

# DigiClips Database Analysis

Pavan Manjunath

## Implementation Plan: Data Retention & Optimization

### Goal Description

The primary goal is to optimize the dc database by enforcing strict data retention policies (deleting records > 90 days), improving search relevance and performance, streamlining database backups, and enhancing data integrity through validation.

### Info

TvVTT and RadioClips store timestamps as VARCHAR. To efficiently delete records older than 90 days, these MUST be converted to DATETIME. This is a potentially blocking operation for large tables. Use of STR\_TO\_DATE in delete queries is possible but performs poorly.

### Proposed Changes

#### Phase 1: Data Retention & Cleanup (Priority)

We will focus on automating the cleanup of old data to maintain database performance and compliance.

##### 1. Schema Updates for Date Handling

- Target Initial Tables( later on extend to others after verification) : TvVTT, RadioClips, RadioClipsB.
- Action:
  1. Add a new TStamp\_dt DATETIME column.
  2. Populate it using STR\_TO\_DATE(TStamp, format).
  3. Rename/Swap to replace the old column. Note: If table size is massive, we may need to do this in batches.( so this is optional)

##### 2. Logging Infrastructure

- New Table: cleanup\_logs (or enhance existing cleanup\_log).
  - Columns: id, table\_name, deleted\_count, execution\_time, status, timestamp.

##### 3. Stored Procedures for Cleanup

- Optimize/Create: Cleanup\_SRT\_OldRecords and Cleanup\_Radio\_OldRecords.

- Logic:
  - Delete simple records > 90 days.
  - Log the number of deleted rows to cleanup\_logs.
  - Handle SRT, SRT21, TvVTT, RadioClips.

#### **4. Automation(to be reviewed)**

- Create MySQL Event Daily\_Data\_Purge to call the cleanup procedures every 24 hours.

#### **Phase 2: Database Dump Management( can be modified acc post Phase 1 results)**

Improve the backup process to be more robust and organized.

- Structure: Create specific directories
- Split Dumps:
  - Schema Dump: Full structure of all tables (including massive ones like SRT).
  - Data Dump (Config/Ref): Data for smaller tables (Stations, Users, Config).
  - Exclusion: Do NOT dump data for SRT or RadioClips tables (too large).

#### **Phase 3: Search Optimization ( based on results from indexing)**

Improve the quality and speed of search results.

- Relevance Ranking: Modify Bool\_Search (or creating Relevant\_Search) to order results by:
  1. Match frequency (count of hits).
  2. Recency (newer items first).
  3. Category weighting (if applicable).
- Indexes: Review execution plans for search queries and add specific secondary indexes where full-text search is insufficient.

#### **Phase 4: Data Validation**

- Validation Procedures:
  - Update Insert procedures to validate Email format (regex).

- Validate State codes against the State reference table.
- Error Logging:
  - Create Validation\_Errors table to capture rejected inserts/updates for review.

## Verification Plan

### Automated Tests

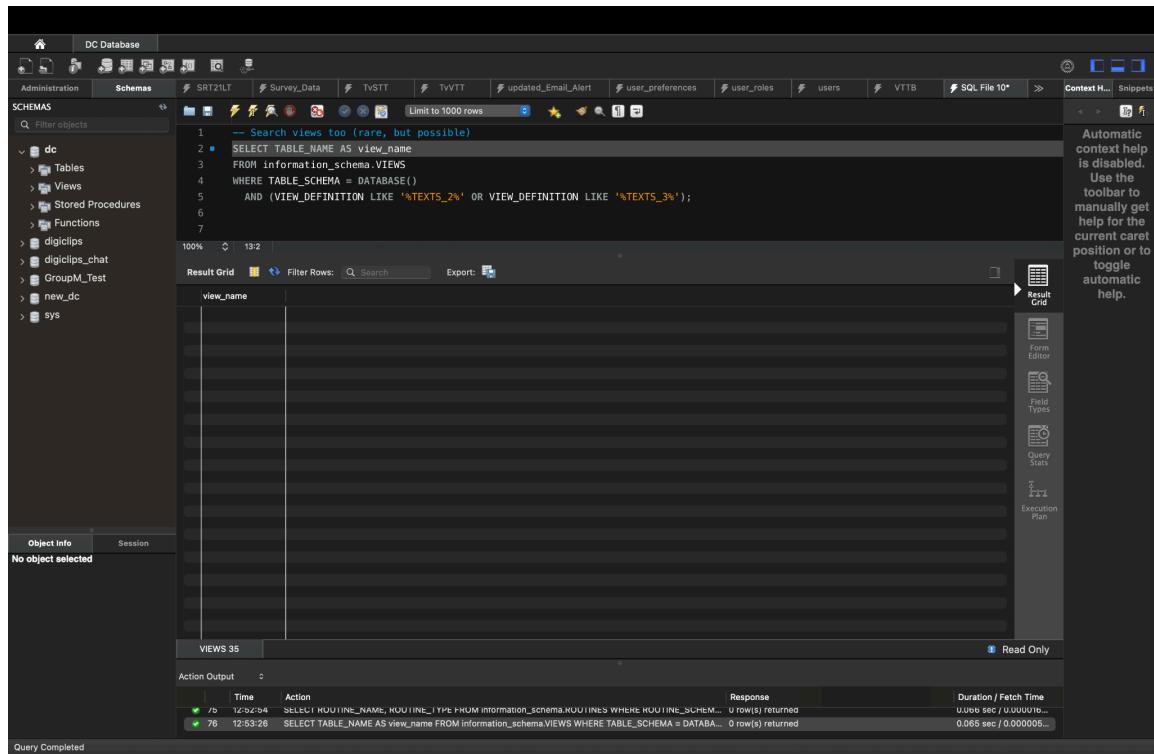
- Retention Test: Insert dummy records dated 91 days ago, run the procedure, verify they are gone and logged.
- Validation Test: Attempt to insert invalid emails/states and verify rejection/logging.

### Manual Verification

- Search Quality: Run sample queries before and after optimization to compare result relevance.

All the additional tables to be added can have different table names, just mentioned them for an example.

Queries:



The screenshot shows the Oracle SQL Developer interface. The left sidebar displays the schema structure under 'DC Database' with the 'Tables' node expanded, showing 'dc', 'Tables', 'Views', 'Stored Procedures', 'Functions', 'digiclip', 'digiclip\_chat', 'GroupM\_Test', 'new\_dc', and 'sys'. The main pane shows a query results grid titled 'Result Grid' with the column 'view\_name'. Below the grid, the status bar indicates 'VIEW 35' and 'Read Only'. At the bottom, the 'Action Output' section shows two rows of execution details:

Action	Time	Response	Duration / Fetch Time
75	12:52:04	SELECT ROUTINE_NAME, ROUTINE_TYPE FROM information_schema.ROUTINES WHERE ROUTINE_SCHEMA = 'DC'	0 row(s) returned 0.000 sec / 0.000016...
76	12:53:26	SELECT TABLE_NAME AS view_name FROM information_schema.VIEWS WHERE TABLE_SCHEMA = DATABASE()	0 row(s) returned 0.065 sec / 0.000005...

The status bar at the bottom left says 'Query Completed'.

The screenshot shows the MySQL Workbench interface with the following details:

- Toolbar:** Includes icons for Home, Database, Schemas, Tables, Views, Stored Procedures, Functions, and various search and export options.
- Schemas:** Shows the current schema is "dc".
- Search Bar:** Contains the query: "SELECT ROUTINE\_NAME, ROUTINE\_TYPE FROM information\_schema.ROUTINES WHERE ROUTINE\_SCHEMA = DATABASE() AND (ROUTINE\_DEFINITION LIKE '%TEXTS\_2%' OR ROUTINE\_DEFINITION LIKE '%TEXTS\_3%');".
- Result Grid:** Displays the results of the query, which are currently empty.
- Object Info:** Shows "No object selected".
- Action Output:** Shows the execution log:

Time	Action	Response	Duration / Fetch Time
12:41:08	SHOW INDEXES FROM digiclip	10 row(s) returned	0.001 sec / 0.00012...
12:52:23	SELECT ROUTINE_NAME, ROUTINE_TYPE FROM information_schema.ROUTINES WHERE ROUTINE_SCHEMA...	0 row(s) returned	0.005 sec / 0.00063...
- Help:** A context help message states: "Automatic context help is disabled. Use the toolbar to manually get help for the current caret position or to toggle automatic help."

The screenshot shows the MySQL Workbench interface. The left sidebar displays the 'SCHEMAS' tree, which includes the 'dc' schema and its objects: Tables, Views, Stored Procedures, Functions, digiclip, digiclip\_chat, GroupM\_Test, new\_dc, and sys. The main area shows a query editor with the following SQL command:

```
SHOW INDEX FROM RadioClips;
```

The results grid displays the index information for the 'RadioClips' table. The columns in the grid are: Table, Non\_unique, Key\_name, Seq\_in\_Index, Column\_name, Collation, Cardinality, Sub\_part, Packed, Null, Index\_type, Comment, Index\_comment, Visible, and Expression. The data in the grid is as follows:

Table	Non_unique	Key_name	Seq_in_Index	Column_name	Collation	Cardinality	Sub_part	Packed	Null	Index_type	Comment	Index_comment	Visible	Expression
RadioClips	0	PRIMARY	1	ID	A	30				BTREE			YES	
RadioClips	1	TEXTS	1	TEXTS	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS	2	FName	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS	3	Categories	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS_2	1	TEXTS	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS_2	2	NULL	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS_2	3	Categories	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS_3	1	TEXTS	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS_3	2	FName	NULL	30				YES	FULLTEXT		YES	NULL
RadioClips	1	TEXTS_3	3	Categories	NULL	30				YES	FULLTEXT		YES	NULL

The bottom status bar indicates 'Result 31' and 'Read Only'. The Action Output pane shows the following log entries:

Action	Time	Response	Duration / Fetch Time
SHOW CREATE TABLE IVV1	12:39:41	1 row(s) returned	0.008 sec / 0.000010...
SHOW INDEX FROM RadioClips	12:40:08	10 row(s) returned	0.071 sec / 0.000012...

The right side of the interface features several floating toolbars: Context Help, Automatic context help is disabled. Use the toolbar to manually get help for the current caret position or to toggle automatic help.; Result Grid; Form Editor; Field Types; Query Stats; Execution Plan.