

## \* Software Installation \*

- \* Go to control panel and uninstall java from your machine.
- \* Install jdk 7.
- \* unzip the software TOS-BD-20141207-1530-V5.6.1.zip
- \* Install oracle.
  - a) unzip the oracle XE112-win64.zip
  - b) click on setup to install oracle.
- \* While installing oracle  
you have to set :pwd:manager  
: re enter the pwd:manager
- \* Go to command prompt and start typing. Sql plus
  - user id : system
  - pwd : manager
- \* unzip the sql developer - 17.2.0.188.1159 - x64.zip
  - Go to unzipped folder and click on the sql developer icon.
  - Create a new connection by clicking on + button.
    - user name : system
    - pwd : manager
    - host : local host
    - port : 1521
    - sid : xe

### Notes:-

first time when you open Talend  
make sure your system is connect to the  
Internet.

\* goto unzipped folder of TOS-BD-20141207-1530-V5

click on 3 icons (.exe files)

first time when you open talend

→ give project name : AMEX

→ Select project click on open button

→ when you open talend make sure popups will come and you clicked on accept all and install option

\* Install notepad ++.

30/07/19

→ Data base

Creating oracle db connection :-

① → Go <sup>to</sup> repository

② → Select : Meta data

→ click on small arrow (→)

→ And choose Db connection.

→ Right click on select "create connection  
(oracle - local host)

→ Give some name (oracle) \* click on next button

Now choose Db type as oracle with service name

And fill the form.

Login : system

Pwd : manager

Service : local host

Port : 1521

Service name : xe.

→ click on check button you should see connection

as successfull. click on finish button.  
→ you can see the connection which you created  
in Repository

↳ Media data

↳ Db connection

↳ connection name  
(oracle-local-host)

Scanning the file

\*\* Scanning the dev Metadata of Delimited of  
CSV file (comma separated values file):

→ Go to reporting \* metadata

\* Rig file delimited

↳ right click.

→ any say Create file delimited.

→ Give some Name mostly file Name.

→ Click on Next Button.

→ Browse for file using Brower Button

By default it shows only csv files.

→ you can change the filter form.

• csv to \*.\* and choose the file.

→ Click on Next Button and fill the form.

→ By choosing proper field separator.

→ Proper row separator.

→ If the file has heading. choose set  heading row has column names option.

\* click on refresh Review.

You should see the data in tabular format

- Click on Next Button the Talend guess the Metadata as per the source data.
- you can change <sup>the</sup> metadata as per the Business aspect.

Click on finish Button).

→ file with pipe delimited and No header :-

→ follow the same as while choosing Separated choose "custom ANSI".

- And change the delimited as per the source.
- If the file has no header, make sure you are not selected set heading row as Column Names option.

31-07-19

\*\* loading Data from file to Table :-

- Go to repository Job design.
- right click and say create job.
- Give some name like job - file - to - table
- Click on finish Button.
- Now you can see a talend.

↳ Go to metadata

↳ delimited section.

- Select the file which we can scan earlier (Metadata) drag and drop to the job.

And select T file input Delimited.  
→ Drag and drop the oracle connection from meta data tree to the job.

→ Select T-oracle output.

→ Select the input icon  
↳ right click  
↳ select row, select main.  
↳ And click on the O/P component.

→ Now double click on the T-oracle O/P  
give the table name Double coats.

→ And select action on table as create  
table if doesn't exist.

→ Run the job using RunTab (or) F6.

### Retrieving table Data :-

→ Go to Repository

↳ select oracle connection

↳ right click

↳ And select day Retrive schema

↳ click on Next Button.

↳ select the user from the list it.

→ Click on small arrow beside the system.

and select the required tables by selecting check box.

→ Click on Next button, click on finish.

→ you can see the table metadata in the Oracle connection.

→ By clicking on small arrow.

\* Table to ~~third~~ file

\* Table to file :-

→ Drag and drop the table & meta data from the scanned tables to the job.

→ On the select toolcal input

→ Now go to pallet and search for o/p delimiter

→ Drag and drop the component to the job.

→ Click on the designer and start typing delimited.

→ And choose T-file o/p delimited.

Connect your i/p throw the o/p using now double click on T-file o/p delimited.

→ Choose the file path where the file need to be created by using Browse Button.

→ Run the job.

- By default talend create the file, without header, you can choose include header option and you can choose the delimiters for row & column.
- As per the requirement.
- for each and every run, talend will overrides the target file.
- if you want to append the data to the bottom of the file, choose append option in the o/p delimited.

01/08/19

#### \* SQL Query:-

- Select,
- empno,
- sal,
- dept name,
- loc,
- From system.emp - 08/07/2019 where deptno = 10 order by sal desc".

if you want remove duplicate.

Query : \* select distinct

empno

Sal

emp name -----

from system.emp - 08/07/2019 where deptno = 10 order by sal desc"

## \*\* T-oracle input :-

By default T-oracle I/P comes with select query in \* double coats.

You can change the query as per the requirement like adding where clause.

order by clauses.

distinct clause. ---

→ When you modify the query, make sure you hit the query schema.

## \*\* Applying filter on + file input delimited :-

→ Take file as input

→ connect into + filter row component

→ Double click on it.

click on + button to add a column.

Input column	function	operator	value
dept no	empty	equals	10
Sal	empty	Greater (or) equal to	15000

## \*\* T-Sort Row :-

→ Sort the records based on some key column either ascending (or) descending.

→ We go with t sort, Row.

→ click on + button to add a column. ← → options in the + sort row

→ Select type

Schema column	sort num or alpha	order asc (or) des.
Sal	Num	asc

\* Removing duplicates:-

→ we go with + uniq row.

→ options w/ in the + uniq row.

→ select all the columns as  key attribute.

column	key attribute	case sensitive
empno	<input checked="" type="checkbox"/>	<input type="checkbox"/>
empname	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Sal	<input checked="" type="checkbox"/>	<input type="checkbox"/>
HRA	<input checked="" type="checkbox"/>	<input type="checkbox"/>
dept NO	<input checked="" type="checkbox"/>	<input type="checkbox"/>
dept Name	<input checked="" type="checkbox"/>	<input type="checkbox"/>
loc	<input checked="" type="checkbox"/>	<input type="checkbox"/>

05/08/19

\* \* t-filter row → t-row row

↳ t-file output delimited - 1

Rejected.

t-unique → oracle

↳ t-file output delimited - 1

65/08/19 \* t-filter column → To remove unnecessary column and select only required columns from the flow we can go with t-filter column.

→ Go to edit schema & copy only required columns from left to right.

### \* t-replicated :-

→ To copy the same data into multiple o/p go with t-replicate.

→ No options available t-replicate.

→ you can take n-number of from t-replicate

### \* t-Sample Row :-

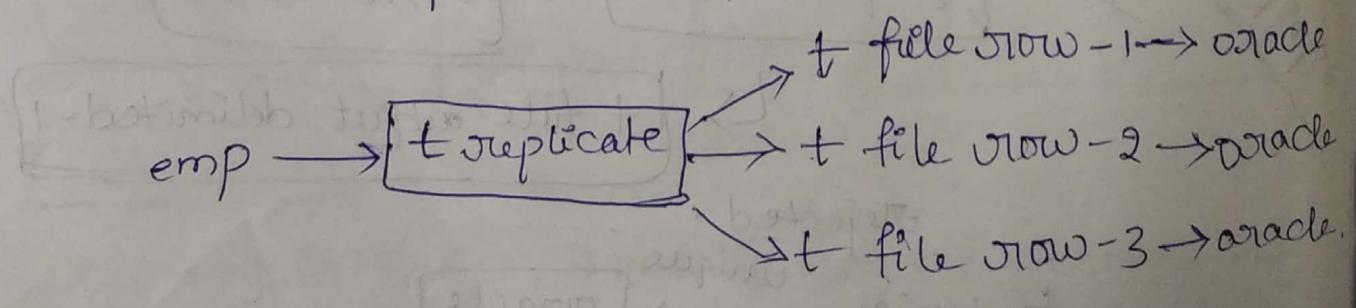
→ you can specify the rows.

example : 1.. 100 (1 to 100)

101.. 200 (100 to 200)

1.5.7 (1 and 5 and 7)

### \* Multiple → t Replicate



## header, footer and limit :-

\* header :- In i/p component if you specify the header as "n".

→ it skips first 'n' no.of records from the file while reading.

\* footer :- it skips "n" no.of records from the bottom of the file.

\* limit : once after start the reading file it stops at nth line.

\* +unite :- works like union all operation in database.

→ if you have same metadata flow, if you want to combined them as a single o/p, you can go with +unite.

options in the +unite :

click on sync columns.

+unite will not remove duplicates.

to get union functionality use +unique

row along with +unite.

## \*\*\* Multithread execution:-

→ To run the sub jobs in the parallel we need enable multithread execution <sup>option</sup> available at job tab.

## \*\* Pre job and Post job:-

→ Some sub jobs you want run at first ~~at~~ <sup>100%</sup>

Can go with "Pre job" on component OK.

if you some sub jobs → if you want to Run at last : "Post job" on component

OK.

07/08/19

## \*\* t-file list:

→ folder containing multiple files with same metadata. need to load into a single target.

Metadata

The const count of the files will keep on changes → The count of files in the i/p folder will keep on changes.

### Step:- 1

→ Scan the metadata of single from the file & a folder.

→ Drag and drop job. (make it as i/p)

### Step:- 2

→ connect it to target.

- Now take the component called t-file list
- Double click on ~~the~~ table it.
- and choose the directory ~~for~~ the files from where you to pick the files.
- click on "(+)" button to a file.

example: `*.csv`  
`*.txt`  
`*.*`  
`emp*.csv`

- connect t-file list to t-file i/p delimited using now iterate link.
- Double click on i/p delimited component.  
change the property type from Repository to [Built-in]
- remove the file location use ~~ctrl~~ and ~~space~~ button to the choose global variable variable called t-file list current file path.
- Run the job.  
if you are using file as an op, make sure you select "append" option.  
if your target is table don't choose ~~to~~ options drop and downcate.

## \* Reading files from subdirectories:-

↳ use an option called include sub directories in the t-file list component.

## \* Improving the performance:-

↳ Select the iterate link

↳ Right click

↳ Go to setting

↳ Select enable parallel

option. and give some Number (4/5).

08/08/19

## \* Normalize:-

I/P :-

= empno, ename, Skills

100, Mahesh, informatica; informatica; datastage; oracle

104, mani, informa; tale; dataart; ora

101, Kalyan, informatica; c;c++

O/P :-

100 Mahesh informatica

100 Mahesh informatica

100 Mahesh datastage

100 Mahesh oracle

104 mani mani - - -

### Step 1:-

- ↳ Scan the metadata of the file.
- ↳ Drag and drop the job. → input
- ↳ Connect it to t-normalize.
- ↳ Double click on
- ↳ Select the normalize column.

Ex:- Skills.

- ↳ And select the delimited column.

Ex:-

- ↳ Connect into target.

- ↳ By default remove duplicates if you want  
remove duplicate use t-unique row after  
normalize.

- ↳ Else we have an option called get rid of  
duplicate rows from o/p, in t-normalized  
component.

### \* Denormalize:-

Input :- empno, ename, skills

100, Mahesh, Informatica

100, Mahesh, talend

100, Mahesh, oracle

101, Kalyan, A

101, Kalyan, B

101, Kalyan, C

101, Kalyan, D

Output :-

empno, ename, skills

100 Mahesh, Informatica; talend; oracle

101 Kalyan A;B;C;D

↳ Scan the metadata file make it as i/P.

↳ Connect it to denormalize.

↳ Double click on denormalize.

↳ Go to edit Schema.

↳ copy the columns from left to right.

↳ using double arrow.

↳ And increase the denormalize length.

↳ Click on 'OK'.

↳ Click on (+) Button.

↳ select the column to do denormalize.

Ex:- Skills.

↳ select the delimited which you want in target.

Ex:- (;

↳ Run the job.

↳ By default if you want remove use t-mq Row.

↳ Before the denormalize (or) use an option called  merge same value in t-denormalize.

### \* Buffer :-

Buffer is a component will help you to keep the data in Memory for some memory purpose and reuse it Multiple times in the same job using Buffer.

- input and Buffer output.

### \* Hash :-

Hash is a Memory component which keeps the data in the memory for some memory purpose and reuse it Multiple times. using hash i can store different Metadata into the Memory.

data

### \* Scenarios :-

Folder containing Multiple files with same Metadata need to load in a single target using sorting operation.

13/8/19

\* t Map :-

- Add output table option (Right hand side - top most corner (+) button).
- Drag and drop columns from left.
- Remove unwanted columns from right hand-side using edit schema.
- Remove unwanted output tables using (-) button on the right hand side topmost corner.
- Activate, deactivate expression filter.

Ex:- dept No = 10.  
and f f.

\* Filter condition in the t-Map on string columns

row.1.empname.equals ("mahesh")

Addition operation +

row.1.empname + row.1.dept

Full name

row.1.sal + row.1.HRA

Full Salary.

for space :- Row.1.empname + " " + Row.1.deptname  
→ Full Name.

to change date column → ctrl space.

⇒ Adding extra columns in the target.

→ column filter.

→ (header) Row filter.

→ (fff) Multiple filter condition.

[By using expression Builder]

→ Common condition

(left side added)

[Ex: salary - - -]

Row.1.sal > 20k.

→ filter condition in string column.

→ Added extra columns [Ex: full Name]

if we want space Row.1.emp.name + " " +  
Row.1.dept.

→ data

tallend date

get current date.

"dd-MM-YYYY". [ctrl space] to change date.

Ex:- EEE dd-MMM-YYYY HH:mm--

(→ pd output)

~~8 14/8/19~~ \* if condition :-

Ex:- if we want more than 15,000 sal - Grade A  
- Grade B

Add a column.

→ Grade

row.1. Sal >= 15,000? "Grade A": "Grade B":

condition ? true : false

\* multiple if condition:

row.1. Sal < 15,000 && row1. Sal >= 5000? "Grade B":  
"Grade C".

row1. Sal < 15,000 && row1. Sal >= 5000? "Grade B": "Grade C".

Emp. name:

All need in upper case function

String handling only using on strings.

String handling → upper case.

\* Alpha: Checks whether data is sorted in Alphabetical order not.  
[True / False → Boolean]

\* BTRIM:- Remove extra space from the Back side

\* CHANGES: to replace value in string.

("row.1. emp.name, ". ", "-")

↓

(• replace by →)

to remove replace with empty value. " "  
count:- count the no. of occurrences.  $\rightarrow$  repeated values (integer.)

REPLACE:- Replace value

FTRIM:- Remove extra space from front side.

INDEX:- Position. (Count start from 0). integer.

(\* the hollow world i. " " hollow")  $\rightarrow$  0:  
it counts with "0".

if doesn't found you will get "-1".

IS-ALPHA:- True/False (Boolean)

To check value is in alphabetics or not.

(Order doesn't matter).

\* LEFT:- How many characters we want from left side.

$\rightarrow [$  " hello world 1.7 ]

\* LENGTH:- Mokno. integer.

\* RIGHT:- How many char. from right side

[ " hello world 1.8 " ]

\* SPACE:-

row1. dept.name.

row1. emp name + string

\* STR:- (" ", 6)

\* TRIM:- front and last

\* Upper Case:

SPACE:

row1.deptnam

row1.empname + string

STR: ("row1", 6)

front & last

TRIM: front and last

upper case;

upper case:

22/08/19

Scenario :- 1:-

Mobile Numbers

Robby

String Handling. LEN(row1.mobilenum) = 12 ? "Yes"

Scenario :- 2 mail id - count.

22/08/19

Scenario :- 3

~~Scenario~~ \* Add extra column in target call

ename - valid - flag.

columns  
column

\* ename contains only Alphabits say "true"

Else "false". "No".

\* String Handling. LEFT (row1.mobilenum, 2). equals ("

\* Scenario :- 1 22/08/19

① Mobile Numbers.

String handling. LEN (row1. mobilenum) == 12 ? "Yes": "No"

\* Scenario 2 :- { Tlog Row  
(Table View)

Mail.id - count

\* Scenario 3 :- Add an extra column in the target  
called employee Name - valid - flag.

→ if the employee Name contains only Alphabets.  
say "true" else "No".

    → Mobile Nums.

\* String handling. LEFT (row1. Mobile num, 2). equals  
("91").

→ Need 4 output files based on the Mobile  
Number country code.

→ if No. starts with '91' it should go to  
India.txt.

→ If it starts with '01' it should go us.txt  
01 → USA

02 → Canada.

\* Add an extra column in the target called  
"country" if mobile No. starts with '91' it  
should be india.

01 → USA

02 → Canada

04 → UK

- \* valid email id's in one output.
  - \* invalid email id's into another O/P.
  - \* valid mobile nos in one O/P.
  - \* Invalid mobile nos in another O/P.
  - \* Valid emp. Names → one output.
  - \* Invalid emp. Names → another O/P.
- Need 2 extra columns in the target.  
 → called Mail Name and domain Name.  
 → So split the Mail Id into 2 parts:  
 Before @ → one output → Mail-name  
 after @ → another output → domain-name
- abaaa@gmail.com | ; abaaa; | @gmail.com

### \* Padding: [Expression]

- In my target the length of the emp.no. is fixed to "5". If IP column length is < 5, then pad "0" to the employee number to the left.
  - In the same way, employee name length is fixed to 15.
  - If the length is not 15, add space as padding character.
- 5-stringHandling. LEN(crow1; end).
- String Handling. STR('0'; 5-stringHandling. LEN(crow1 + row1; end))

- \* name space [in the place of 0]
- \* StringHandling · STR (" ", 15 - StringHandling · LEN)

23/08/19

### Date files:-

#### Compare date function:

- \* Talend Date · compareDate (row1 · Shipment · Start - date, row1 · Shipment · end - date)  $\leq 0$  · valid

#### Difference date:-

String  $\rightarrow$  Long

extra { No. of days  
Column }

Add date function:- some thing to  
date/month/year.

Proper  $\rightarrow$  <= 0  
improper  $\rightarrow$   $> 0 / = 1$  in  
used in 2nd func

String  
Days  $\rightarrow$  Integer  
only 10 digits  
Long.

compare Date :- Compare 2 dates and gives you  
which one is bigger and which one is smaller by giving 0, +1, -1.

Diff date :- Difference b/w 2 dates 0.

format Date :- converts from date to string.

Get current date [in string] :- get date.

Get first day of month :- first date.

Get last day of month :- last date.

Get part of date:

TalendDate · getPart of Date ("WEEK\_OF\_YEAR", ) TalendDate,  
get

6/08/19

get part of Date:-

get Random Date:

\* Is Date:-

Scenario:-

Date pattern:

date-id	date-pattern	month	days of year	days month
1/1/2010				
2/2/2010				
3/3/2010				
4/4/2010				
5/5/2010				
	Week of month	Week of year	dayname	holiday

Bank holiday.

→ java.lang.NullPointerException.

when we do operations on a null column

Value - Null pointer exception

How to handle null:

→ Relational.ISNULL (row1.Sal) ? 100 : row1.Sal + 100

→ Relational.ISNULL (row1.HRA) ? null : row1.HRA - 50

→ get date to date ("MSEEKDE AYER", "TOMORROW")

def

Numeric:

Numeric



Sequence

} S1

-2 min.

Mathematical:

↳ ABS [Always returns positive value].

MOD :- [Always returns remainder value].

\* Scenario 2 :-

Take 10 records file.

1, 3, 5, 7, 9 → 1 file.

2, 4, 6, 8, 10 → 1 file.

Key : seq + Mod

Scenario 3 :-

Take 10 records file.

I/P → 100

O/P :-

Sub	100	101	102	103
100				
101				
102				
103				
104	01			
105	02			
106	03			
107	04			
108	05			
109	06			
110	07			

Sub	100	101	102	103
100	100	0001	A	001
101	021	0002	B	101
102	001	00001	C	601
103	120	0021	D	201
104	041	00000	E	001
105	021	00001	F	201
106	001	00000	G	601
107	121	00101	H	001
108	001	00000	I	201
109	100	00000	J	601

Scenario 4 :-

first 4 records }  
last 4 records }  
file

## Joins :-

↳ 5 types.

- ↳ \* Left / Left Outer join
- \* Right Outer join
- \* Inner / Equal join
- \* Full outer join
- \* Cross (X) join.

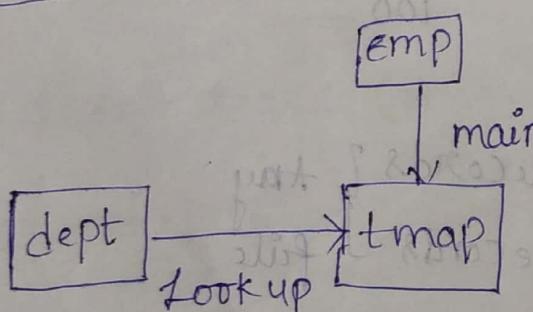
### \* Left Outer join :-

- All the records from left table will go to target.
- If there is a match, you will be getting a matched record.
- If there are multiple matches you will get multiple records.
- If there is no match, still the record goes to target with joint columns target as Null.

eno	ename	sal	HRA	DNO.
100	A	1000	100	10
101	B	1800	180	20
102	C	1600	160	10
103	D	1500	150	40
104	E	1400	140	80
105	F	1200	120	90

Dept (R)

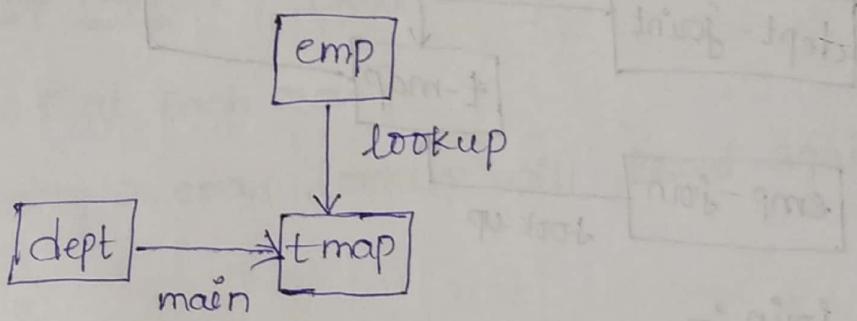
D-Dno	Dname	loc
10	HRA	Hyd
20	MKT	Che
40	Fin	Blr
20	Lab	Pune
20	Prod	Del
60	R&D	Noida



unique match - all matches.

## \* Right / Right Outer join:

- All Records from right table will go to target.
- if there is a match, you will get matched records.
- if there are multiple matches, you will get multiple records.
- if there is no matches still the grievance will go to target with joint column value as well.



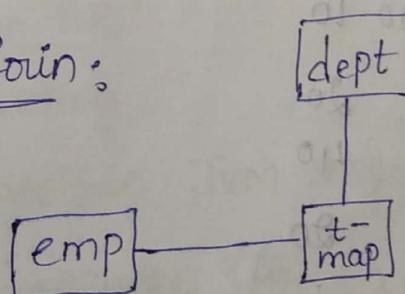
## Inner / Combi join:

- Only matched Records from Both the table.
- if there is a match, matched - Record.
- if there is multiple matches, multiple Records.
- if there is no match, we won't consider in the target.

lett [inner]

unique [All matches.]

## Full Outer join:



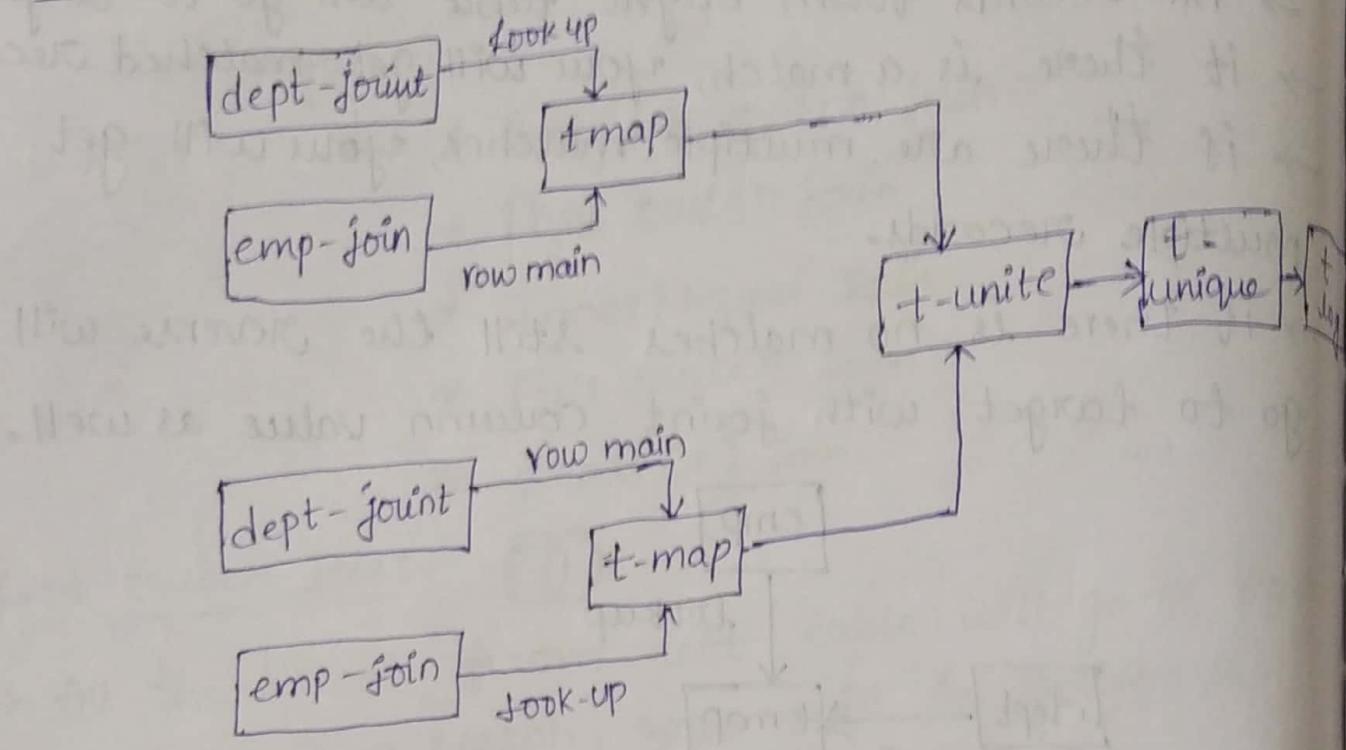
1st Left Outer  
2nd Right Outer

Next

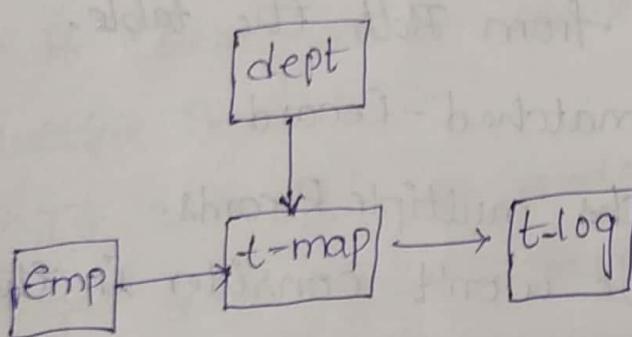
+ unique ↴

+ unique.

\*\* Full outer join:



\*\*\* cross - join :-



\* No Key column

\* Copy all i/p records to o/p.

100	A	1000	100	10	10
100				20	
:					
100					
101	B	1800	180	20	10
101				20	
:					
101				40	
:					
101				20	
:					

Match model :- first match

↓  
Only first

unique match → last record

Default : left outer joint } latest record will taken.  
: unique match }

Look up model : [load once]

[reload at each row]

for every record lookup will load again and again.

\* Reload at each row (cache) :-

\* Keep memory issue / Memory out of Bound exception:

Advance settings.

use specific JVM arguments.

if we want 5GB

-XMS 5120M	max
-XMS 3072M	min.

\* Steps :

→ Go to Run tab.

→ Advance setting.

→  use specific JVM arguments.

And increase the  
 $\begin{cases} \text{xms (min)} \\ \text{xmx (max)} \end{cases}$

→ Store temp data [true]

\* Click on t map component  
temp data [temp folder].

→ then it stores data in hard disk.  
[BIN file].

\* Inner join records:

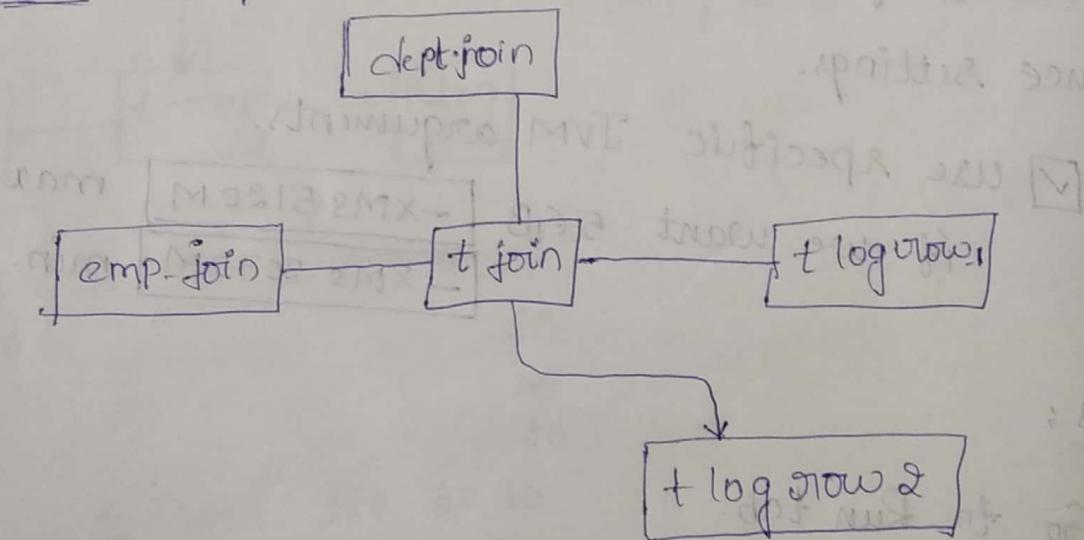
→ if we want to join inner join reject records.

→ Catch look up inner join rejects [true].

→ Catch o/p records [true] → reject records on output.

29/08/19

\* t-join component:



→ Allows only two sources get join.

→ By default left outer join unique match.

→ By subtyping right outer join unique match.

→ by selecting option called with reject o/p.

→ I can get inner join with unique.

and rejected records the main flow.

\* Steps :-

→ connect your sources to the t-Join

→ Double click on t-Join.

→ Go to schema

→ copy columns left to right from both main and look up.

→ Click on "OK".

→ Select on option called include look up columns in O/P.

→ Click on "+" button to add look up columns and O/P columns.

→ come to key definition section.

→ Click on "+" button to add a key definition

Scenarios :-

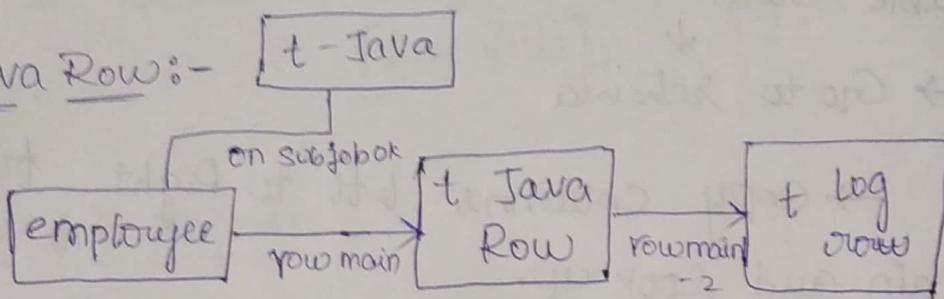
① Take some file <sup>need</sup> first 100, next 200 next 300 record target using [t-map]

② using t-Join need achieve first match

## t-Java :-

→ System.out.println("The job " + jobName + " started...  
and project name " + projectName);

## t-Java Row :-



→ Double Click



→ Edit Schema.

→ Here we have chance to add extra columns  
→ click on generator code.

→ Automatic script was written.

### t- Java :-

- \* t- Java is the component where you can write the java code
- \* which will get execute directly on the JVM.
- \* This component will not have edit Schema.
- \* So we can't connect to the follow.

### t- Java Row :-

- \* t- Java Row is the component. will have edit Schema.
- \* By you can add extra columns on the edit schema and we can write expressions in the o/p.

\* Same like t-map.

### t- Java flex :

- \* we have start code, main code and end code.
- \* The code which we are writting in the start code first.

\* The code which we are writting in the main code will ~~execute for~~ each and every row. execute at the last.

\* To reset the Numeric sequence we use the t-java component.

Numeric sequence (S1, 1, 1);

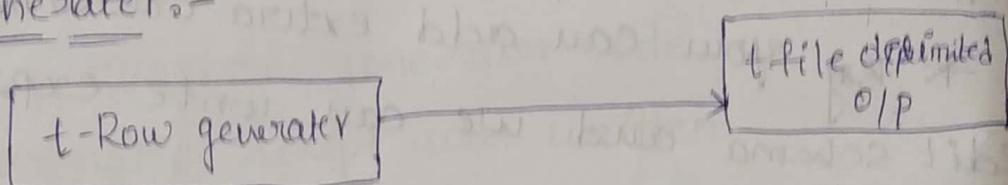
30/08/2019

\* Advanced mode t-filter row:-

input-row.deptno == 10 && input-row.sal > 15000 ||

input-row.deptno == 20 && input-row.sal > 15000.

\* t-Row generator:-



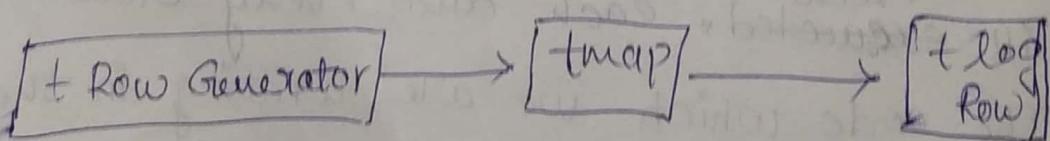
→ Double click on it.

→ Give the number of rows you want generate

→ Click on + button add required columns  
an corresponding data and corresponding  
function.

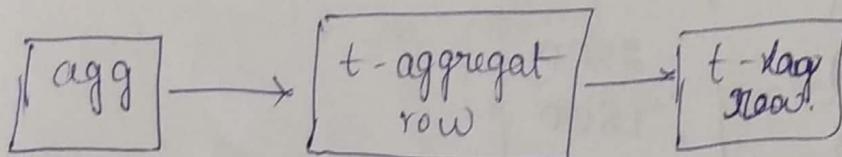
→ connect to target.

→ Run the job.



→ To push the Record from t-map you can we  
can use t-row.

t-aggregate row:-



\* Go to  
Advanced  
setting  
Change the  
column.  
[;]

Scenario 1:-

- ① Need semi column in the list column.

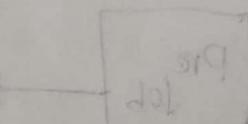
Scenario 2:-

I/P :-

	Col1	Col2	Col3	Col4	Col5
1	x	x	x	x	x
x	1	x	x	x	x
x	x	1	x	x	x
x	x	x	1	x	x
x	x	x	x	1	x
O/P :-	1	1	1	1	x

Scenario 3:-

e.no	ename	skill
100	maheish	informatica
100	maheish	Talend
100	maheish	oracle
101	mani	Talend
101	mani	Java
101	mani	MS SQL
101	mani	unix



e.no	ename	skill
100	maheish	i,t,o
100	mani	Ta,Ja,m,un,i,x

e.no	ename	Sal
100	A	1000
100	A	2500
101	B	1800
101	B	1300
100	A	6000
101	B	1400

%p:

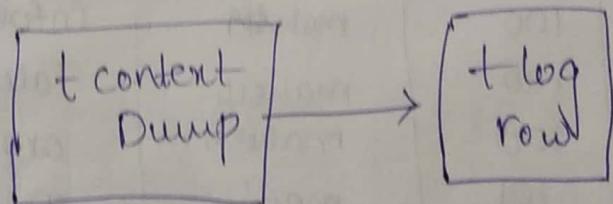
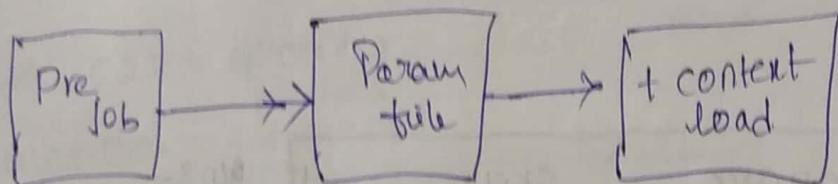
e.no	ename	Sal
100	A	6000
101	B	1800

03/09/19

### t- context load:-

- \* VAT → user acceptance testing
- \* QA → Testing cor) Quality Analysis

04/09/19



## Project Setting

↳ Job setting

↳ implicit context load.

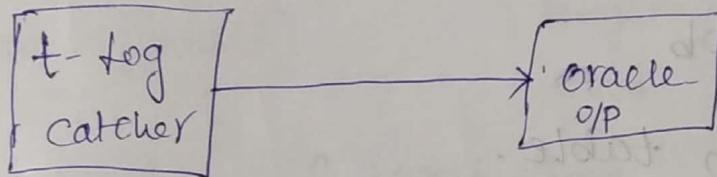
### \* Implicit context load :-

↳ from file

↓  
Browse

05/09/19

t-log catcher :- log catcher shows where errors occurred.



t-start catcher : [start & end job status]

columns :-

moment

Pid

father Pid

job :- log-start

2 job :- Run

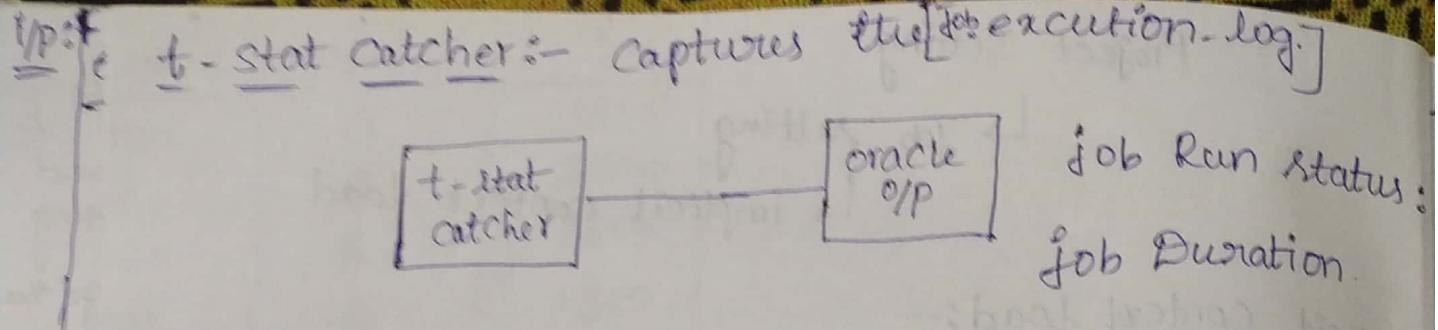
job [ ] → oracle

job [ ] → oracle

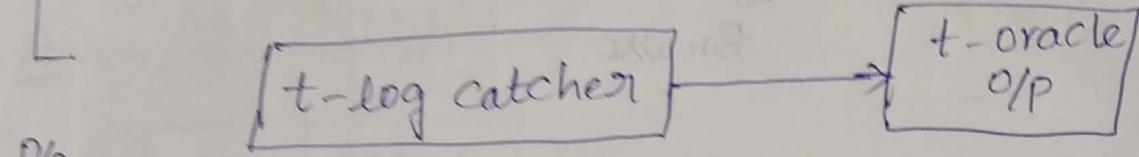
if user name / file name wrong

die on error

Job failed



- \* Where error exactly occurred: "exception".



- \* captured on table.
- \* Give table name
- \* Run Job

select \* from table.

- \* it shows errors.

log-catcher :- Captured the job failures, with moment when, reasons.

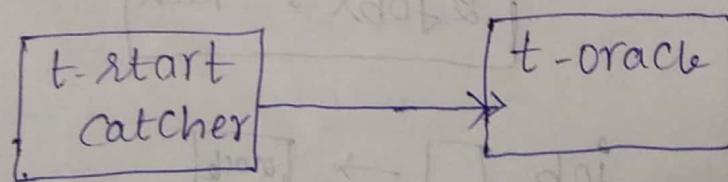


table Name.

Begin entry [start time - End time]  
 end entry [Job - Job]

if local host given as wrong?

Audit info : To Capture audit information.

\* t-log catcher: error log.

\* t-stat catcher: Begin, End → failure

process id. [job execution log]

job Run status: success / fail.

job Duration.

Bottle Neck: which component took time to run the job. Each component → Adv. Setting.

t-start catcher statistics.

origin [null]

Begin }  
End }

Component - { start time  
end time

Entire project: [project level setting]:

file :-

→ edit project prop

stats and logs

use stat

use logs

on data base

status table.

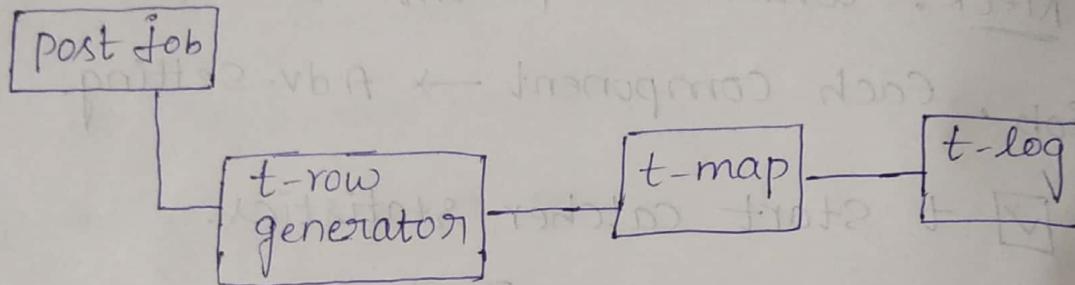
logs table.

Each job :- [job level setting] :-

status & logs :  status table  
 logs table.

transient catcher statistics

Audit into count:



\* Add Columns :-

job name

pro name

Pid

emp-file - count

tgt 1 - count [table name]

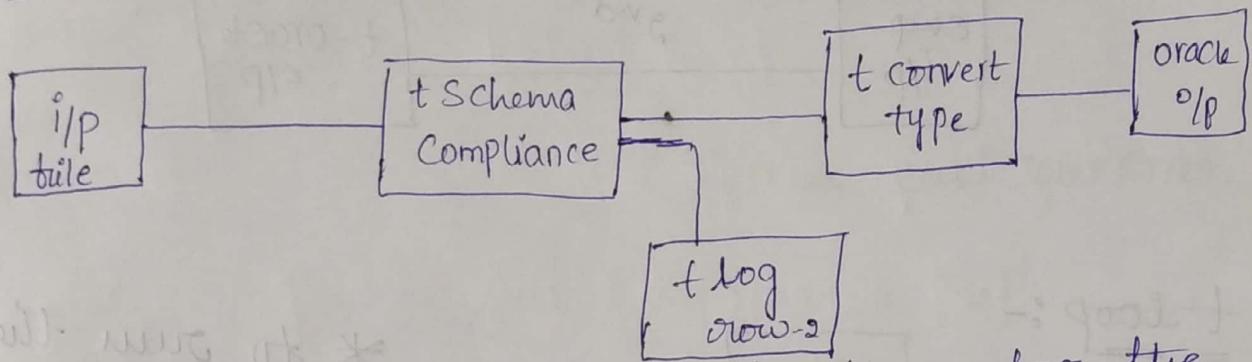
date file - count

NB-line global Variable:-

↓  
Gives you the number of records processed by particular component.

06/09/19

### t- Schema Compliance :-



- \* Read the i/p file as string datatypes for the columns.
- \* Connect to schema compliance. You can define options in the schema check. ~~auto~~ compliance and defined your rules.

Example's data type

date pattern

Nullable

max length.

\* Schema Compliance

\* convert type

\* Generate DOC-HTML

\* Built-in

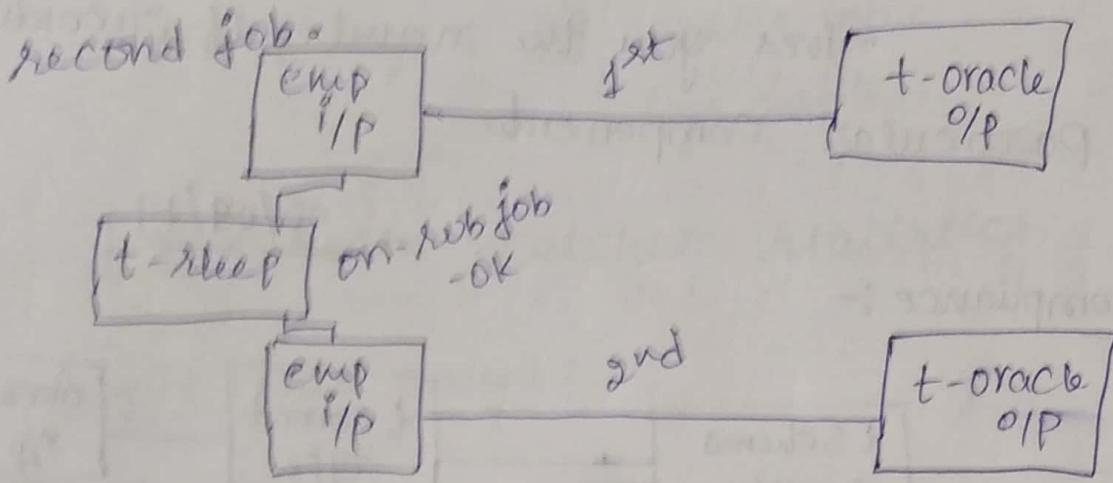
\* memory issue

\* Scheduling

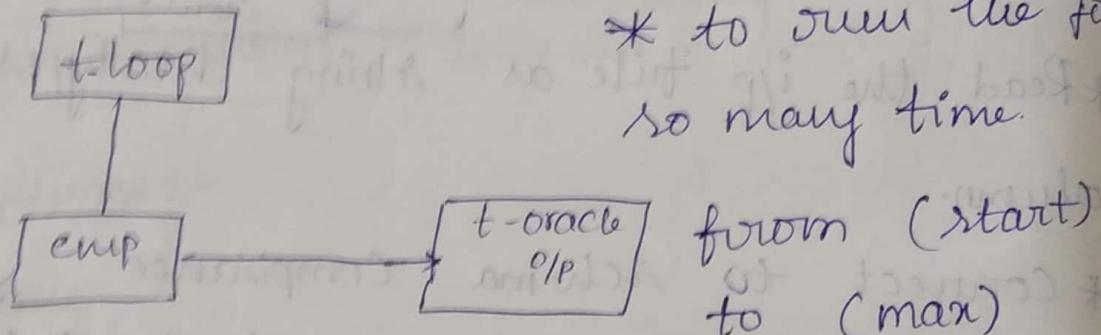
### Repository vs Built-in:

↳ centralised  
(Best) ↓  
Local

t-sleep :- It gives some gap to run [09/09/19]

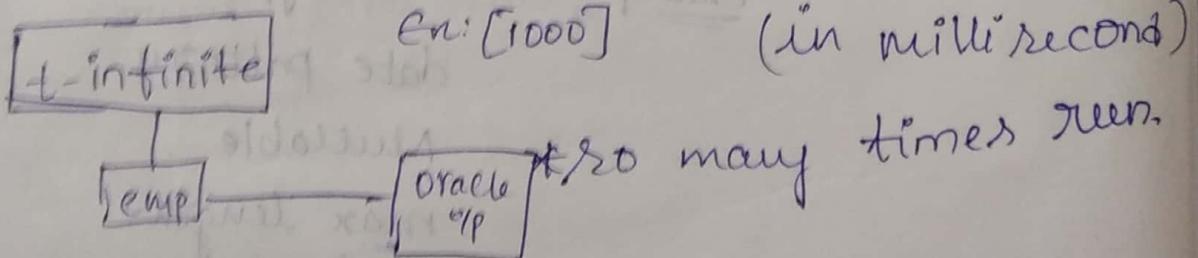


t-loop :-

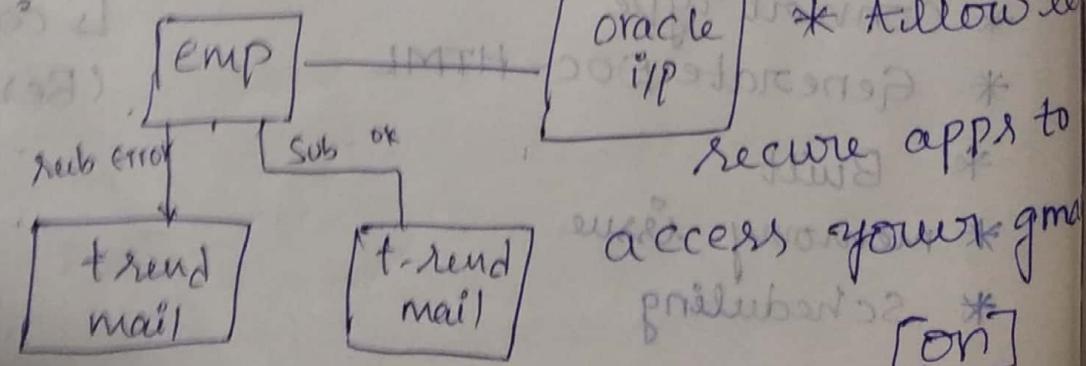


t-infil

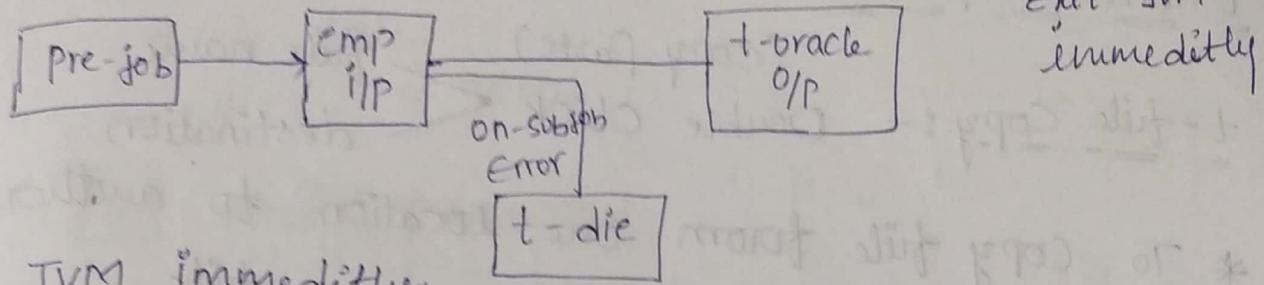
t-infinite loop :- \* wait at each iteration



t-l Send mail :-



t-die:-



Exit JVM  
immediately

Exit JVM immediately:-

\* to kill the job any time on some trigger.

\* like on-sub-job-ok (or) on-sub-job-error.

\* we go with t-die.

t-system:-

\* to execute unix, windows unique, (or) das commands to talend.

\* How to execute batch file from command prompt?

t-SSH:- [Secure shell] → not in a local machine.  
→ executes unix/das command remote box.

t-Run job:-

\* Create job-workflow.

\* drag all job. what you want to run.

## t-file Components :-

10/09/19

t-file Copy :- Double click < file name  
(COPY Part)

\* To copy file from one location to another.

\* If we enable  remove source file.

[cut Parte].

\* COPY only single file

Rename  $\Rightarrow$  Name change

copy directory  $\rightarrow$  full directory / file mode

## t-file delete :-

To delete file.

single [folder/file].

fail on error. [empty file].txt error

delete file / folder.

## t-file Archive :-

Dir - file  $\rightarrow$  COPY

Archive here. (Rename)

format: ZIP.

compress level: best

encrypt files [password]  
if enable

We have a chance to password.

## t-file unarchives

Arch: file  
dir: - ↘ copy

Need password.

## t-file properties: [contains edit schema]

file: —

it give properties.

t-file →  logrow

## t-file touch:

Create empty file.

[0-bytes]

[unia-touch

create empty file]

file name touch.csv

file name touch.txt

→ .x1

## Scenarios:-

folder contain multiple files.

2 log files:

log file 1

total files total

size

log file 2

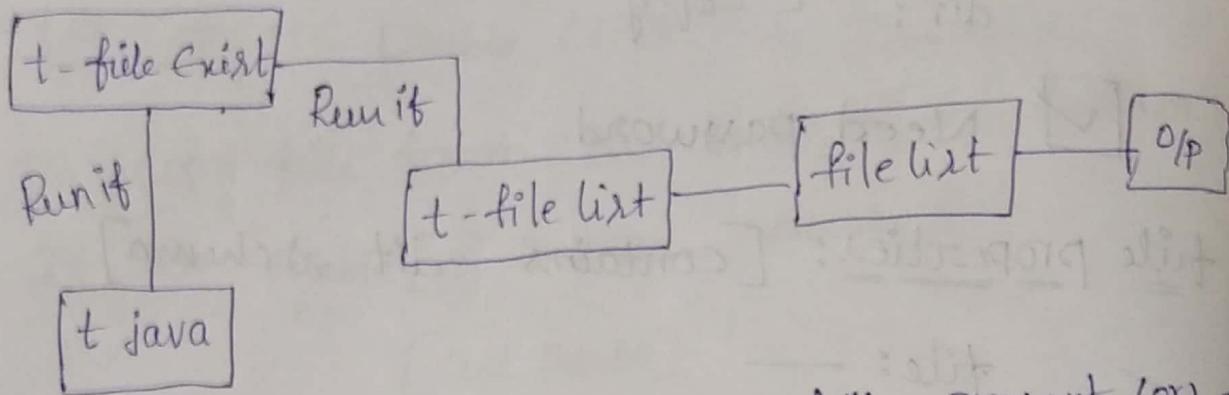
file names time

stamp

asc.  
order.

11/10/19

## t-file Exist :- - Runif



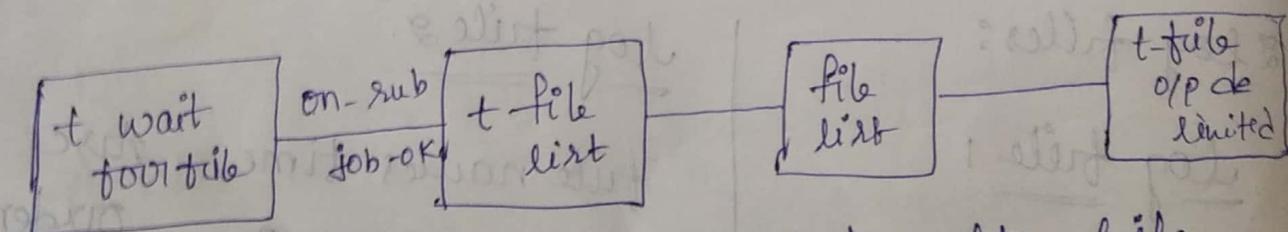
- \* it will check whether the file present (or) not. if present there otherwise false.
- \* if we file exit component we will use "Runif" condition.

\* Runif → right click give condition.  
 (it is a boolean type).

## t-wait for file :- file mask

choose Inculde present file

exit loop → Then



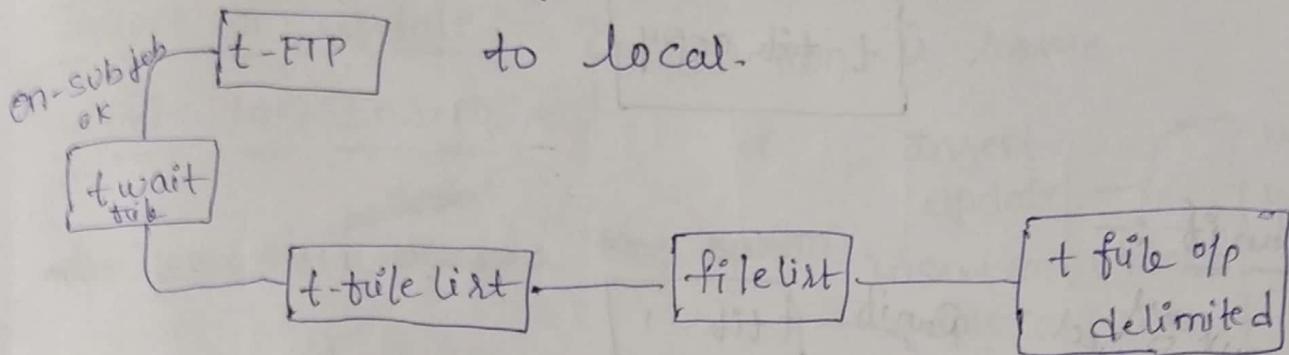
- \* it will keep on wait for the file.
- \* when the done file appends it will run the job.

\* In the t-wait for file :- we need to enable "include present".

winscp: (windows secure copy):  
unix box

\* To copy the file from <sup>unix box</sup> to windows, we go with "winscp-tool" (or) file "file".

t-FTP get: \* it get the files from FTP Server to local.

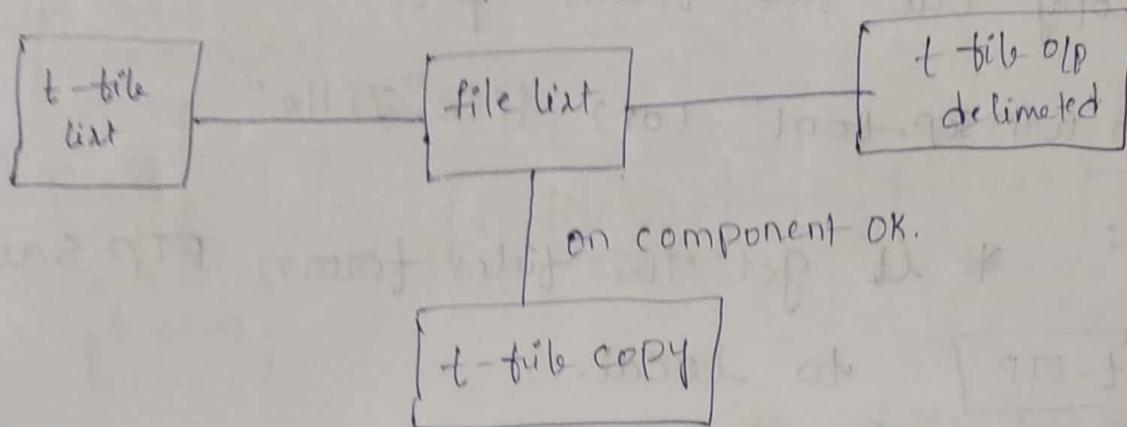


t-FTP PUT :- It puts the files from local save to ftp.

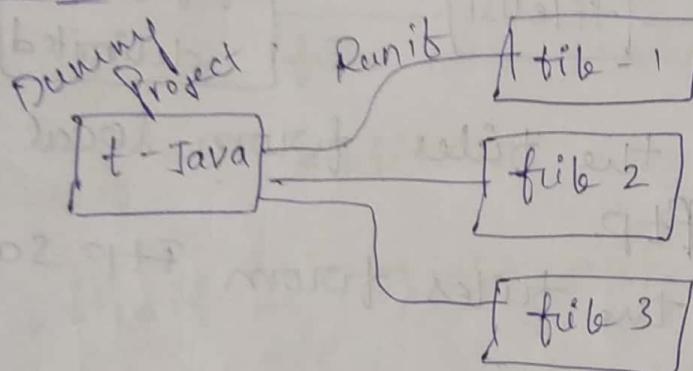
t-FTP Delete :- To delete the files from ftp save

(e) following + " \ backslash \ (d) following + " \ backslash \ file

On subjob OK :-



Run if :-



Condition: context.job.id = 1

context.job.id = 2

context.job.id = 3

~~13/09/19~~  
\* t oracle input:

Improving Performance t-oracle input :-

- \* Go to Advance setting.
- \* and select use cursor option.
- \* parallel hint in the oracle input.

~~Select /\*+parallel(6)\*/ query; Select /\*+parallel(5)~~

\* Check whether Indexs created on Key Column (or) Group by Column. Else if not create them.

## t-oracle outputs

### Insert (or) update :-

### update (or) Insert :-

~~update schema option~~

→ Go to Advanced setting

increase the commit size

and Batch size.

Insert  
update

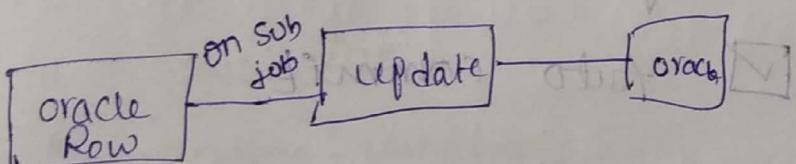
Insert (or) update  
update (or) Insert

Delete

use  
key  
column

t-oracle Row :- \* It is used to truncate.  
ion delete (or) update by  
execute -

Query: "truncate"



to execute some statement directly

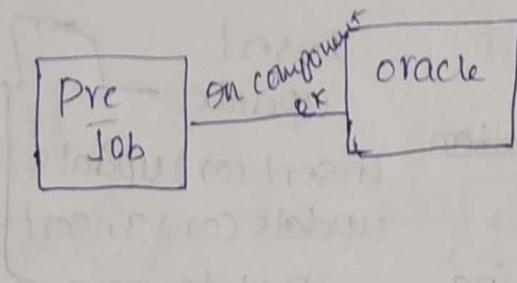
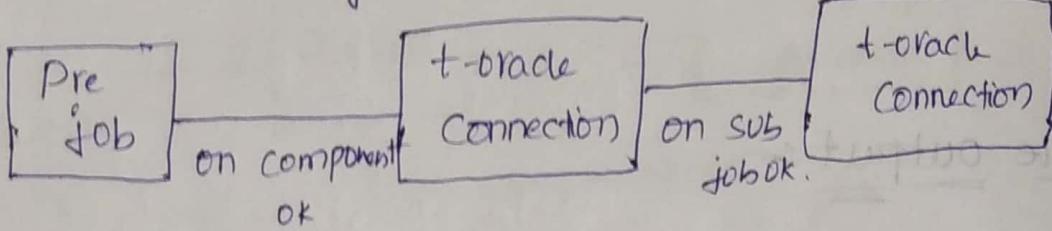
on a database.

Multiple jobs: To Run by only one connection.

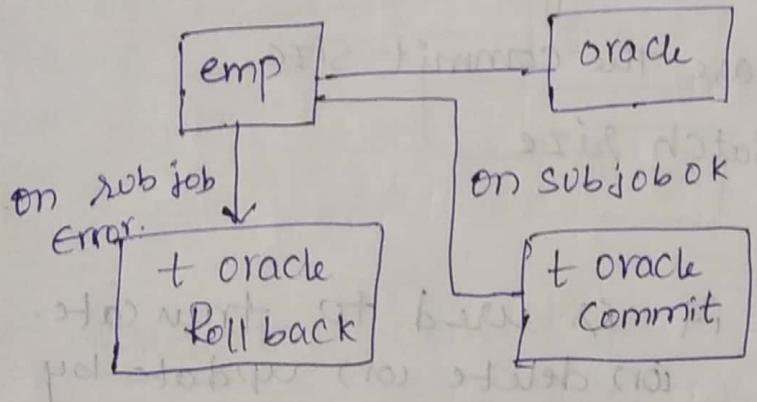
t oracle Connection: }  
to run multiple jobs

## t-oracle connections

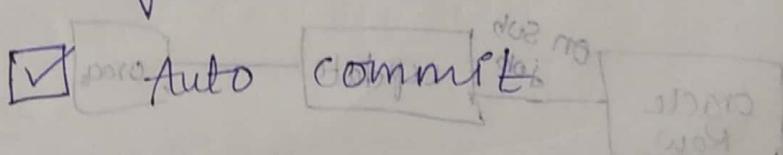
- use an existing connection.



\* t-oracle commit is used in t-oracle connection only.



\* Go to Advance settings



When use connection.

must use t Commit / t Roll Back.

- \*  use an existing connection.

Component to every component

## t oracle output Bulk execution:

Data loaded speedy.

- \* This component contain some rules.
- \* we didn't use this component as t oracle o/p.
- \* it should be plain table.  
we use this when it's only truncated.
- \* To use this component target table should not contain any constraints like primary key, foreign key, unique key and triggers and indexes.
- \* Target table should be always truncate and load.

for this we need to create table.

create table bulk...

( "empno" int,

    "empname" varchar2(50),

    "Sal" int,

    "hra" int,

    "deptno" int,

    "dept name" varchar2(30),

    "loc" varchar2(10),

)

\* While creating a target table, make sure -  
target table columns should be in double  
quotes.

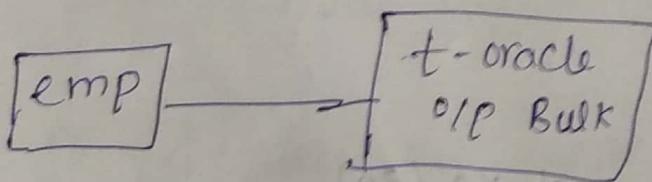
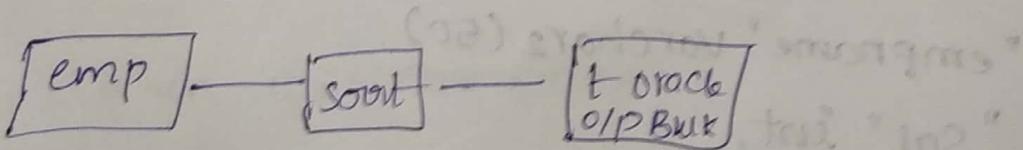
\* Action on table default.

file name : file location.

3 files : total file . csv . (data file)  
ctr. file (control)

log file → time span }  
Begin  
end }

- Bulk load loads data from file to the table with high speed.
- Before loading it creates csv file and ctrl file and loads data from csv to table directly.
- To use this component target table should be empty (or) we can go with truncate option



\* At this time we had another component

t-oracle Bulk creation.

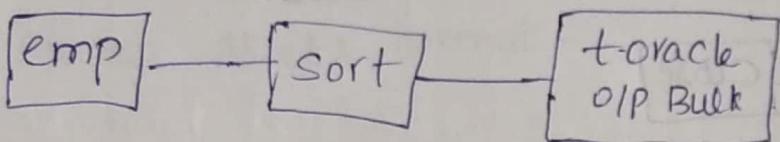
## oracle execution component:

table.

file.

copy schema.

deactivate



t-oracle execution comp

copy file

file

copy schema

\* to know exact error (where it should occurred)

Go to code.

\* Search line no [ctrl L -].

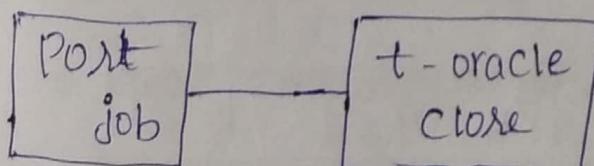
(or)

\* Debug Run.

\* When we use t-oracle connection.

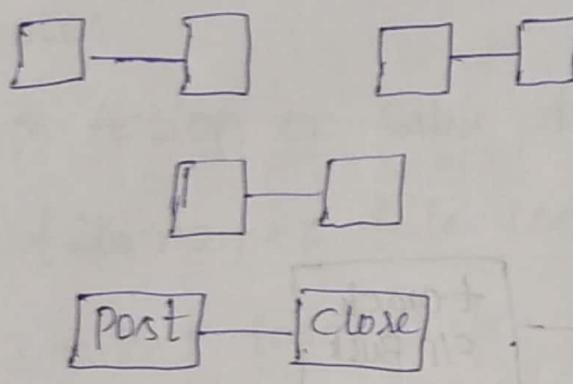
\* if it does not close.

then use t-oracle close.



Pre job

↳ connection.



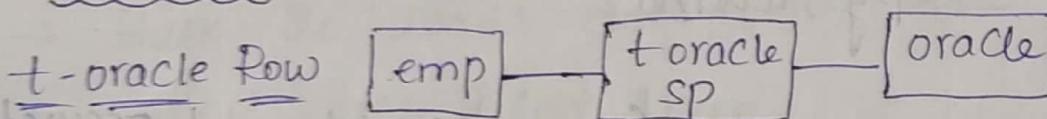
Best practices:

- context create
- tstat catcher
- tlog catcher
- schema compliance check

Stored procedure:

How to call stored procedure in Talend?

t-oracle Sp:-



double click

Sp name : "tax cal".

edit Schema: extra column tax.

Parameters: HRA  
empno }  
Sal }

HRA  
Tax.

we go with t-oracle SP / t-oracle Row



Reading:

19/09/19

Excel file threw Talend:

1st Check properties xlsx, 10MB

Read excel (2007 file format (xlsx))

Read excell arrets format.

Generation Mode:

[less memory consumed for large excel (Event model)]

Select file → Choose file

Next --- finish

Take as input.

if other excell.

\* All sheets have same metadata, select all;

Next --

\* four small file generation mode:

[

\* If we want sheet name also in O/P.

\* Click on t map

\* Add column as River

Select → Current Sheet [global variable]

t-file output excel:

\* if we want excel file as output: (create excel file)

[2007 file format] (Read excel)

file name.

--- /folder/excelout.xlsx & sheet name:

S To put that sheet any where like middle of sheet / corner ---

is absolutely pos. → shift + 10002

(63) first cell u[5|-] --- y[8|-]

to Scan metadata for

create file positional: ]

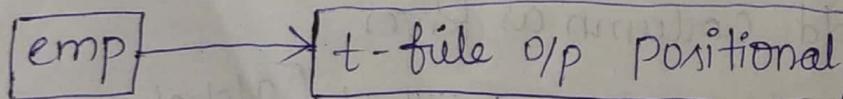
Select file

Select lengths, next ---

take as input → ---

t-file output Positional:

Output Positional: T file output Positional.



~~20/09/19~~

file name ✓

include header.

Reading XML file Input:

Scan meta data

loop limit [-1]

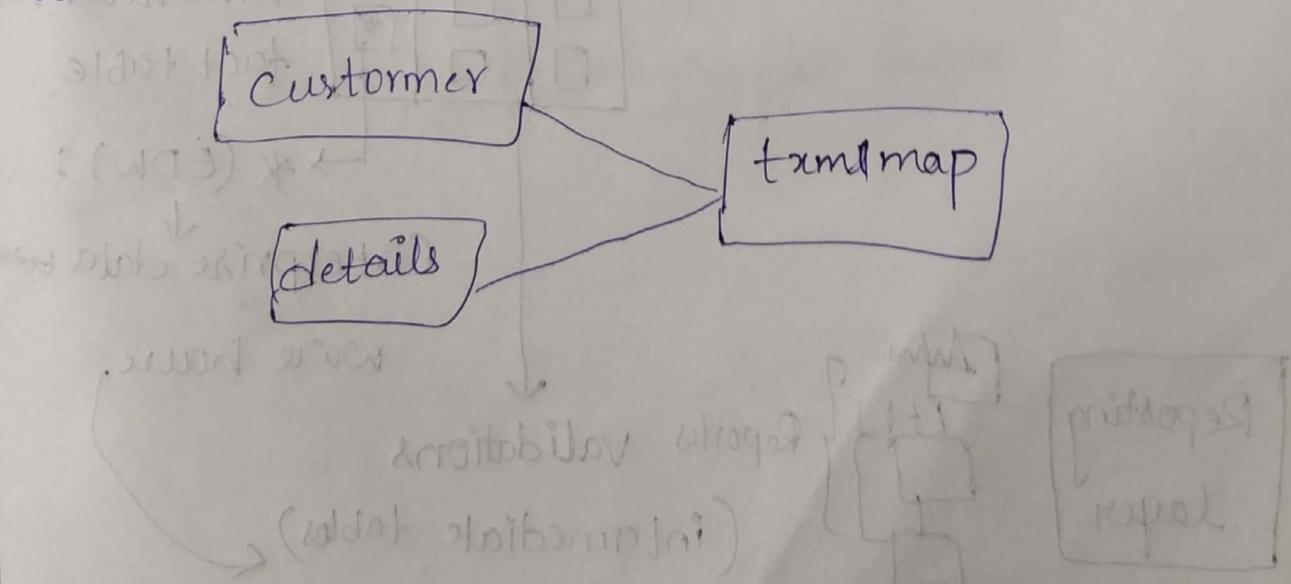
Schema viewer → Next, unlimited Records

→ Absolute x path. Expression [-1].

Ex: / power mart / Repository / folder / source /  
source field.

→ field to extract Expression.  
Ex:- copy → drag here.

XML creation (output) :- txml map component is used.



## JSON file :-

Create

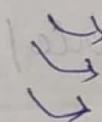
① i/p

Next

Browse file

Abs. path [ - ].

→ Relative / absolute path expression.



## Sales force :

oracle

Excel

XML → Talend

GE Power

Greenplum

□	□	□
□	□	□
□	□	□
□	□	□
□	□	□

\* here tables are  
called as  
dimension table

\* (EDW) :

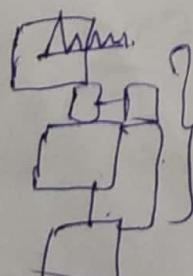
Enterprise data warehouse

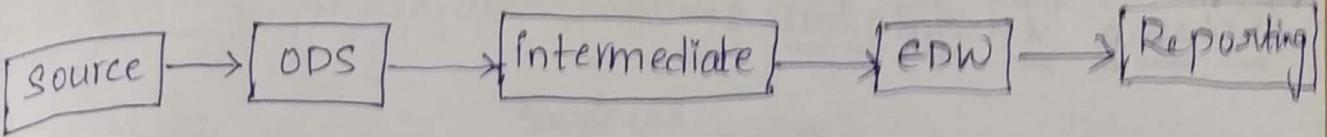
ware house

Reports validations

(intermediate tables)

Reposting layer





Dimensions :-

Facts :-

→ with respect to dimensions only you can derive facts.

~~21/09/19~~

\* Dimensional table will have a detailed information.

\* Dimensional table will have primary key.

Fact tables :-

\* fact table is a foreign keys of all and fact values.

\* We can make an entry only when the recover existing dimensional table

Schema in database :-

\* Star, snowflake.

Star schema :- way of representing fact table :-

dept dim

Shopid  
Prod id  
emp id

Prod.dim

Pro id  
pro name

Talend dim

cal id  
date  
month  
week

shop dim

Shop id  
Shop name  
Shop dimension

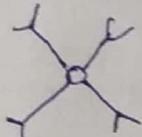
Star type:

shop Prod  
|  
Tal emp

empdim

emp id  
emp name  
sal  
HRA

Snowflake :- → We are mostly used in projects.



Business

4 dept → 10, 20, 30, 40

4K employees

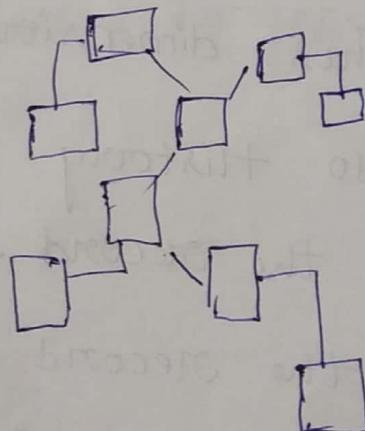
empid	ename	sal	tfa	dept name
:			:	10
1000			:	10
1001			:	20

\* Duplicate are there.

→ then remove all columns.

again ?  
create  
table.

D.no	D name	loc
10	-	-
20	-	-
30	-	-



Speed : star

memory (less) : snowflake.

Galaxy rehema :-

Dimensions :- 3 types of dimensions.

\* Slowly changing dimension :-

{ SCD-I  
SCD-II  
SCD-III

\* SCD-I :-

eno	ename	sal	HRA	Dept.no
100	A	1000	100	10
101	B	1800	180	20
102	C	1500	150	10
103	D	1200	120	20

\* This dimensions always maintain

\* No history

If the record is coming as a new

if the record is coming with update

if the record is coming

ignore it (or) Rejected it.

\* SCD-II \* 4 extra columns is there.

* emp Key	empno	ename	sal	HRA	D.NO	Scdstartdate
1	100	A	1000	100	10	17/08/19
2	101	B	1800	180	20	17/08/19
3	102	C	1500	150	10	17/08/19
4	104	D	1200	120	20	17/08/19

Scd end date	Activeflag
31/12/1999	y

\* update

\* we are used in Primary key.

\* Here primary key is a emp key.

\* SCD to maintains complete history.

\* By adding 4 extra columns in the target. Call Surrogate key (empkey).

\* Sc Started, Enddate

\*

Surrogate key:

\* Which

primary key in SCD-II

table.

which is generated by ETL

Ex: Talend.

\* Using sequence numbers (numeric sequence functions)

- \* If the record
  - \* Insert it in the target with as new sequence number start date as current date.
  - \* End date as max date (31/12/1999).
  - \* Active flag as 'Y' or '1'.
- \* If the record is coming with
  - \* Insert it as a new record with new sequence number and update the old records with end date as current date.
  - \* Active flag as 'N'.
- \* To update the old record we use Surrogate key as the key column.

SCD :-

eno	ename	Pre.Sal	Curr.Sal	Pre.HRA	Curr.HRA	Deptno	Etldate
100	A	1000	1000		100	10	
101	B	1800	-1900		180	20	
102	C		1500	150	250	10	
103	D		1200		120	20	

Etldate
18/8/19
18/8/19
18/8/19
18/8/19

Same - ignore  
change - Pre  
New Record - Current

25/09/17

## SCD3: Passing context parameter

### How to

Passing context variables from parent to child:-

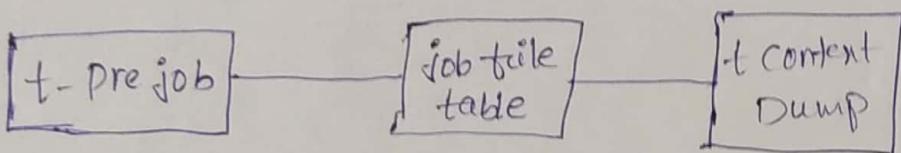
- \*  transmit whole context.

Only some variables changes:

Content parameter →

(+) Button.

Child to parent :-



\* Go to child Job take context Dump component and content to Buffer O/P.

\* Come to parent job.

\* Select tRunJob.

\* Go to component job.

\* Click on <sup>COPY</sup> child job schema.

\* Now tRunJob map to t content load using glow main.

6/9/19 4 Tables (frame work)

\* Job execution log

t - set Global Var.



Key → max key.

row6.max

\* Audit table

\* Error file 1

\* Error file 2.

Job execution log: Execution id

proid

Execution id

job name