1	26.31	JASON BIRCH	1350 WALTON	WAY	30901	
1	26.31	JASON BIRSH	1350 WALTON	WAY	30901	
1	20.42	JOHN SMITH	2047 PRINCE	AVE	30604	
1	10.83	MARY SMITH	2047 PRINCE	AVE	30604	
1	RES A	JOHN SMITH	P.O. BOX 123		30604	
<hr/>						
2	20.42	JOHN SMITH	2047 PRINCE	AVE	30604	
2	10.19	JOHN SMITH	P.O. BOX 123		30604	

- Pass 1 blocked on street name
- Pass 2 found additional matched records in which the street name was different but the names were the same

© 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 w

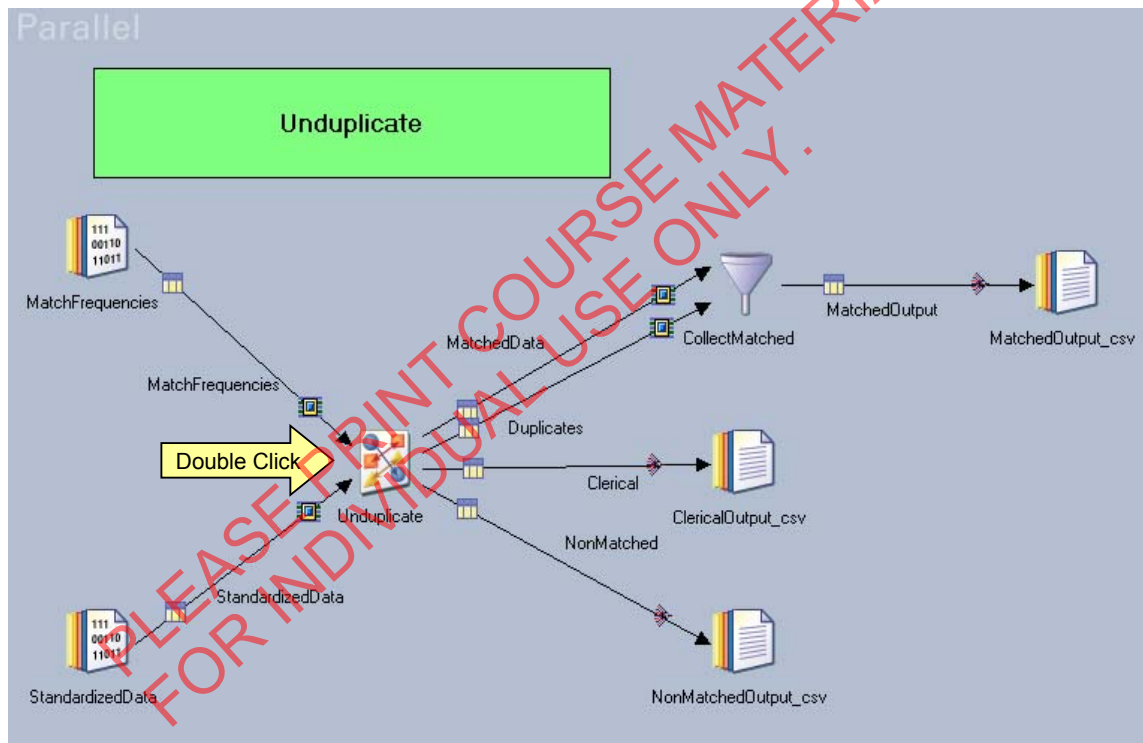
Match Implementation – Unduplicate job

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.

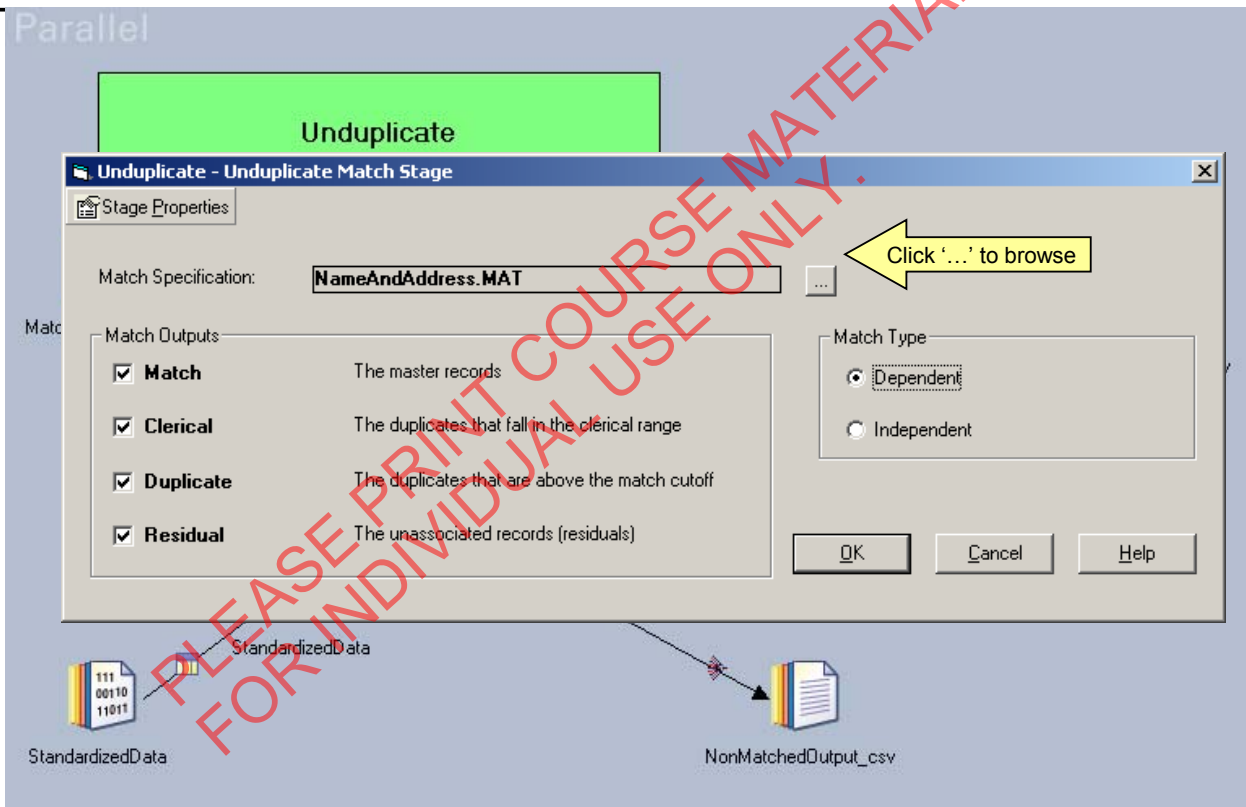
Unduplication Implementation



© 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.

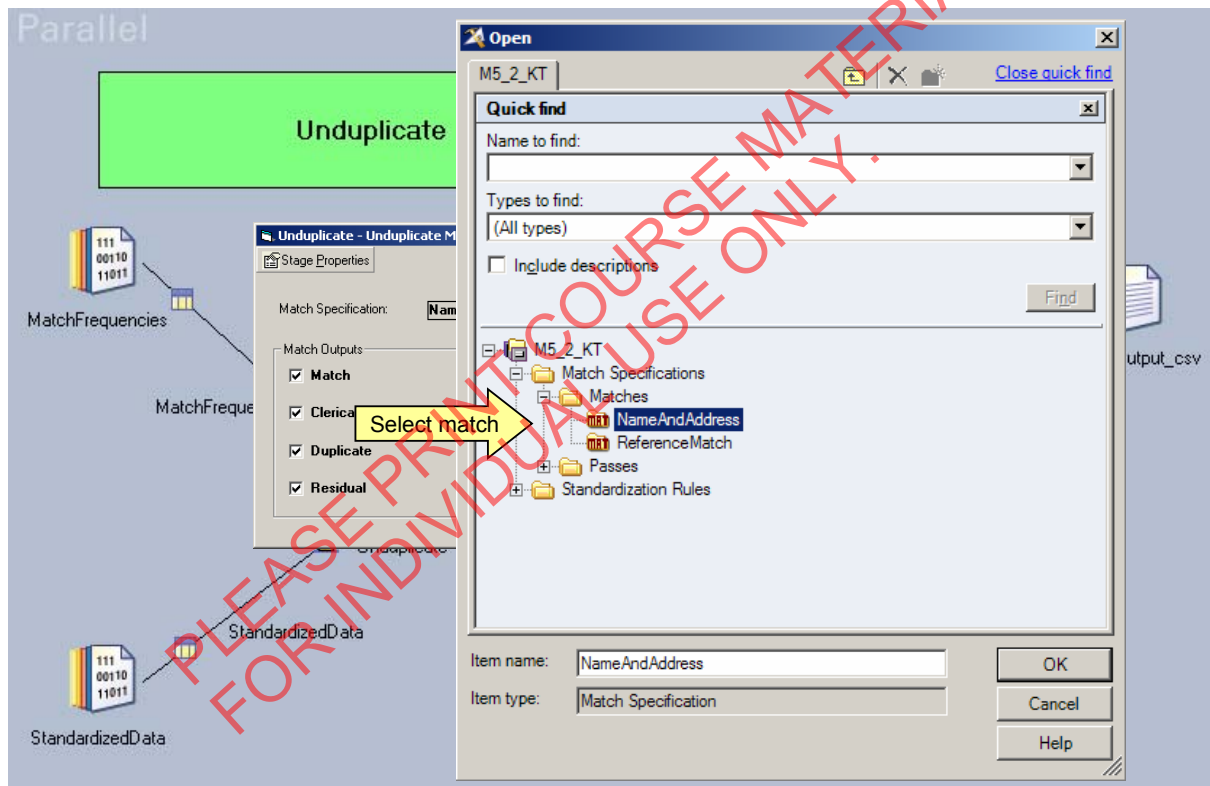
Unduplication Implementation



© 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.

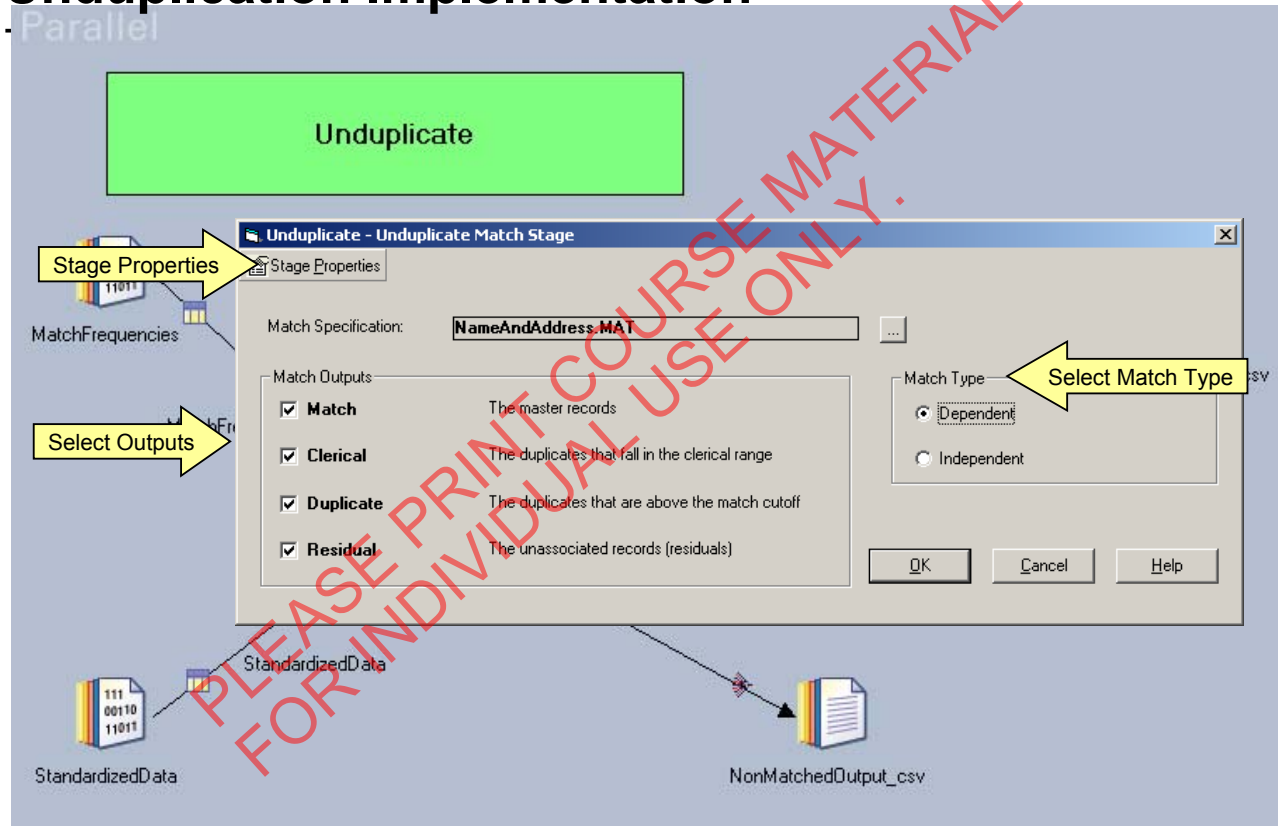
Unduplication Implementation



© 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.

Unduplication Implementation

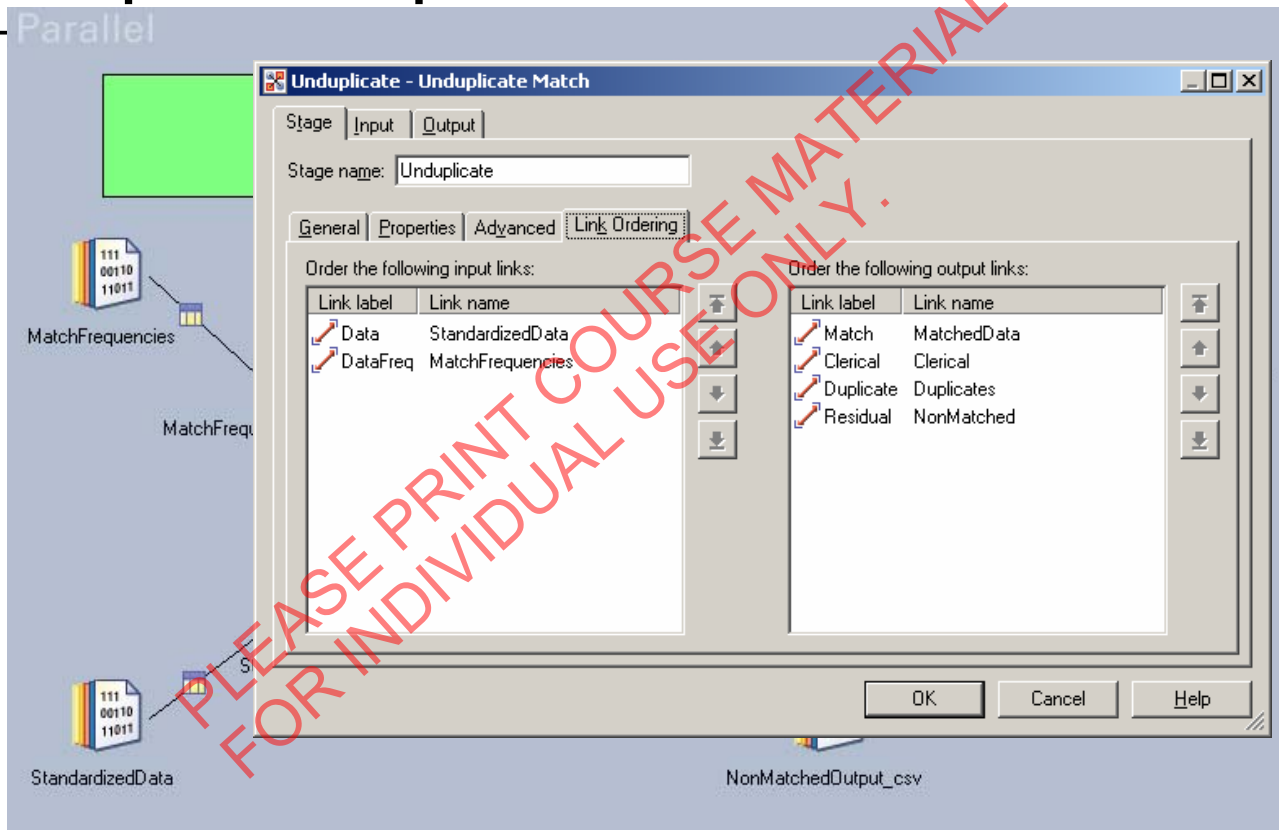


© Copyright IBM Corporation 2007

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.

Unduplication Implementation

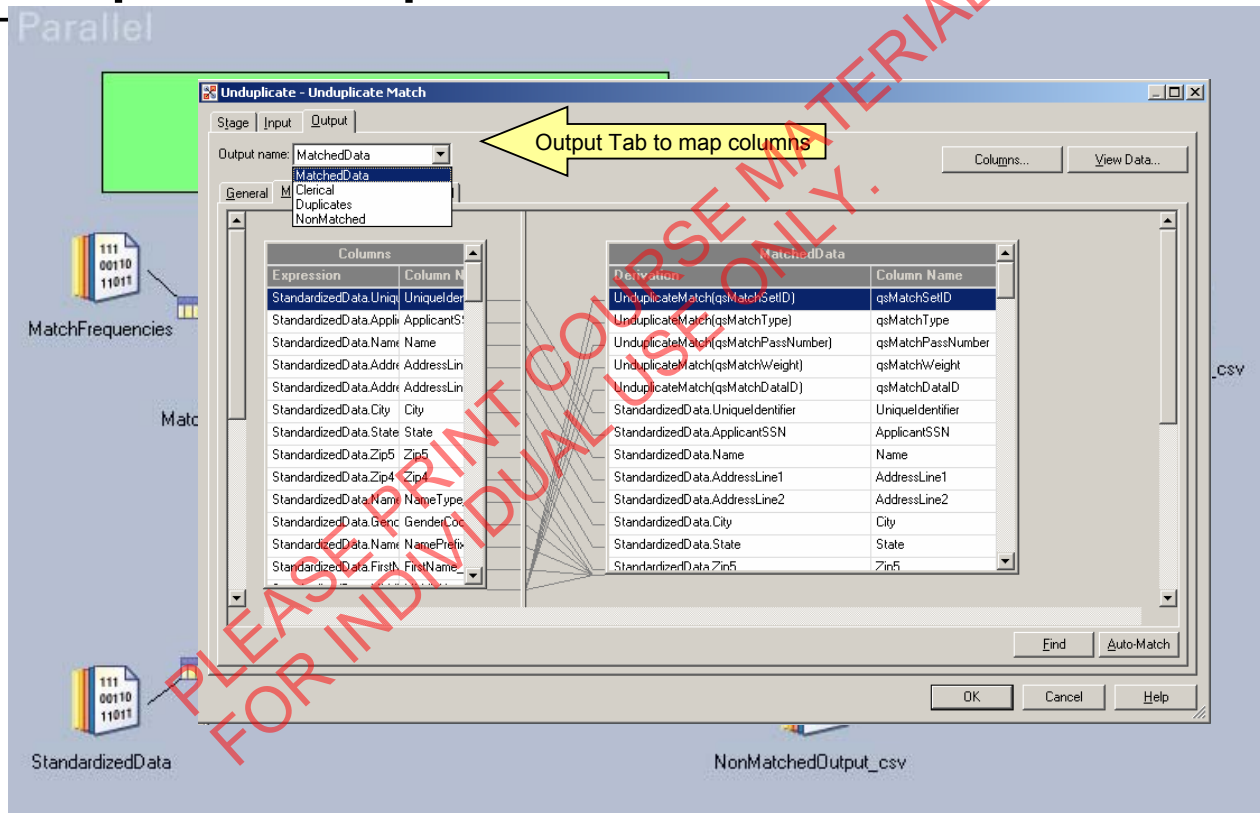
Parallel



© 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.

Unduplication Implementation



© 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) Match specifications are created using Designer.
2. (T/F) An unduplicate match can be used against two files.
3. Which match specification component determines the extent of the clerical review records?

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) Match specifications are created using Designer.
Answer: True
2. (T/F) An unduplicate match can be used against two files.
Answer: False
3. Which match specification component determines the extent of the clerical review records?
Answer: cutoff values

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 where Match fits in the Data Re-engineering Methodology
- Describe QualityStage Match concepts
- Define the type of matching algorithms
- Describe the importance of blocking
- Apply multiple match passes to increase efficiency/efficacy
- Interpret and improve match results

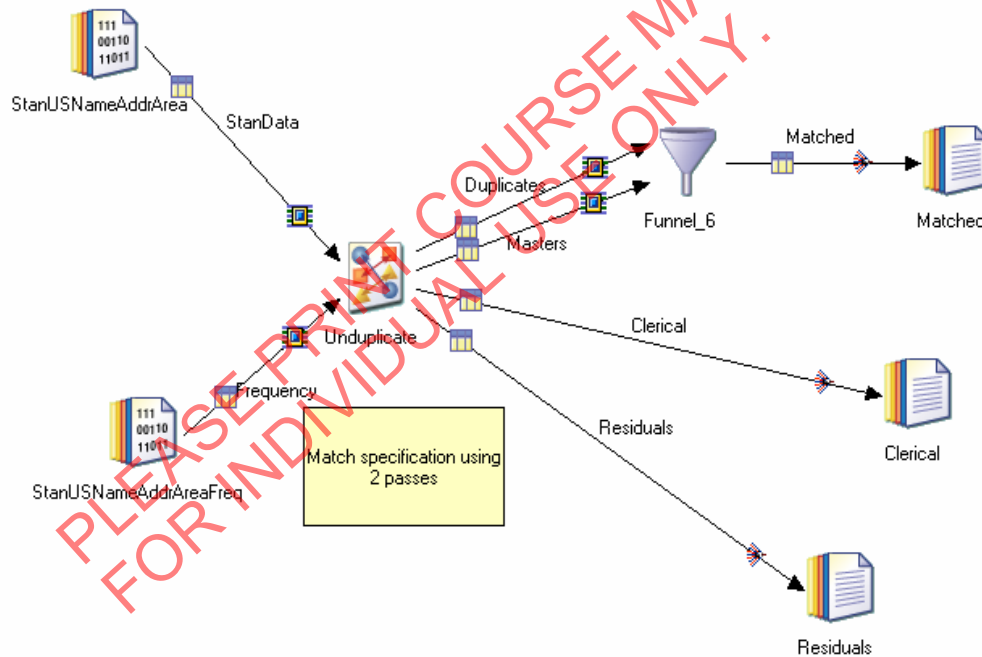
PLEASE PRINT COURSE MATERIALS.
FOR INDIVIDUAL USE ONLY.

© 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 20: Unduplicate

- Build unduplicate job using the match specification



© 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.