

IBM Software Group

DB2 Information Management Technical Conference: F04 – *Building a Test Database System*

David R. Schwartz
dschwar@us.ibm.com

DB2. Data Management Software



 e-business software

IBM DB2 Test Database Generator

- A powerful tool that provides several methods of generating test data for DB2
 - Can generate test data on demand
 - Can create scripts that can be run now, later, and repeatedly
- Supports multiple interfaces
 - ISPF, GUI, Offline
- Supports DB2 across multiple platforms
 - z/OS, Linux, Unix, Windows



Test Database Generator

For when you need to...

- Create test data in new or existing tables
- Debug application failures that are data dependent
- Copy a slice of data instead of all of the data
- Mask and censor sensitive data for testing
- Create a restructured database for testing
- Create test data in a variety of output formats



Agenda

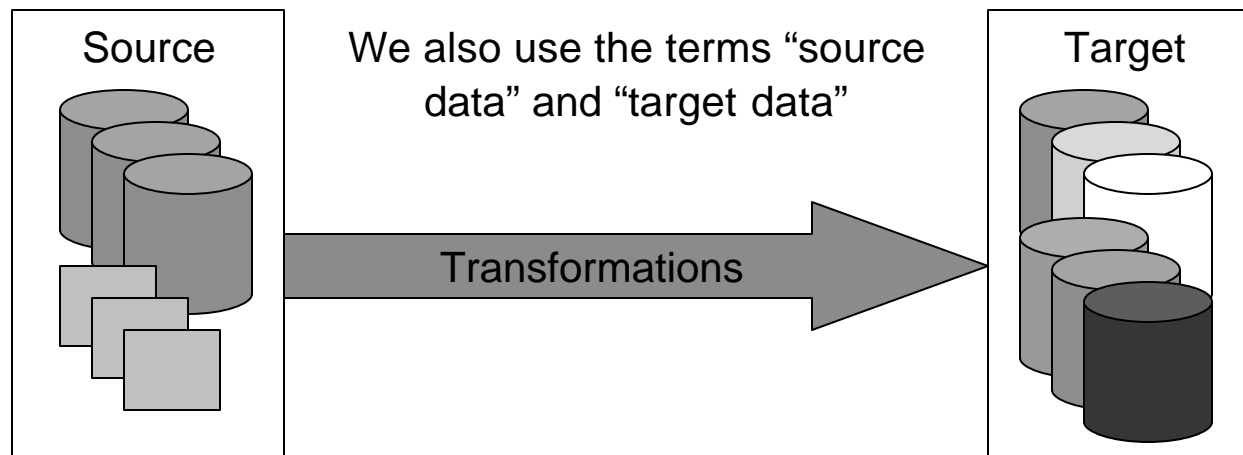
1. An introduction to the Test Database Generator computing model
2. An overview of the Test Database Generator architecture
3. Some examples of generating test data using Test Database Generator

Note: Test Database Generator = TDBG

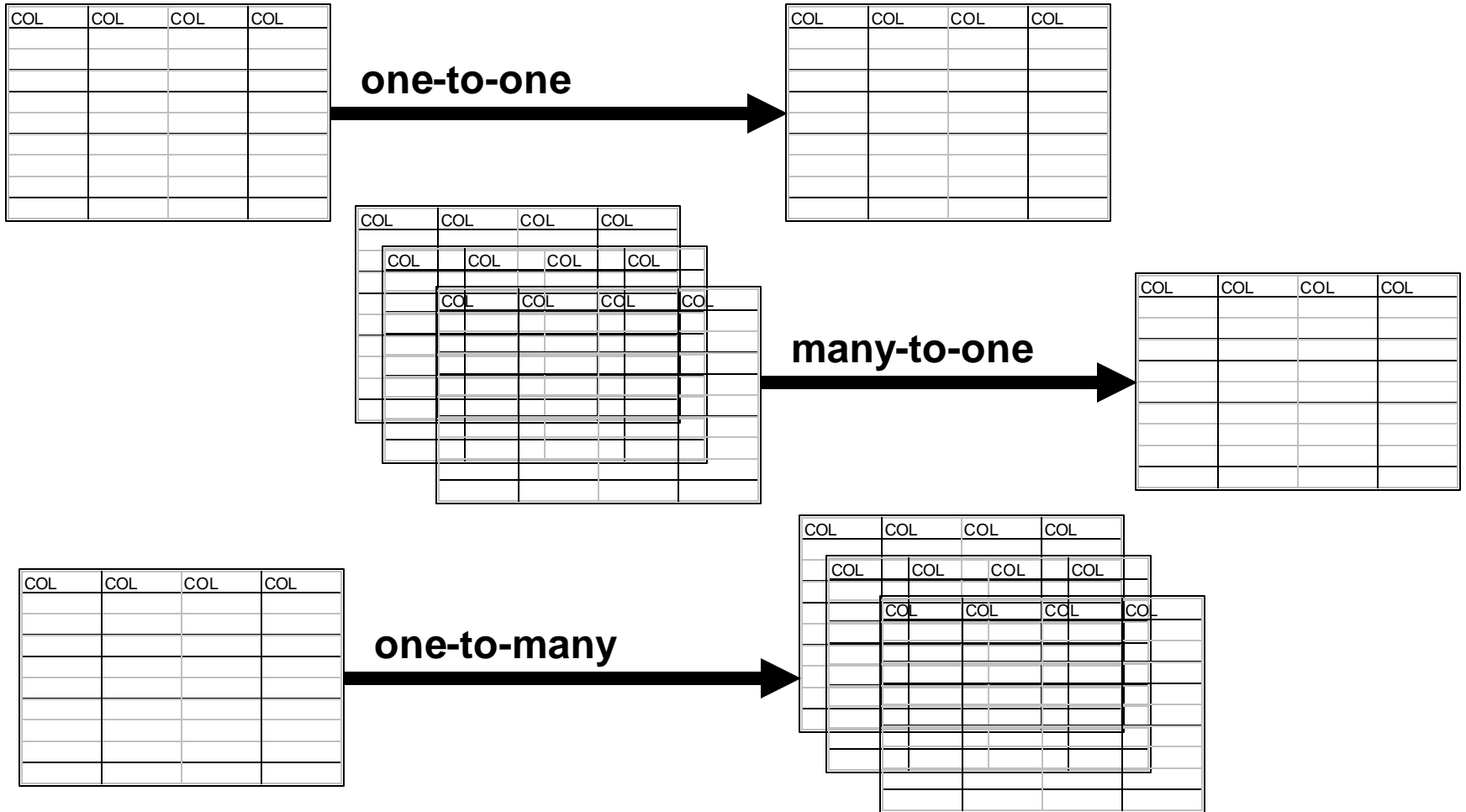


Introducing the TDBG computing model

1. Start with data that exists somewhere in your enterprise
2. Leverage knowledge of data relationships
3. Apply transformation rules
4. Create test data

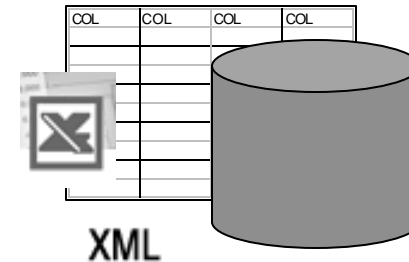


Start with data that exists means more than just make a duplicate copy ...



Data sources

- Data sources can be database tables or files
- Database tables supported include
 - DB2
 - Informix (future)
- File types supported include
 - XML file (future)
 - CSV (comma separated value) file
 - Text file with any delimiter
 - Text file of fixed width



How data is accessed

Various protocols are supported, including...

- JDBC
- HTTP
- FTP
- File

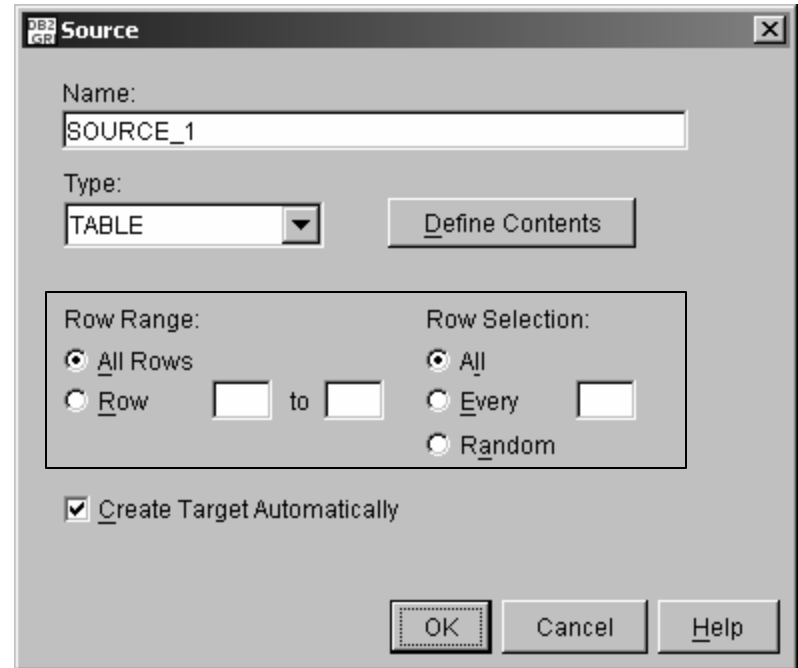
Data can be specified directly

- You can specify data values (rows and columns) directly or copy/paste data directly into TDBG using the GUI



Identifying the data you want to start with

- Multiple levels of filtering
- WHERE clause to limit result set
 - When source is a table
- Range
 - All
 - From/to
- Row selection within the range
 - Every row sequentially
 - Every nth row
 - Every row randomly



The screenshot shows the 'Source' dialog box in DB2. The 'Name' field is set to 'SOURCE_1'. The 'Type' is set to 'TABLE', with a 'Define Contents' button next to it. The 'Row Range' section has 'All Rows' selected. The 'Row Selection' section has 'All' selected. The 'Create Target Automatically' checkbox is checked. The 'OK' button is highlighted with a dashed border.

DB2 Source

Name: SOURCE_1

Type: TABLE Define Contents

Row Range: All Rows Row [] to []

Row Selection: All Every [] Random

Create Target Automatically

OK Cancel Help

You don't have to start with existing data

- Starting with existing data is an optional step
- You can choose to direct TDBG to create test data “from scratch”



Introducing the TDBG computing model

1. Start with data that exists somewhere in your enterprise
2. Leverage knowledge of data relationships
3. Apply transformation rules
4. Create test data



Would you like to include related tables?

- TDBG supports grouping of objects
 - By catalog RI
 - By application
 - And more

- TDBG discovers, through a variety of methods, related tables

- You can also define related sets of tables and edit the discovered relationships



Copy related data

- Copy all rows in a related group
 - Copy all rows that are related across an entire related set of tables

- Copy a slice of data across a related set of tables
 - Start with a specific customer and copy that row and all rows related to that customer across all related tables

- Copy only those rows that are orphaned
 - Help troubleshoot applications that have RI maintained by an application, not DB2
 - TDBG locates children with no parents



Introducing the TDBG computing model

1. Start with data that exists somewhere in your enterprise
2. Leverage knowledge of data relationships
3. Apply transformation rules
4. Create test data



Transformation rules

- Transformation rules define the target test data
 - How to generate test data from source data
 - How to generate test data from scratch

- Examples
 - Create a target column PHONE which is the combination of a country code (derived from COUNTRY file), an area code from TABLE1, and a phone number from TABLE2
 - Create a target column ACCT_BALANCE which is a random number that falls within a specified range
 - Create a target column that is exactly the PIN column with the 3rd and 5th positions replaced (masked) with the letter X



Scopes and sets

- You define your test data one target column at a time
- The scope of a transformation rule set is target column
- Multiple rules can be specified for each target column
- Transformation rules are applied in order
- Each rule can modify, replace, append, or preface the previous value to allow for incremental building of a target column



Transformation rules

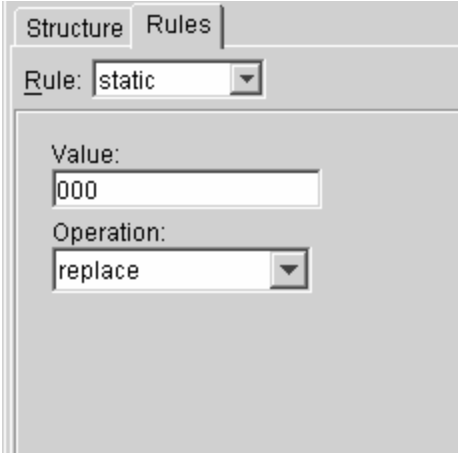
TDBG currently provides seven (7) transformation rules

1. Static Values
2. Source Column Values
3. Data Lookup
4. Data Masking
5. Expressions
6. Random Values
7. Pattern Generation



Rule 1: Static Values

- Specify static data value
- Source data not used as input for this rule
- Examples
 - Set target column STATE to CA
 - Set target column SPEED to 85
 - Set target column EXP_DATE to 2030-12-31



The screenshot shows a configuration window with two tabs: "Structure" and "Rules". The "Rules" tab is active. Below the tabs, there is a "Rule:" dropdown menu with "static" selected. Underneath, there are two fields: "Value:" with a text box containing "000", and "Operation:" with a dropdown menu showing "replace".

Rule 2: Source Column Values

- Generate target column value based on source column value
- No transformation (use source data as-is)
- Copies a column value from the currently selected row in a source object
- Used to perform a “copy” function
- Used to supply initial column value from a data source
- TDBG Auto-Target feature makes use of this rule
 - Source Column Values rule is set automatically for each target column



Rule 3: Data Lookup

- Provides a method to replace data based on table lookups
 - Uses the current value of the generated column as a key to the lookup
 - Specify replacement column
 - Example:
 - Source data has a product code which needs to be represented as a product name
- PROD_NAME = LOOKUP(PRODUCT_CODE in PRODUCT_TABLE)



Rule 4: Mask

- Replace or hide sensitive data
- Masks can be set with static text
 - Replace account number with a string of X's
 - Replace last four digits of License ID with 9999
- Mask can be set using a pattern rule (see rule #7)
 - Replace first character with a letter in the range A-F and then replace the next position with a number between 0 and 9
- Mask can be applied to the entire column or to a substring of the column



Rule 5: Expression

Expression provides a method to call database specific functions

- Supports any expression that can be evaluated by the DBMS in which TDBG is installed
 - String manipulation, calculations, etc.

- Examples
 - Calculate the current date / time / timestamp
 - Evaluate a mathematical expression
 - Target column = source column * 1.1



Rule 6: Random

- Generate a random value
- Allows for creating random date, time, timestamp, integer, and decimal values
- Specify min / max ranges for the generated values
- Randomly generated values are propagated across related tables
 - Only applies to primary/foreign key relationships that are system or user defined
 - A DEPT_CODE that is randomly generated shows up in both the DEPT table and the EMPLOYEE table



Rule 7: Pattern

- Generate data based on a specified pattern
- Numeric pattern
 - [0-9] evaluates to any single-digit number
- Character pattern
 - [A-Z]*3 evaluates to any three-character uppercase string
 - Character patterns are randomly selected at generation runtime
- String pattern
 - (Mrs|Mr|Ms) evaluates to 'Mrs', 'Mr', or 'Ms'.
 - (C[ATO] | A[KLR]) evaluates to 'CA', 'CT', 'CO', 'AK', 'AL', or 'AR'



Introducing the TDBG computing model

1. Start with data that exists somewhere in your enterprise
2. Leverage knowledge of data relationships
3. Apply transformation rules
4. Create test data



What can TDBG create?

TDBG can generate test data as a(n)...

- Comma separated value (.csv) file
- Text file with fixed width columns
- Text delimited file
- File of SQL INSERT statements
- DB2 for z/OS internal load format file
- XML file

TDBG can also...

- Generate test data and directly insert it into DB2



You can also choose...

- How many rows of test data to generate
- The exact structure of the target table(s)
 - Inherit from source (auto-target feature)
 - Specify column definitions
- To automatically generate DDL needed to create the target table(s)
 - Includes primary and foreign keys
 - Can be modified



TDBG architecture

- Components
- Data profiles
- Operating environments
- XML
- Unicode



TDBG components

- TDBG Server
 - Installed as a set of DB2 stored procedures
 - Can run on any DB2 running on z/OS, Linux, Unix, Windows
 - Connects to and reads from data sources
 - What databases can be accessed is governed by the product license
 - Is the component that reads Data Profiles and generates target test data

- TDBG Client
 - ISPF client running on z/OS
 - GUI (Java) client running on Linux, Unix, Windows
 - Connects to TDBG Server
 - Provides a user interface wizard and panels to guide you through the creation of a Data Profile



Data profiles

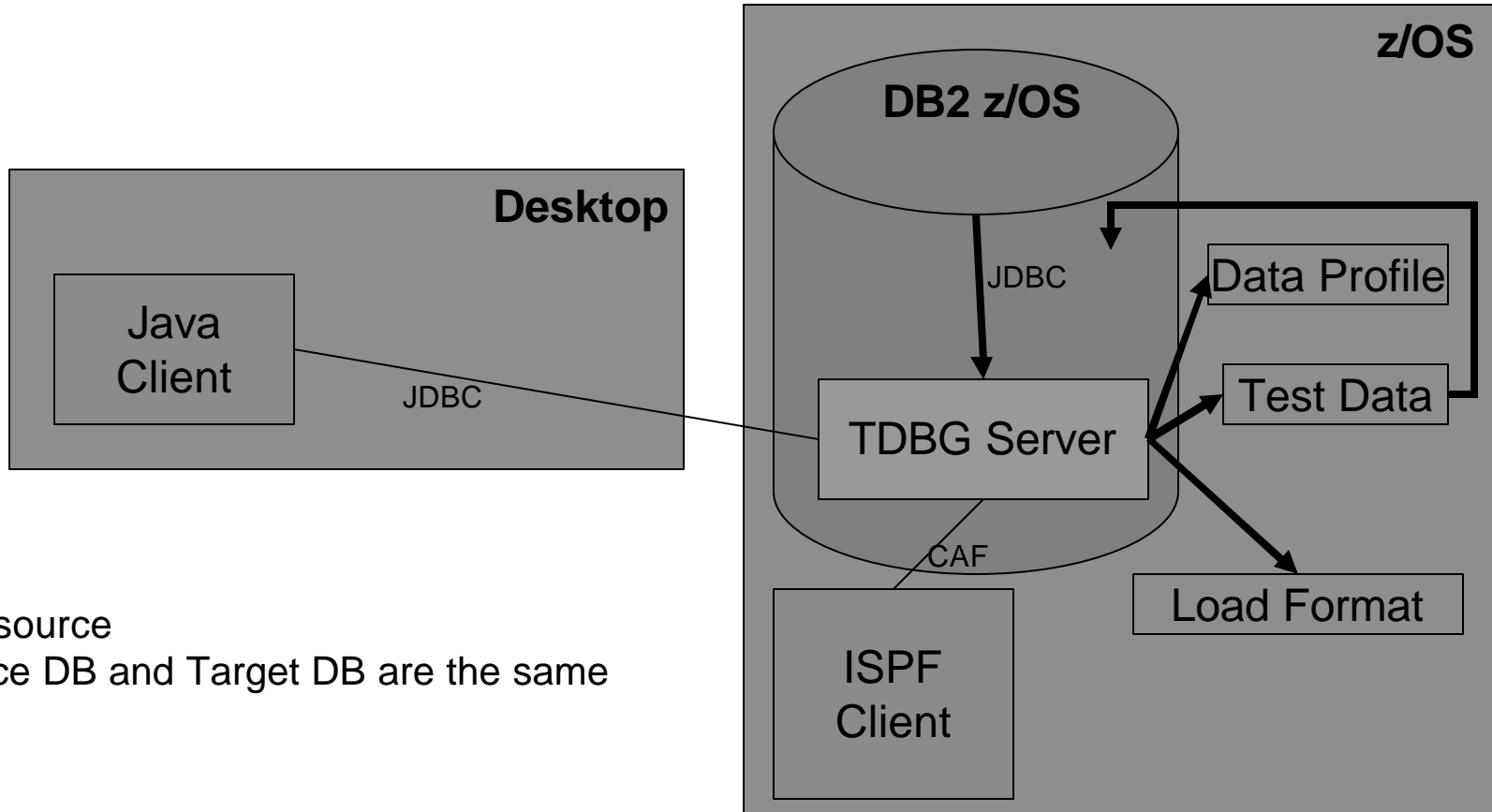
- Data Profiles define the generation process
 - Source data objects used to seed generation.
 - Definition and relational rules of the targets.
 - Describes how data will be copied, filtered, masked, and/or transformed.

- Data Profiles are stored in the HFS by the server

- Data Profiles are written in Test Database Generator Markup Language (a.k.a. GRIML)
 - GRI: Internal product code
 - ML: Markup Language
 - GRIML is an XML-based markup language



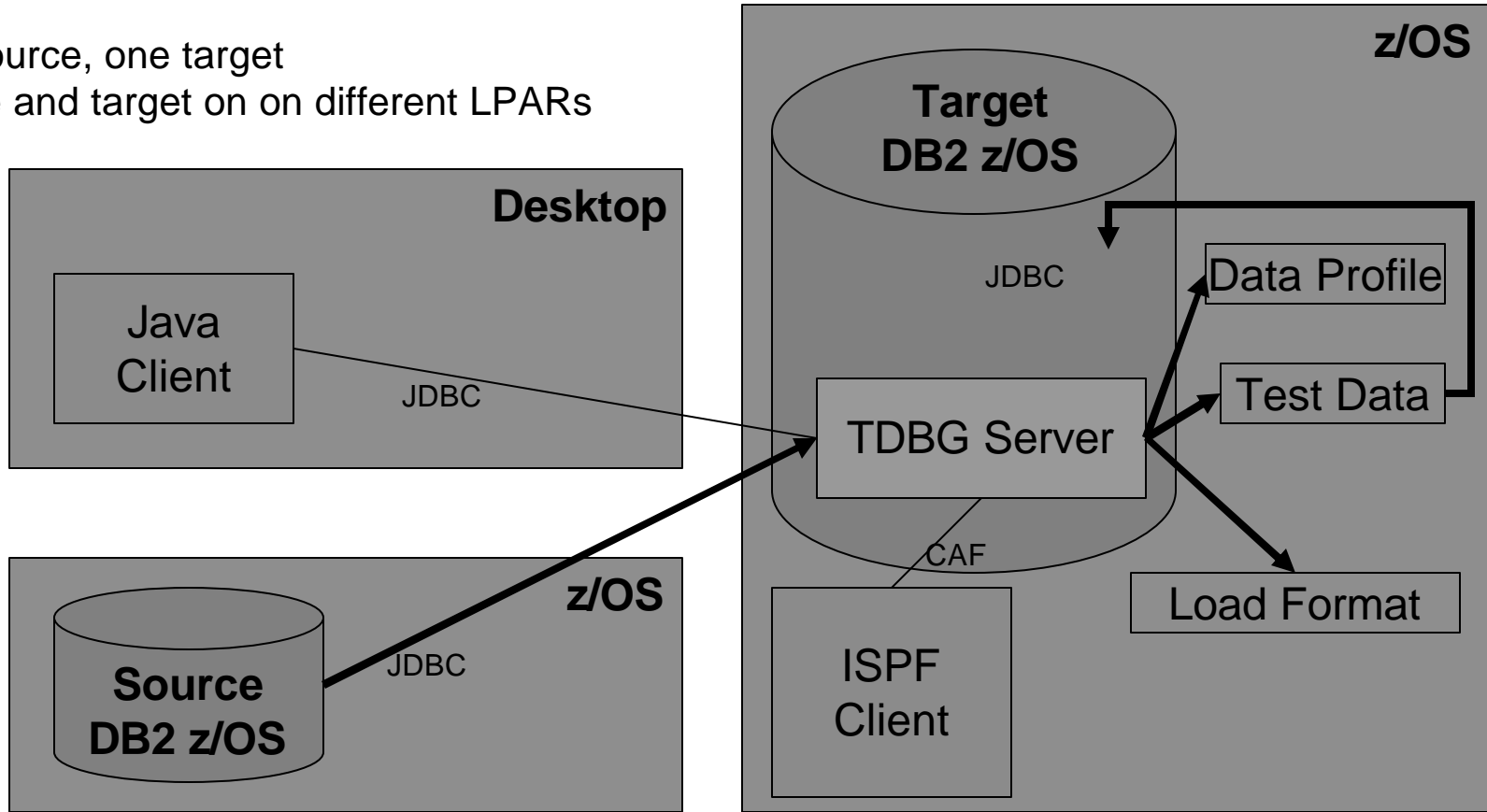
Operating environments: A z/OS example



- One source
- Source DB and Target DB are the same

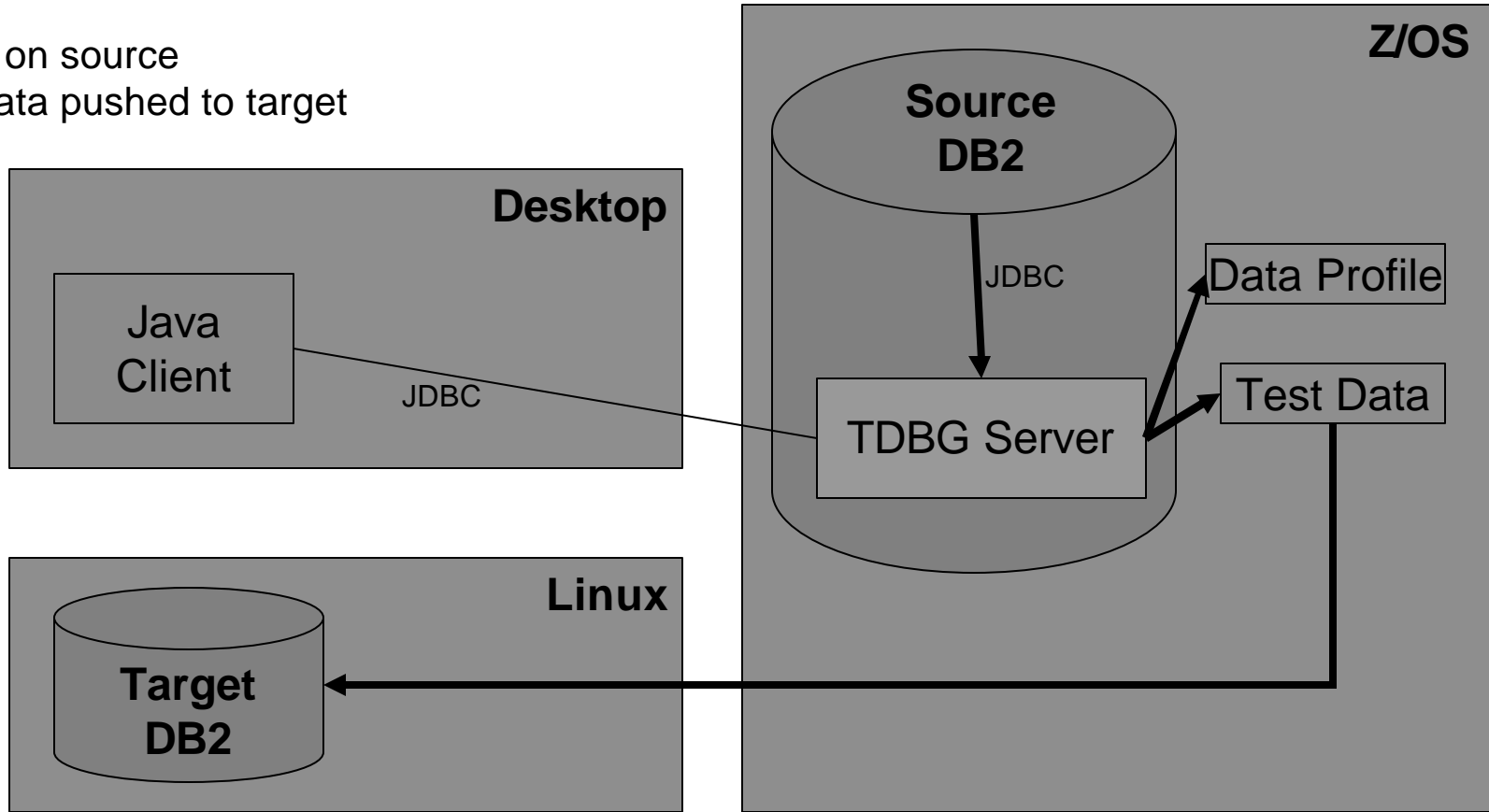
Operating environments: A z/OS example

- One source, one target
- Source and target on on different LPARs



Operating environments: A Multiplatforms example

- Server on source
- Test Data pushed to target



Operating environment: ISPF

```
GRI$MAIN  V2R1  -----  Test Database Generator  -----  2003/05/06
Option ==>
```

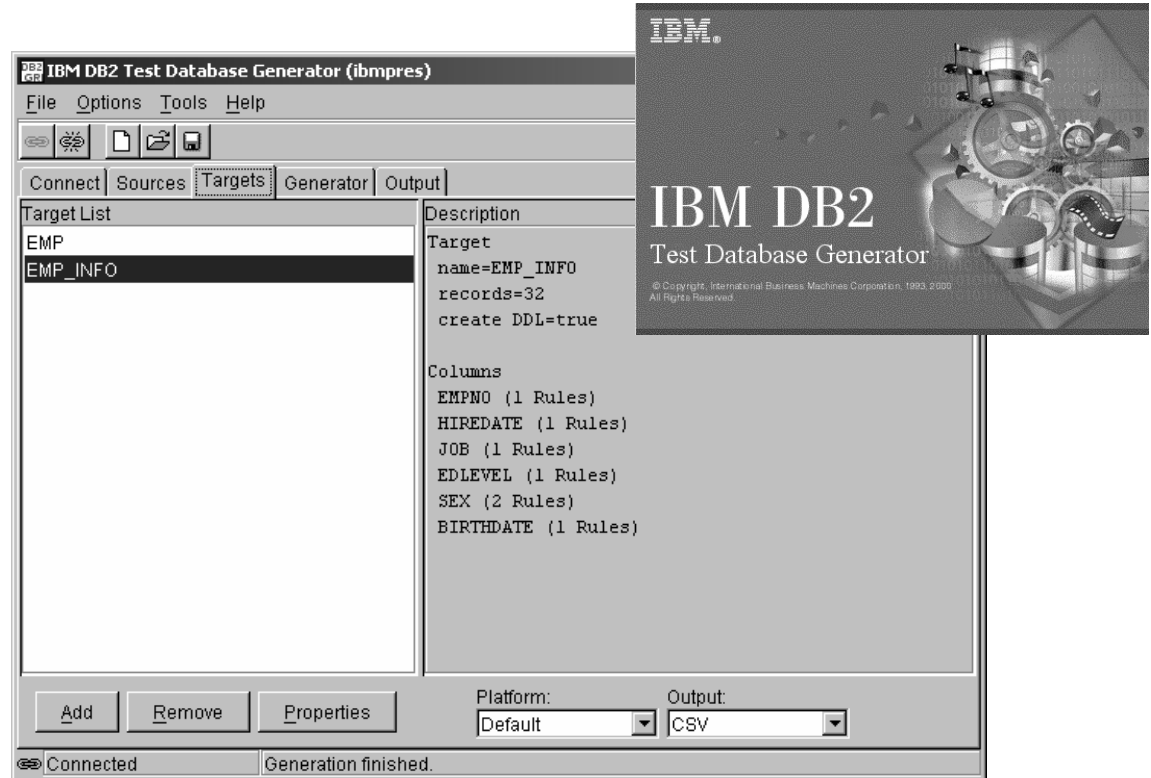
```
Current Server: W32B      Current SQLID XYZHAW      User:  XYZHAW
-----
```

Data Profile:

1	Sources	Sources: 0
2	Targets	Targets: 0
3	Generator	
4	Load Data Profile	
5	Save Data Profile	
6	Reset Data Profile	
S	Setup	
A	About	
X	Exit	

Operating environment: Java GUI

- Java application running on your desktop
- Can be launched from DB2 Control Center



TDBG and XML

- Data Profiles
 - Data Profiles are stored in an XML document
 - You can create a Data Profile manually using an XML editor

- XML data sources (future)
 - You can specify any row-and-column oriented XML document as a source of data
 - For example, a query result from DB2 Web Query Tool

- XML target
 - Test data can be generated into an XML document



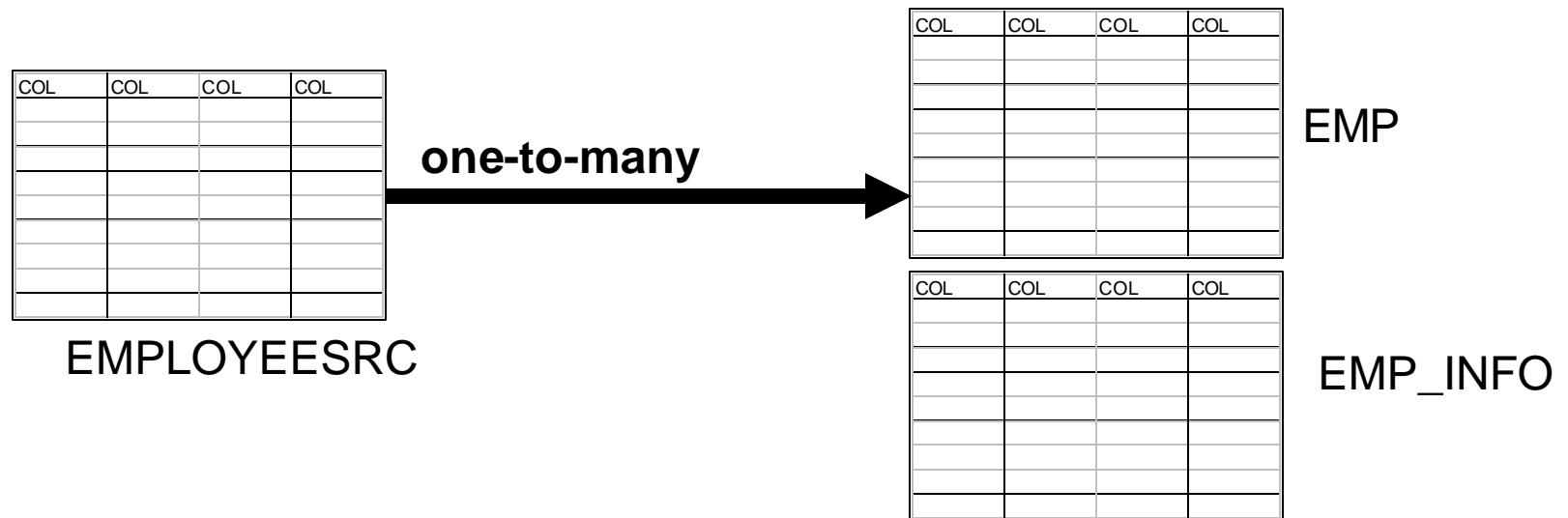
UNICODE

- Data profiles are written in UNICODE in UTF-8
- Targets are written in UNICODE in UTF-8
- Provides platform independence and portability for your test data
- All UTF-8 characters are supported



Use case

- Creating two new tables of test data
- Source data will be read from the EMPLOYEESRC table
- Two target tables (EMP and EMP_INFO) will be created



Use case

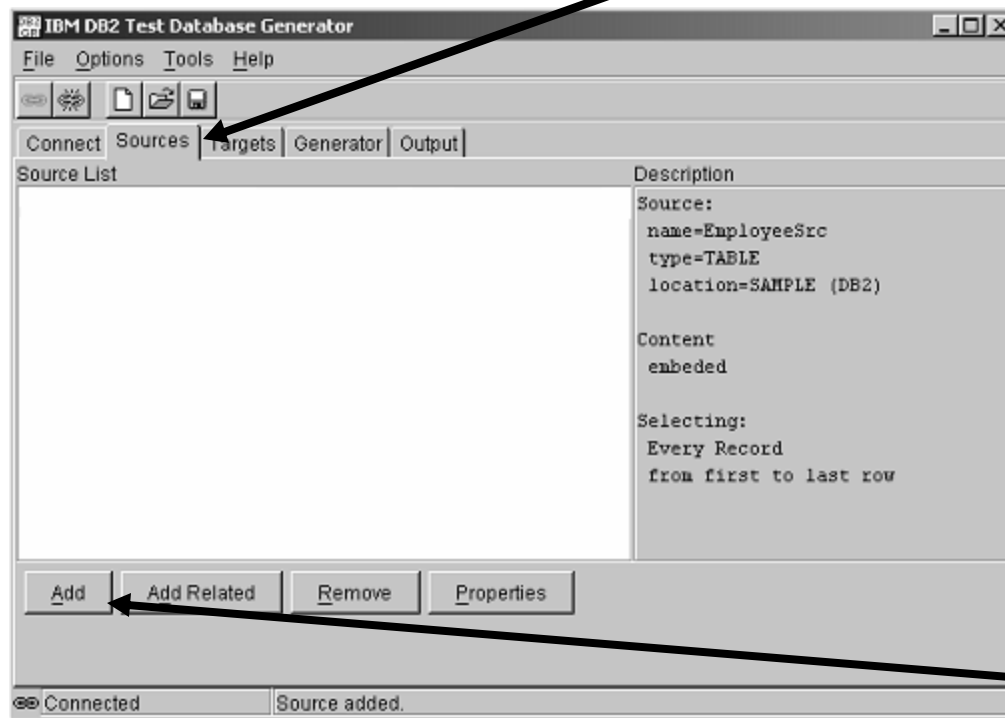
The steps we'll follow are

- Specify source
The EMPLOYEESRC table consisting of 14 columns
- Specify targets
Create two targets (EMP and EMP_INFO)
EMP will have the first 6 columns from EMPLOYEESRC
EMP_INFO will have the next 5 columns from EMPLOYEESRC (plus the column that relates the rows – EMPNO). The remaining columns in EMPLOYEESRC are not needed in the test environment.
- Specify rules
To copy specific columns and to protect sensitive information
- Generate test data
Watch a generation session



Use case : Specify source

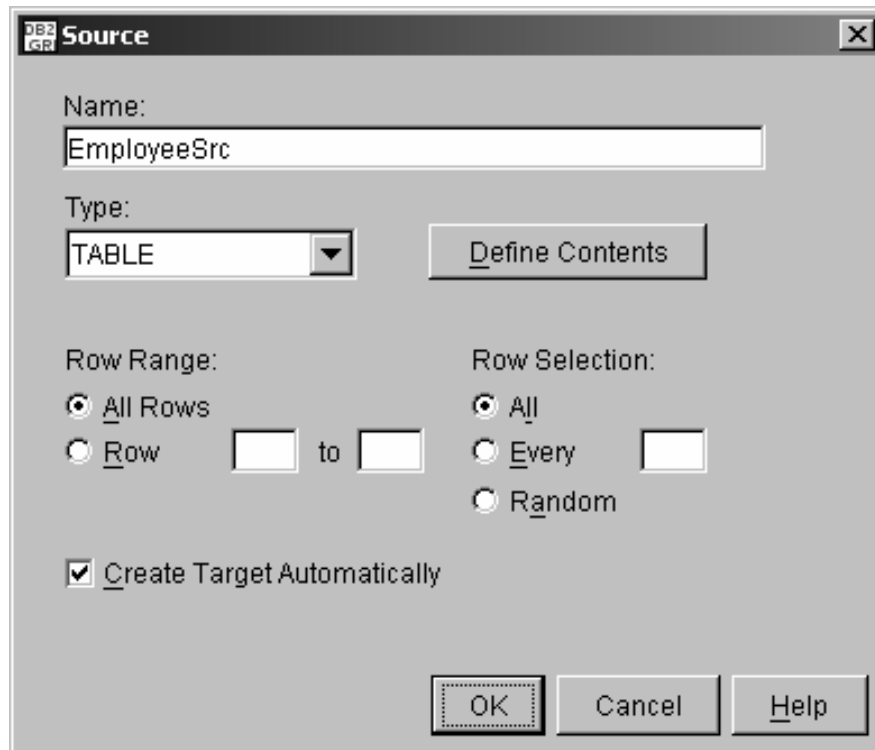
- Start from an empty Source List on the **Sources** tab



Add source using
the Add button

Use case : Specify source

- Define a source EmployeeSrc as “all rows from a TABLE”



The screenshot shows the 'Source' dialog box in DB2. The 'Name' field contains 'EmployeeSrc'. The 'Type' dropdown is set to 'TABLE'. The 'Row Range' section has 'All Rows' selected. The 'Row Selection' section has 'All' selected. The 'Create Target Automatically' checkbox is checked. The 'OK' button is highlighted.

DB2
Source

Name:
EmployeeSrc

Type:
TABLE

Define Contents

Row Range:
 All Rows
 Row [] to []

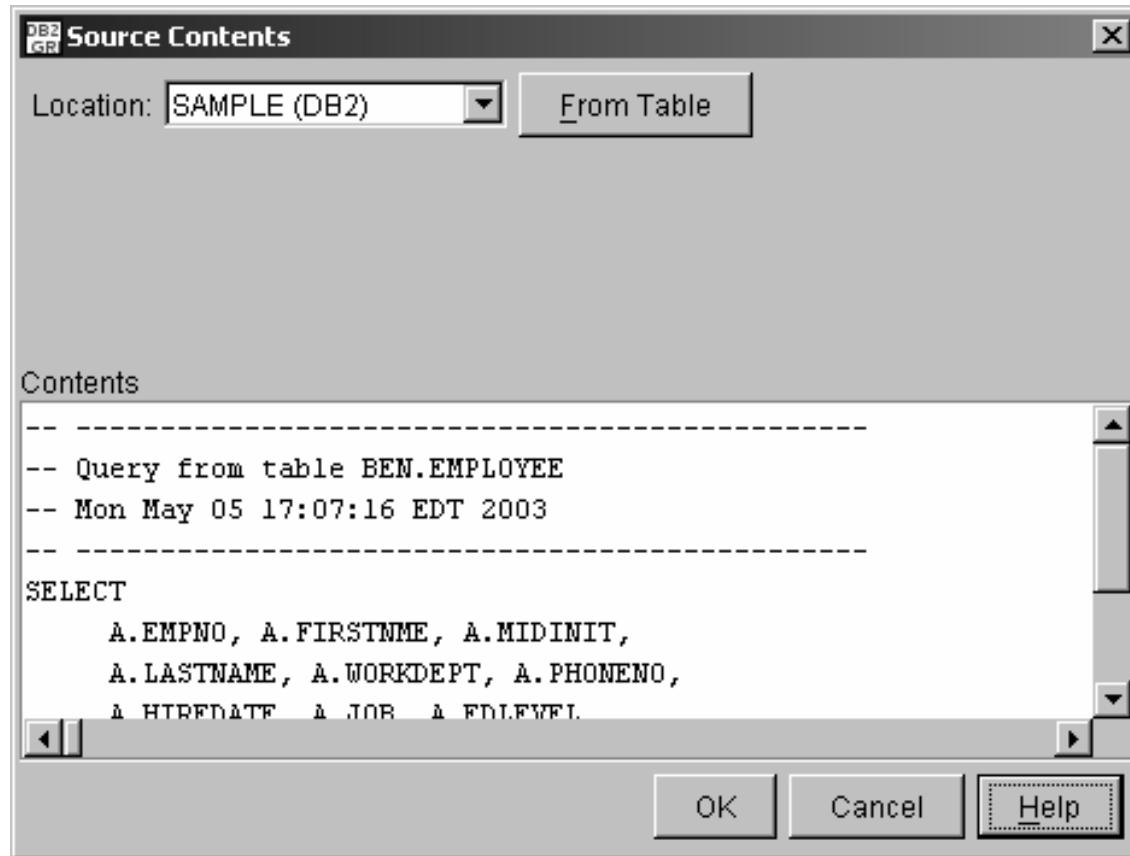
Row Selection:
 All
 Every []
 Random

Create Target Automatically

OK Cancel Help

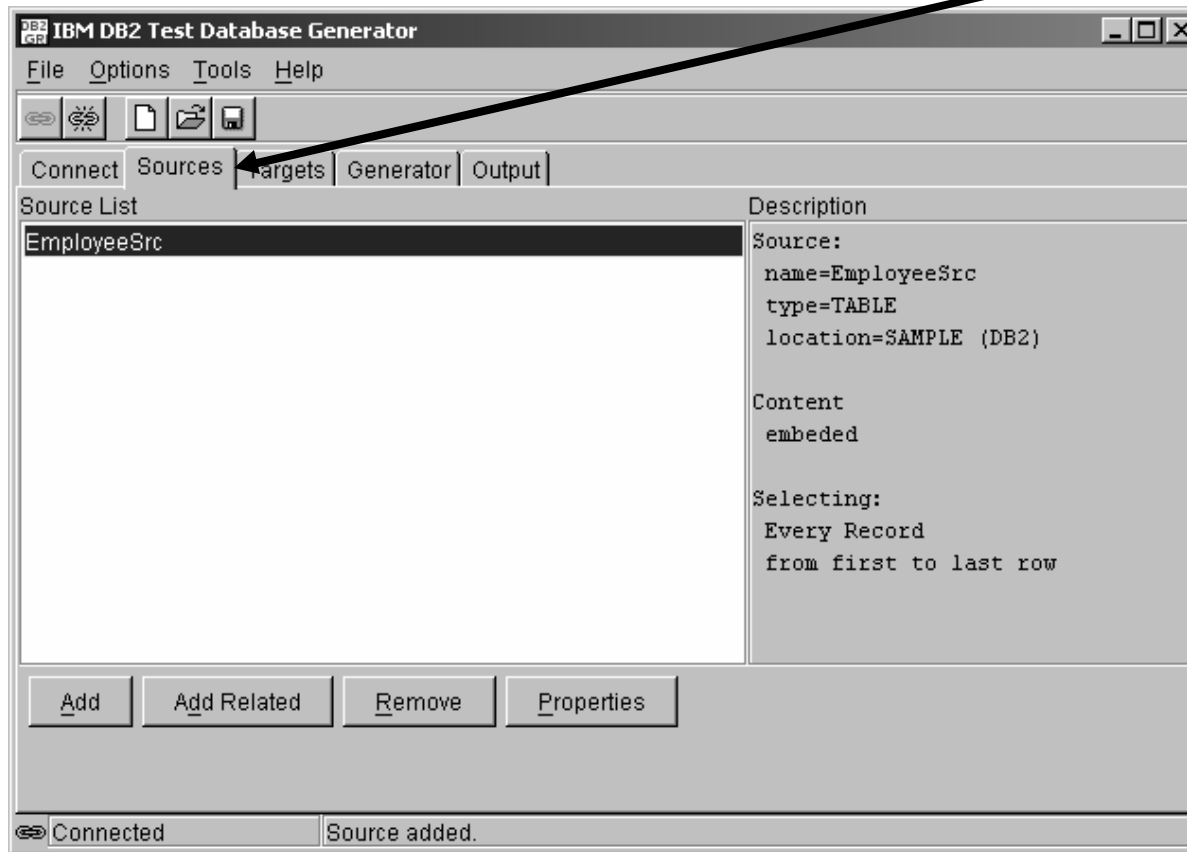
Use case : Specify source

- Here is the query that will be used to obtain the source data



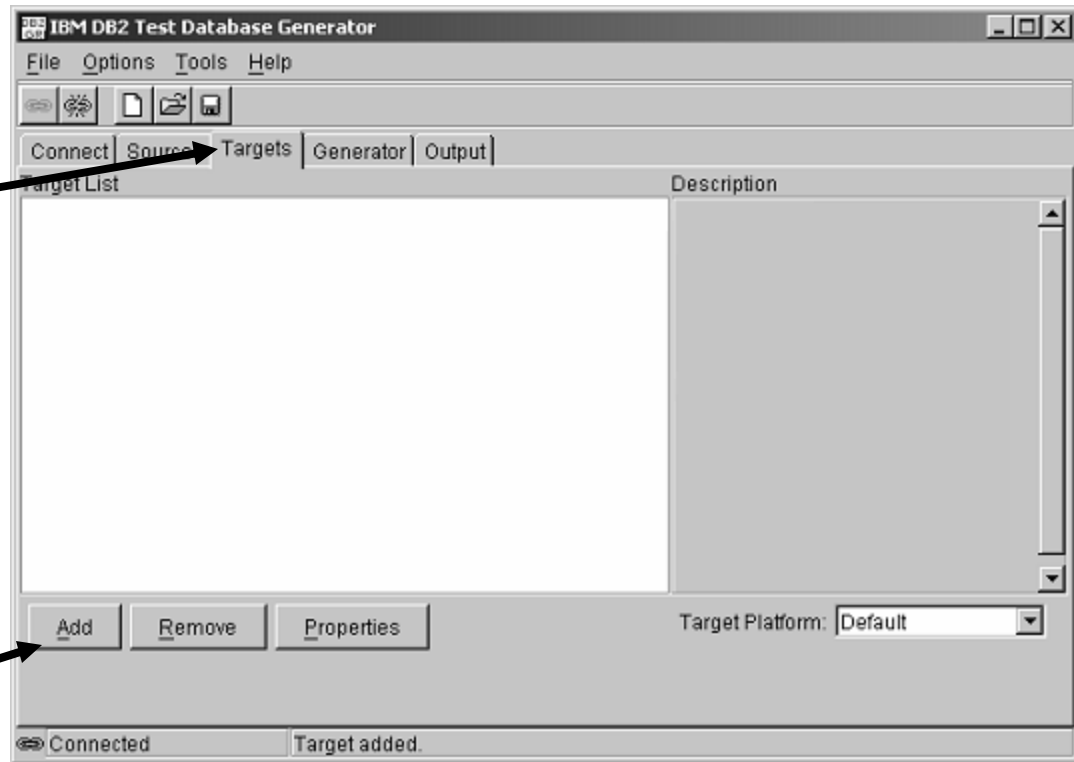
Use case : Specify source

- The source appears in the main window on the **Sources** tab



Use case: Specify targets

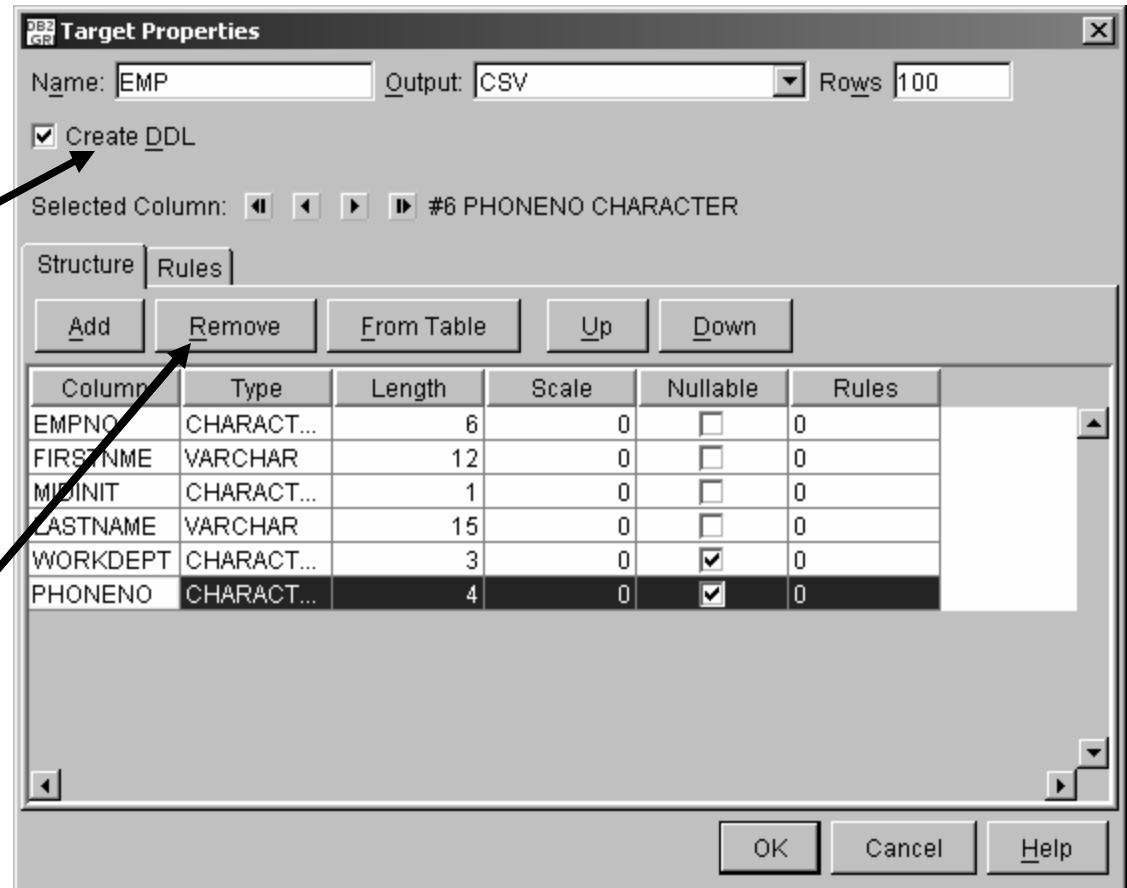
- Start with an empty Target List in the **Targets** tab



- Add first target using the Add button

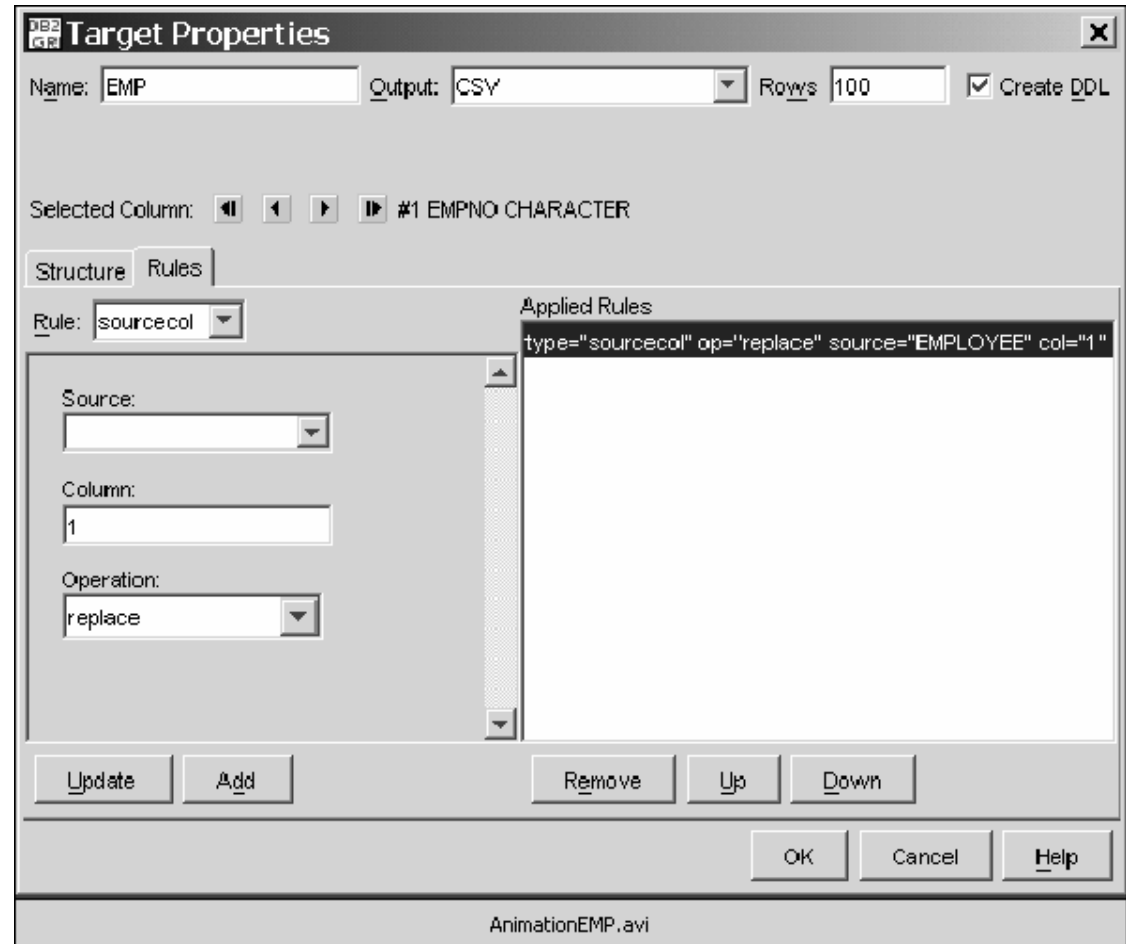
Use case: Specify target 1

- Add a target for the EMP table
- Select the Create DDL check-box so that DDL will be created for this target
- Remove unwanted columns from the source definition



Use case: Specify rules for columns target 1

- SourceCol rule applied to each column
- PHONE is not obtained from source
- Instead, PHONE is generated using a combination of pattern, static, and random rules



Use case: Specify target 2

- Add a target for the EMP_INFO table
- Select the “Create DDL” check-box so that DDL will be created for this target
- Remove unwanted columns from the source definition

Target Properties

Name: EMP_INFO Output: CSV Rows: 100

Create DDL

Selected Column: #2 FIRSTNME VARCHAR

Structure Rules

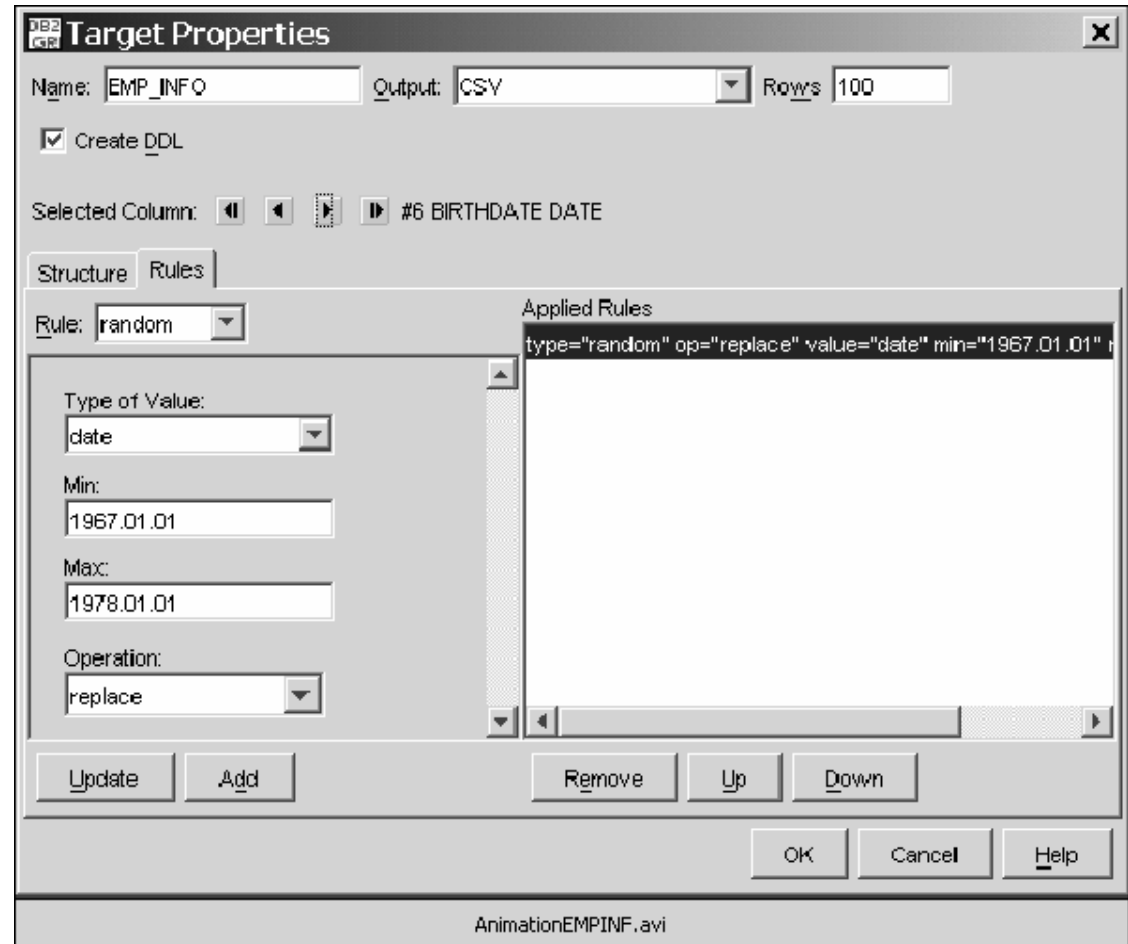
Add Remove From Table Up Down

Column	Type	Length	Scale	Nullable	Rules
EMPNO	CHARACT...	6	0	<input type="checkbox"/>	0
HIREDATE	DATE	10	0	<input checked="" type="checkbox"/>	0
JOB	CHARACT...	8	0	<input checked="" type="checkbox"/>	0
EDLEVEL	SMALLINT	5	0	<input type="checkbox"/>	0
SEX	CHARACT...	1	0	<input checked="" type="checkbox"/>	0
BIRTHDATE	DATE	10	0	<input checked="" type="checkbox"/>	0
SALARY	DECIMAL	9	2	<input checked="" type="checkbox"/>	0
BONUS	DECIMAL	9	2	<input checked="" type="checkbox"/>	0
COMM	DECIMAL	9	2	<input checked="" type="checkbox"/>	0

OK Cancel Help

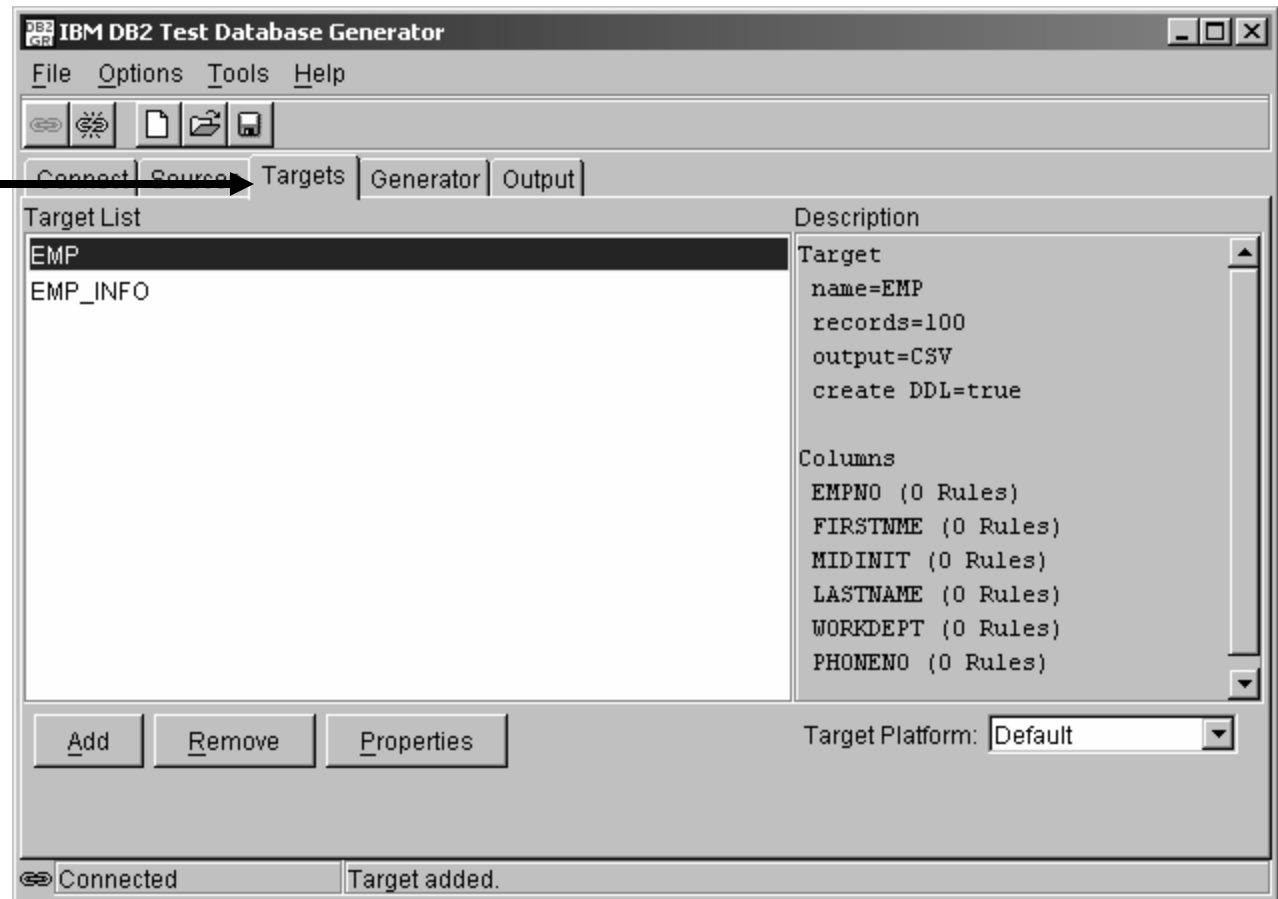
Use case: Specify rules for columns in target 2

- SourceCol rule applied to each column
- BIRTHDATE is not obtained from source
- Instead a random BIRTHDATE is generated in a range starting from January 1, 1967



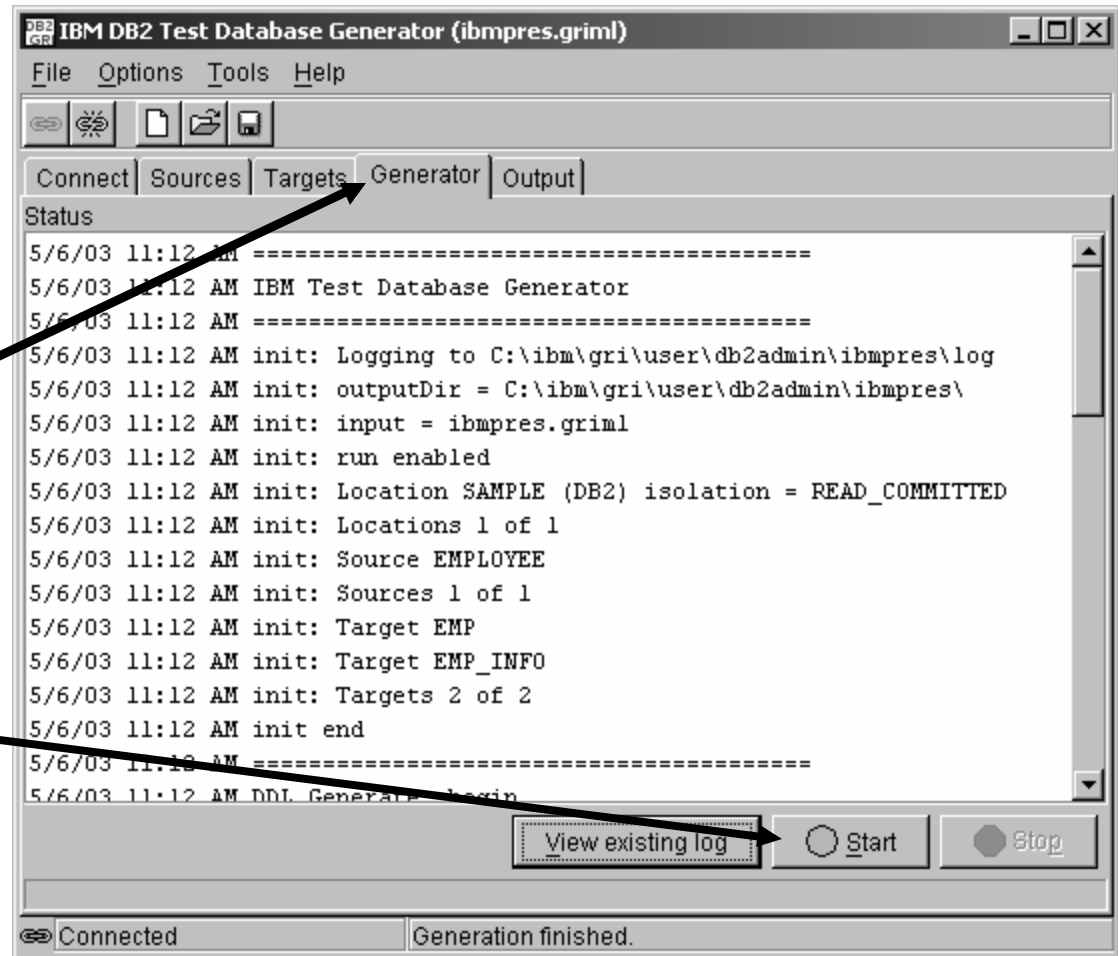
Use case: Specify targets

- The targets appear in the main window of the **Targets** tab
- An at-a-glance description appears in the right-hand frame
- You can save this data profile
- You can generate test data



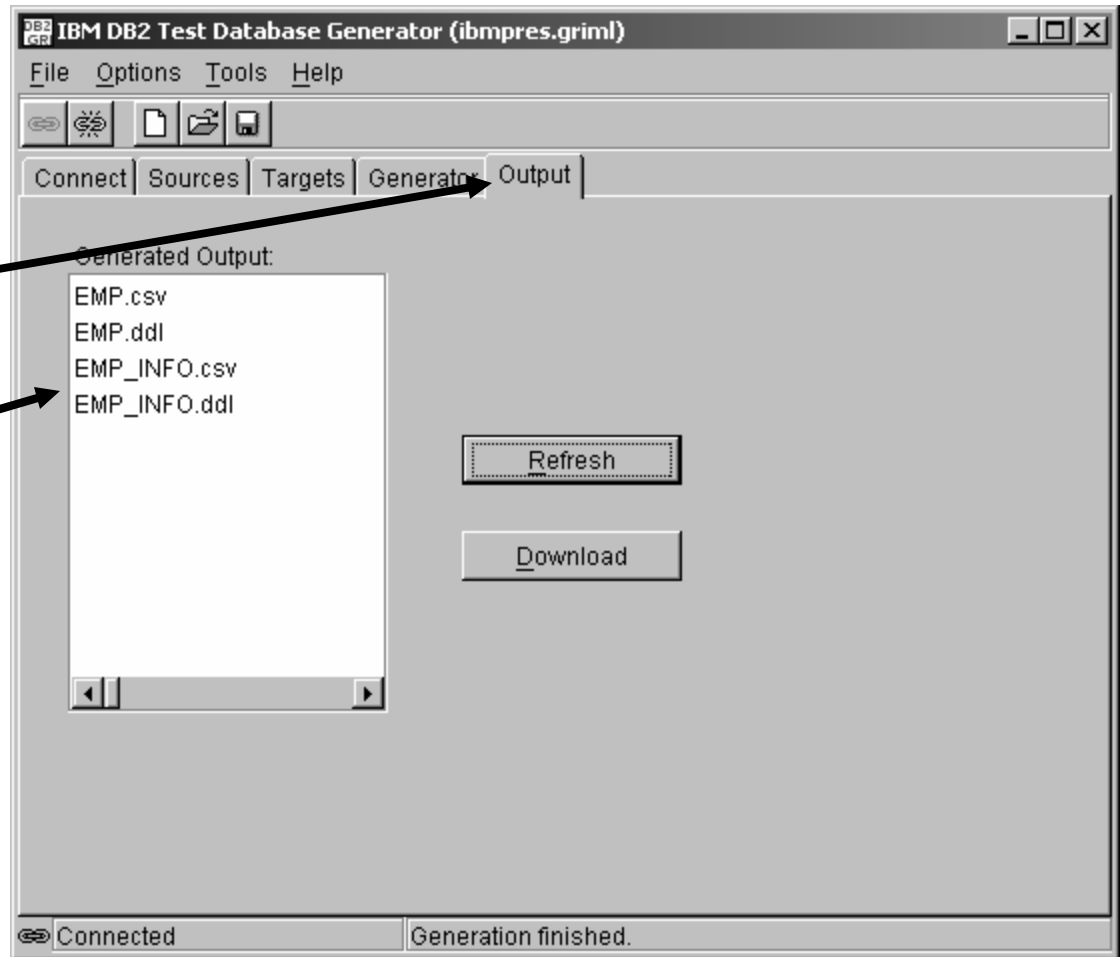
Use case: Generate test data

- Watch the log as your targets are created
- The log appears in the main window of the **Generator** tab
- Sessions can be started and stopped using the buttons at the bottom



Use case: View list of targets

- The generated target files appear in the main window of the **Output** tab
- DDL files and target data files



Use case: View generated test data

The screenshot displays the IBM DB2 Test Database Generator (ibmpres.griml) interface. The main window shows the 'Generated Output' section with a list of files: EMP.csv, EMP.ddl, EMP_INFO.csv, and EMP_INFO.ddl. A 'Refresh' button and a 'Download' button are visible below the list. The 'EMP.csv' file is selected, and its contents are displayed in a separate window titled 'C:\ibm\gri\user\db2admin\ibmpres\EMP.csv'. The data is presented as a list of rows, each containing a unique ID and various attributes including names, initials, and phone numbers. The status bar at the bottom indicates 'Connected' and 'Generation finished.'.

Generated Output:

- EMP.csv
- EMP.ddl
- EMP_INFO.csv
- EMP_INFO.ddl

Refresh

Download

Connected Generation finished.

C:\ibm\gri\user\db2admin\ibmpres\EMP.csv

```

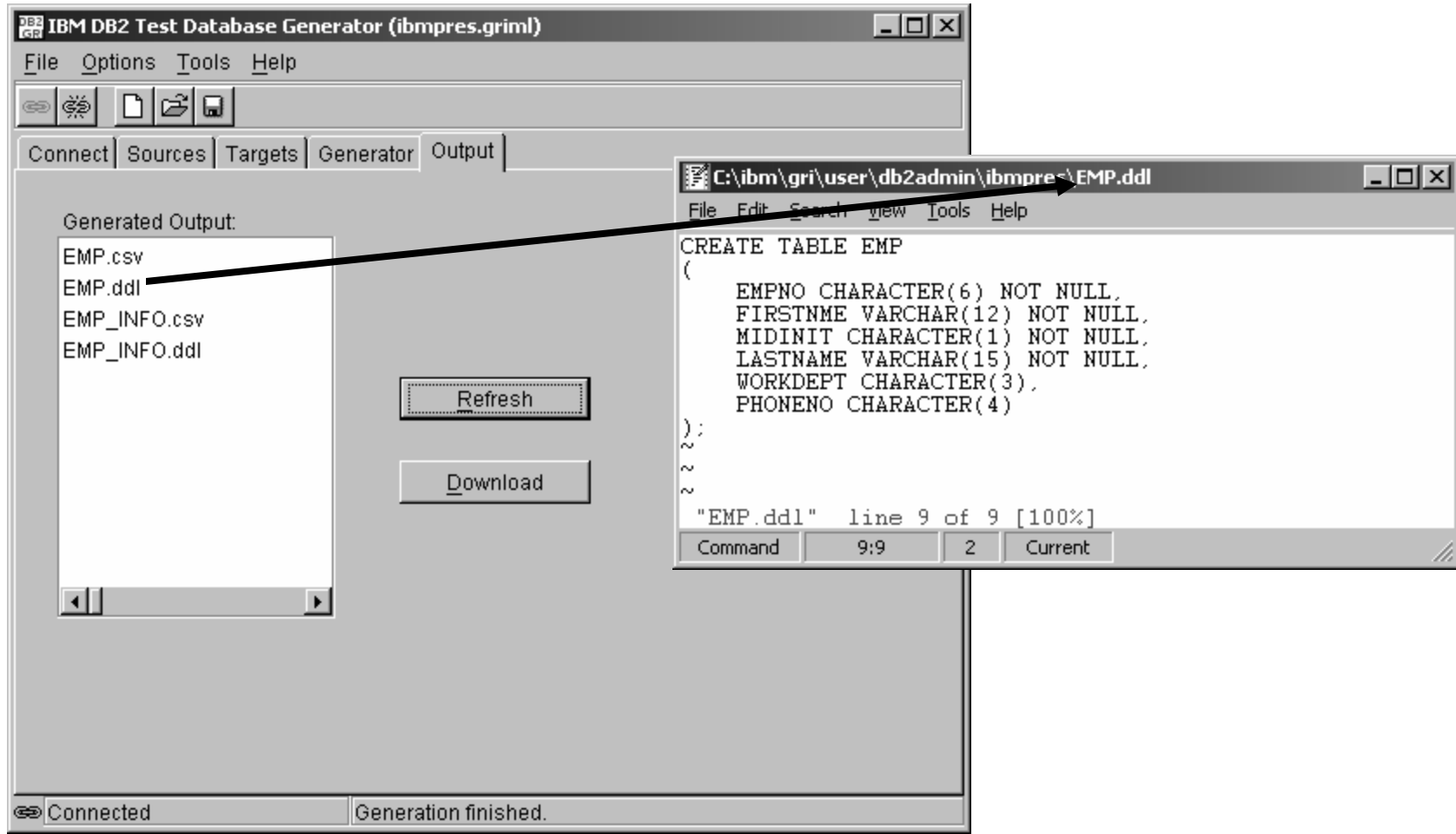
'000020','MICHAEL','L','THOMPSON','B01','946-432-9888'
'000030','SALLY','A','KWAN','C01','534-432-5419'
'000050','JOHN','B','GEYER','E01','764-657-3950'
'000060','IRVING','F','STERN','D11','922-432-6580'
'000070','EVA','D','PULASKI','D21','743-657-3790'
'000090','EILEEN','W','HENDERSON','E11','495-657-1047'
'000100','THEODORE','Q','SPENSER','E21','685-432-2819'
'000110','VINCENZO','G','LUCCHESI','A00','727-657-2158'
'000120','SEAN','O','CONNELL','A00','534-657-1850'
'000130','DOLORES','M','QUINTANA','C01','548-657-1530'
'000140','HEATHER','A','NICHOLLS','C01','743-432-6225'
'000150','BRUCE','ADAMSON','D11','923-233-7962'
'000160','ELIZABETH','R','PIANKA','D11','488-432-3336'
'000170','MASATOSHI','J','YOSHIMURA','D11','346-657-7492'
'000180','MARILYN','S','SCOUTTEN','D11','537-233-4768'
'000190','JAMES','H','WALKER','D11','529-432-7737'
'000200','DAVID','BROWN','D11','326-657-2293'
'000210','WILLIAM','T','JONES','D11','899-233-9823'
'000220','JENNIFER','K','LUTZ','D11','255-233-5407'
'000230','JAMES','J','JEFFERSON','D21','283-432-3239'
'000240','SALVATORE','M','MARINO','D21','858-657-6831'
'000250','DANIEL','S','SMITH','D21','478-432-5731'
'000260','SYBIL','P','JOHNSON','D21','648-432-2857'
'000270','MARIA','L','PEREZ','D21','522-432-7544'
'000280','ETHEL','R','SCHNEIDER','E11','842-432-7254'
'000290','JOHN','R','PARKER','E11','544-657-4743'
'000300','PHILIP','X','SMITH','E11','539-233-4597'
'000310','MAUDE','F','SETRIGHT','E11','432-432-3055'
'000320','RAMLAL','V','MEHTA','E21','822-657-6058'
'000330','WING','LEE','E21','989-657-3660'
'000340','JASON','R','GOUNOT','E21','886-233-3993'
'000010','CHRISTINE','I','HAAS','A00','494-657-2882'

```

"EMP.csv" 32 lines, 1709 chars

Command 6:32 7 Current

Use case: View generated DDL



The screenshot displays the IBM DB2 Test Database Generator (ibmpres.griml) interface. The main window shows a list of generated output files: EMP.csv, EMP.ddl, EMP_INFO.csv, and EMP_INFO.ddl. A black arrow points from the 'EMP.ddl' file in the list to a separate window titled 'C:\ibm\gri\user\db2admin\ibmpres\EMP.ddl'. This window displays the SQL DDL code for creating the EMP table:

```
CREATE TABLE EMP
(
  EMPNO CHARACTER(6) NOT NULL,
  FIRSTNME VARCHAR(12) NOT NULL,
  MIDINIT CHARACTER(1) NOT NULL,
  LASTNAME VARCHAR(15) NOT NULL,
  WORKDEPT CHARACTER(3),
  PHONENO CHARACTER(4)
);
~
~
~
"EMP.ddl" line 9 of 9 [100%]
```

The status bar at the bottom of the main window indicates 'Connected' and 'Generation finished.'

Use case: View generated test data

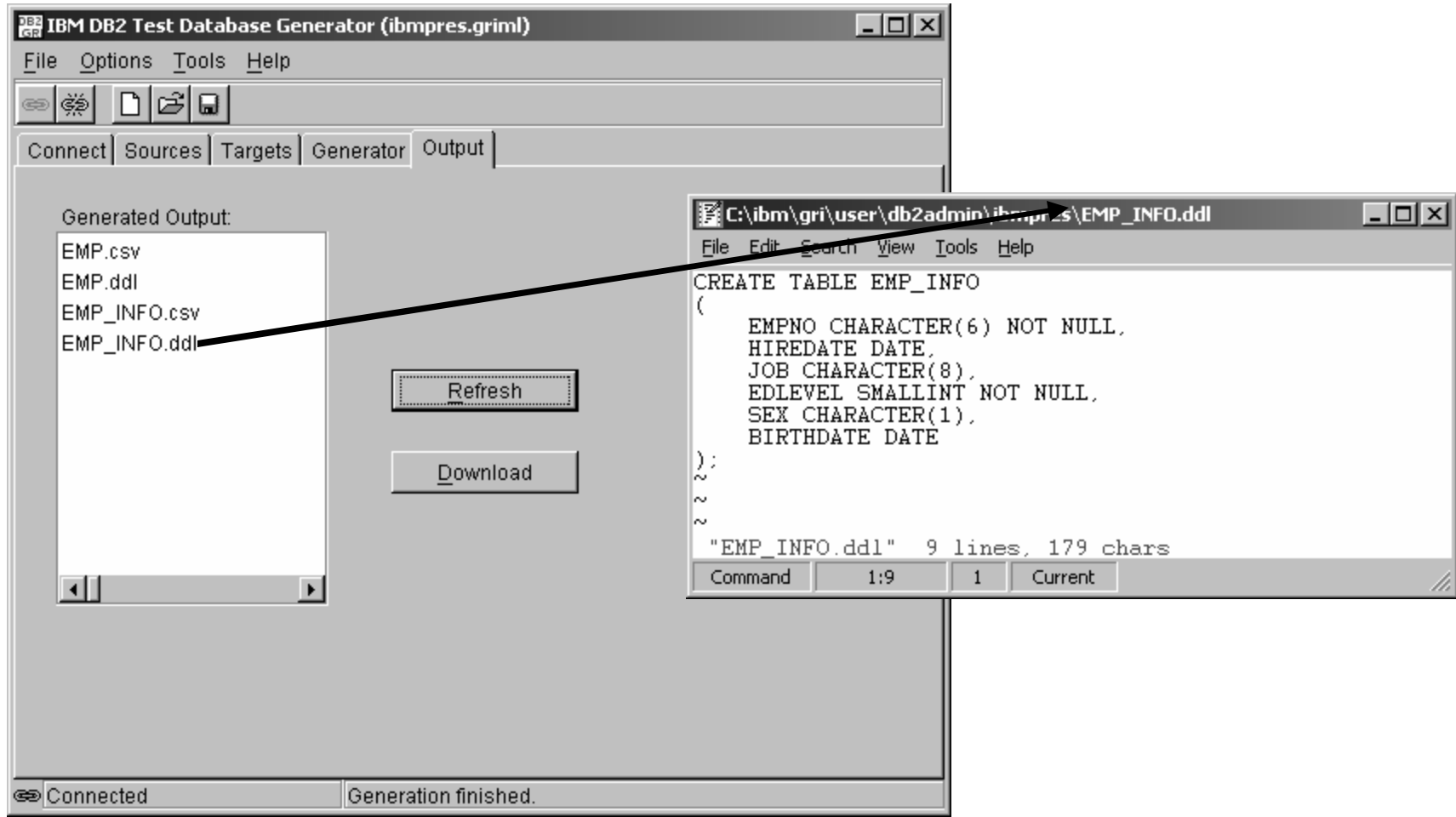
The screenshot displays the IBM DB2 Test Database Generator (ibmpres.griml) interface. The main window shows the 'Generated Output' section with a list of files: EMP.csv, EMP.ddl, EMP_INFO.csv, and EMP_INFO.ddl. A 'Refresh' button is visible below the list. The 'Output' tab is selected, and the status bar indicates 'Generation finished.' and 'Connected'.

An arrow points from the 'EMP_INFO.csv' file in the 'Generated Output' list to a preview window titled 'C:\ibm\gri\user\db2admin\ibmpres\EMP_INFO.csv'. The preview window shows the following data:

```
000020', '2003-11-10', 'F', 8, 'X', '1968-01-20'  
'000030', '2003-04-07', 'C', 3, 'X', '1978-01-01'  
'000050', '2002-05-28', 'D', 12, 'X', '1972-09-10'  
'000060', '2002-05-01', 'F', 8, 'X', '1971-05-17'  
'000070', '2001-07-09', 'D', 6, 'X', '1974-02-27'  
'000090', '2001-09-24', 'A', 12, 'X', '1974-06-16'  
'000100', '2002-06-30', 'D', 12, 'X', '1969-09-17'  
'000110', '2001-12-10', 'F', 12, 'X', '1968-12-30'  
'000120', '2001-12-14', 'A', 5, 'X', '1978-01-01'  
'000130', '2003-07-26', 'D', 12, 'X', '1976-11-06'  
'000140', '2001-05-13', 'F', 9, 'X', '1971-02-18'  
'000150', '2003-03-17', 'A', 9, 'X', '1977-02-27'  
'000160', '2002-04-09', 'C', 10, 'X', '1973-04-18'  
'000170', '2003-10-31', 'D', 5, 'X', '1977-04-30'  
'000180', '2001-08-13', 'D', 9, 'X', '1967-11-29'  
'000190', '2003-06-09', 'D', 11, 'X', '1978-01-01'  
'000200', '2003-11-11', 'A', 8, 'X', '1970-02-01'  
'000210', '2003-05-22', 'A', 12, 'X', '1976-12-21'  
'000220', '2001-11-06', 'F', 3, 'X', '1975-06-03'  
'000230', '2001-05-20', 'E', 12, 'X', '1977-11-22'  
'000240', '2001-11-04', 'D', 3, 'X', '1967-11-01'  
'000250', '2003-07-16', 'C', 4, 'X', '1968-04-06'  
'000260', '2002-07-21', 'C', 6, 'X', '1976-08-29'  
'000270', '2002-01-20', 'F', 9, 'X', '1972-09-03'  
'000280', '2003-04-12', 'A', 8, 'X', '1967-08-27'  
'000290', '2001-12-11', 'C', 11, 'X', '1976-02-03'  
'000300', '2002-03-25', 'B', 10, 'X', '1967-03-16'  
'000310', '2003-10-24', 'D', 10, 'X', '1972-03-01'  
'000320', '2002-12-12', 'E', 5, 'X', '1969-07-18'  
'000330', '2003-05-06', 'F', 12, 'X', '1975-05-05'  
'000340', '2003-04-09', 'B', 4, 'X', '1973-11-01'  
'000010', '2003-01-24', 'F', 8, 'X', '1976-10-11'
```

The status bar at the bottom of the preview window indicates: "EMP_INFO.csv" 32 lines, 1485 chars. The command line shows: Command 1:32 1 Current.

Use case: View generated DDL



The screenshot displays the IBM DB2 Test Database Generator (ibmpres.griml) interface. The main window shows the "Generated Output" list with the following files:

- EMP.csv
- EMP.ddl
- EMP_INFO.csv
- EMP_INFO.ddl

The "EMP_INFO.ddl" file is selected, and its content is displayed in a separate window titled "C:\ibm\gri\user\db2admin\ibmpres\EMP_INFO.ddl". The DDL content is as follows:

```
CREATE TABLE EMP_INFO
(
  EMPNO CHARACTER(6) NOT NULL,
  HIREDATE DATE,
  JOB CHARACTER(8),
  EDLEVEL SMALLINT NOT NULL,
  SEX CHARACTER(1),
  BIRTHDATE DATE
);
~
~
~
"EMP_INFO.ddl" 9 lines, 179 chars
```

The status bar at the bottom of the main window indicates "Connected" and "Generation finished.".

END

