

DavisBase

File Format Guide

Storage Definition Language (SDL)

Version 1.2 (2002-07-19)

Chris Irwin Davis

1. Introduction

DavisBase Nano is a relational database management system (DBMS) that implements a subset of the SQL specification. DavisBase is not intended for commercial industrial applications—but instead, for use in an academic setting for use by students to gain an understanding of relational database internals. This document assumes that the reader already possesses a basic understanding of (1) relational database theory, (2) basic SQL, and (3) hexadecimal file viewers/editors such as the built-in Unix command `hexdump`.

This File Format Guide describes only the Storage Definition Language (SDL) of DavisBase. It is intended to be a subsection of the overall DavisBase Specification, which also includes a data query language (DQL), data definition language (DDL), data control language (DCL), and data manipulation language (DML). Additionally, an API will be included in a later version that is based on a subset of the JDBC specification.

DavisBase SDL file architecture borrows liberally from that of several open-source databases including SQLite, MySQL, and PostgreSQL. In fact, one of the primary goals of DavisBase is to be an educational platform for learning the internals of relational databases. Implementing file structures that mimic those of widely known DBMSs provides insight into the mechanics of relational databases in general.

We note that DavisBase is not fully ACID¹ compliant! There is neither a system log nor a formal mechanism for rollback or recovery of failed transactions.

2. Files

Despite its similarity in name to the real world physical object called a “file”, which is made of paper that can be touched and held, a computer file is a virtual construct. The term “file” can be defined in various ways since it is a logical construct that may be stored in discontinuous subsections or blocks called *pages*, either in volatile memory (RAM, “main memory”) or non-volatile memory (disk, SSD, “thumb drive”).

A file may be a virtual construct in volatile storage (i.e. main memory or RAM). It may also be a construct that is physically stored in non-volatile memory like a hard disk, SSD, or thumb drive.

DavisBase uses two kinds of files: **table files** and **index files**. We note that fully-featured commercial databases utilize additional file types—e.g. view files, script files, macro files, etc. These additional file types are not included in v1.0 of DavisBase.

DavisBase uses a *file-per-table* approach. That is, every table and every index are each stored in a single, separate OS file. Each table file is stored in an OS file whose name is `<tablename>.tbl`, and each Index File is stored in a file whose name is `<tablename_columnname>.ndx`. Both Table Files and Index

All DavisBase files are implemented as page-based files. That is, files are subdivided into same-sized logical sections called *pages*. DavisBase supports page sizes of 2^n bytes, where n is an integer in the range 9-15. All pages of a given file are required to be the same size. All files within the same table space must share the same page size. Once a table space has been initialized, the page size of all of its constituent files is fixed static.

Page size is represented in bytes. In this project, all pages are 0.5KB (512 bytes).

All page numbers are 4-byte integers

¹ [https://en.wikipedia.org/wiki/ACID_\(computer_science\)](https://en.wikipedia.org/wiki/ACID_(computer_science))

2.1. Offsets

Offsets are the mechanism used to locate elements within a file.² Offsets are expressed as non-negative integers that represent the number of bytes from a given point of reference. There are two kinds of offsets used by DavisBase, *file offsets* and *page offsets*.

File offsets indicate the location of an element expressed as the number of bytes from the beginning of a file. File offsets are most commonly used to identify the location of a page within a file—by multiplying page number time page size. For example, a file with a page size of 512 bytes, page 7 will begin at $512 * 7 = 3584$ bytes from the beginning of the file.

Page offsets indicate the location of an element expressed as the number of bytes from the beginning of a page.

2.2. Table Files

A table file stores the records (i.e. rows) of a given table in a single file.³

Table files are implemented as Bplus-trees where each node of the tree is a page of the file. Table Files store all record data for a given table. Table records reside solely in leaf pages of the B+ tree.

Table Files are not required to have any Index Files, i.e. they may have zero-to-many indexes. However, every Index File must have *exactly one* Table File that it indexes.

Each record in a Table File has an automatically generated unique identifier that is used by DavisBase for internal housekeeping. This is like a hidden extra column that is included in every table. Various databases use different names for this unique identifier: **rowid**, **row_id**, and **uid** are all common. DavisBase uses the name **rowid**. Each **rowid** is automatically generated based on the following strategy—The first record inserted into a given table is numbered “1”, then each subsequent **rowid** increases monotonically in integer increments with no gaps. If a record is deleted, its **rowid** is never reused. Each **rowid** is represented by a 4-byte two’s complement integer. Therefore, no table may contain more than $2^{31}-1$ (2147483647) unique records total over the lifetime of the table.

If a table schema is created without an explicit PRIMARY KEY, users may insert records that contain duplicate values for all user-defined columns. This would seem to violate relational theory which requires that all tuples (records) be unique. However, every record is still unique since each contains a unique value for this hidden **rowid**. Most relational databases use this strategy of a hidden unique identifier for internal housekeeping.

When a table file is created with a CREATE TABLE command it does not have any associated index files by default. However, if a column has been explicitly designated as a PRIMARY KEY, then an index file shall be automatically created that indexes this column, and that PRIMARY KEY column shall additionally be automatically designated to be UNIQUE.

All other columns, whether unique or not, shall not have an index implicitly created. Except for a PRIMARY KEY index, all other indexes must be created via an explicit CREATE INDEX command. All indexes may have at most a single column, including the PRIMARY KEY. Thus, neither composite indexes nor composite keys are supported.

² Offsets are the mechanism used by almost *all* relational databases to locate elements within a file.

³ In the relational model, the “tuples of a relation”.

2.2.1. Row IDs

Each record in a table file is uniquely identified by an internally assigned key called the *Row ID*. The Row ID is a column that is created automatically when a table is created with the column name **rowid**.

The **rowid** column *does not display* when a query is executed that references the **SELECT *** column wildcard. Only user defined columns will display. However, a user may explicitly add **rowid** to the attribute list in the **SELECT** clause to force it to display, e.g. **SELECT rowid, ***

By nature, **rowid** is unique, therefore it is a key—one used internally by DavisBase for housekeeping. The first **rowid** begins at 1, increments each new inserted record's by 1, and is always monotonically increasing.

2.3. Index Files

~~Index files provide efficient access paths to records within table files. All DavisBase indexes provide access paths to records based only on *single columns*. That is, multi-column indexes are not supported by DavisBase.~~

~~Index Files are implemented as B-tree files where each node a B-tree is a page of the file. Index Files store indexing information that provide an access path to a related Table File.~~

~~There is a single index entry in an index file for each distinct value in an indexed table column. However, a single index entry may have one to many **rowids** depending on how many records in its associated table share the same column value.~~

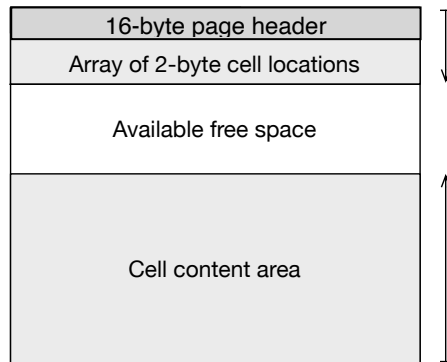
2.4. Database Catalog

The database catalog (i.e. meta-data) for DavisBase is stored in specially designated tables. These are also referred to as **system tables**.⁴ The Table Files associated with these are structured exactly the same as Table Files generated with user-defined tables.

⁴ MySQL groups these system tables together in what it calls the INFORMATION_SCHEMA.

3. Page Formats

All pages use the basic format of a fixed size header of 16-bytes followed immediately by an array of 2-byte integers that indicate the location of cells within the page. Cells contained in the cell area are written from high address upward.



3.1. Page Headers

Page Offset	Element Size	Description
0x00	1	<p>The one-byte flag at page offset 0 indicates the b-tree page type.</p> <ul style="list-style-type: none"> • A value of 2 (0x02) means the page is an index b-tree interior page. • A value of 5 (0x05) means the page is an table b-tree interior page. • A value of 10 (0x0a) means the page is an index b-tree leaf page. • A value of 13 (0x0d) means the page is a table b-tree leaf page. <p>Any other value for the b-tree page type is an error.</p>
0x01	1	Unused
0x02	2	The two-byte integer at offset 2 designates the number of cells on the page.
0x04	2	The two-byte integer at offset 4 designates the page offset for the start of the cell content area. A zero value for this integer is interpreted as 65536.
0x06	4	<p>The four-byte integer page pointer at offset 0x06 has a different role depending on the b-tree page type:</p> <ul style="list-style-type: none"> • Table or Index interior page - page number of rightmost child • Table or Index leaf page - page number of sibling to the right (This number is 0xFFFFFFFF if no pointer exists)
0x0A	4	The four-byte integer page pointer at offset 0x0A references the page's parent. If this is a root page, then the special value 0xFFFFFFFF is used.
0x0E	2	Unused
0x10	2 x cells on page	An array of 2-byte integers that indicate the page offset location of each data cell. The array size is $2n$, where n is the number of cells on the page. The array is maintained in key-sorted order.

Figure 1 - Page Header Format

The number of cells on any type of page is represented by a 2-byte two's complement integer. Therefore, no page may contain more than $2^{15}-1$ (32767) cells.

Page numbers for both types of file are represented by a 4-byte two's complement integer. Therefore, no file may contain more than $2^{31}-1$ (2147483647) pages.

3.2. Cell Formats

The **cell** is the basic unit of content for all types of pages. The cell content area of each page is located at the high offset values, i.e. cell content grows from the end of each page upwards.

The array of cell locations in the page header is an array whose order is maintained to be sorted according to criteria within the cells, depending upon file type. Table file cells are maintained in rowid sorted order.

File Type	Page Type	Cell Header			Cell Body
		4-byte int	2-byte int	4-byte int	N-byte array
		Left Child Page#	Bytes of Cell Payload	Rowid	Payload
Table	Leaf		✓	✓	✓
	Interior	✓		✓	
Index	Leaf		✓		✓
	Interior	✓	✓		✓

Figure 2: Cell Formats

Note that cells in interior pages have child pointers, since leaf pages don't have children.

The following two tables indicate the contents of the cell body (i.e. the “payload”) of record cells and index cells, respectively.

	Record Header		Record Body
	Number of Columns	List of Column Data Types	List of Column Data Values
Data Size	1-byte INT	Array of 1-byte INTs	N-bytes

Figure 3: Record Format (i.e. Payload of Table Leaf Cell)

3.3. Data Types

All data types supported by DavisBase have a 1-byte integer that designates their type. These are shown in both 1-byte hexadecimal (Hex) and equivalent decimal (Dec) values.

Hex	Dec	Data Type Name	Content Size (bytes)	Description
0x00	0	NULL	0	Value is a NULL (i.e. it takes no memory in the record body)
0x01	1	TINYINT, BYTE	1	Value is a signed, 8-bit twos-complement integer.
0x02	2	SMALLINT, SHORT	2	Value is a signed, big-endian, 16-bit twos-complement integer.
0x03	3	INT, INTEGER	4	Value is a signed, big-endian, 32-bit twos-complement integer.
0x04	4	BIGINT, LONG	8	Value is a signed, big-endian, 64-bit twos-complement integer.
0x05	5	FLOAT	4	Value is a single precision, big-endian, IEEE 754-2008 32-bit floating point number.
0x06	6	DOUBLE, REAL	8	Value is a double precision, big-endian, IEEE 754-2008 64-bit floating point number.
0x08	8	YEAR	1	Value is an 8-bit twos-complement integer. Both positive and negative numbers are supported in the range -128 to 127. This indicates a year with respect to the year 2000, i.e. 1872-2127
0x09	9	TIME	4	Value is a big-endian 32-bit twos-complement integer. Indicates time of day in milliseconds since midnight, i.e. "millis". Note that only values of [0-86,400,000) are valid. i.e. Hex [0x00000000-0x005c00)
0x0A	10	DATETIME	8	A big-endian unsigned LONG integer that represents the specified number of milliseconds since the standard base time known as "the epoch". It should display as a formatted string: YYYY-MM-DD_hh:mm:ss, e.g. 2016-03-23_13:52:23.
0x0B	11	DATE	8	A datetime whose time component is 00:00:00, but does not display.
0x0C + <i>n</i>	12 + <i>n</i>	TEXT		Value is an ASCII string of length <i>n</i> . C-style string null terminators are not used or needed.

Figure 4 - Data Types, Their 1-byte Code, and Implementation

Note that only strings of size 0-115 ASCII characters are supported. Any value 0x0C or above represents an ASCII string data type. Each length string 0-115 has its own unique data type. For example, 0x41 data type is a string of length 53 characters.

4. Data Storage Directory Structure

All data and meta-data shall be stored in a directory named **data**. This directory shall have the following structure:

```
data/
|
+---- catalog/
|
|          +-- davisbase_tables.tbl
|          +-- davisbase_columns.tbl
|
+-- user_data/
|
|          +-- table_name_1.tbl
|          +-- table_name_1.column_name_1.ndx
|          +-- table_name_1.column_name_2.ndx
|          +-- table_name_2.tbl
|          +-- table_name_2.column_name_1.ndx
|          +-- table_name_3.tbl
```

etc.

5. Database Catalog (Meta-Data)

The DavisBase Catalog consists of two tables containing meta-data about each of the user table. You may optionally choose to include meta-data about the two catalog files in the catalog itself. These two tables (and their associated implementation files) have the following table schema, as if they had been created via the normal **CREATE** command.

```
CREATE davisbase_tables (
  rowid      INT,
  table_name TEXT,
  root_page  SMALLINT,
  last_id    INT,
  record_count INT
);
```

```
CREATE davisbase_columns (
  rowid          INT,
  table_name      TEXT, -- optionally table_rowid INT
  column_name     TEXT,
  data_type       TEXT,
  ordinal_position TINYINT,
  is_nullable     TEXT
);
```


If you choose to include these two tables in the catalog itself, their content would initially be:

```
SELECT * FROM davisbase_tables;
```

rowid	table_name	root_page	record_count	avg_record
1	davisbase_tables	0	0	0
2	davisbase_columns	0	0	0

For the **data_type** column, you may choose to store the data type hex byte code found in Table 2.

```
SELECT * FROM davisbase_columns;
```

rowid	table_name	column_name	data_type	ordinal_position	is_nullable
1	davisbase_tables	rowid	INT	1	NO
2	davisbase_tables	table_name	TEXT	2	NO
3	davisbase_tables	root_page	SMALLINT	3	NO
4	davisbase_tables	last_id	INT	4	NO
5	davisbase_tables	record_count	INT	5	NO
6	davisbase_columns	rowid	INT	1	NO
7	davisbase_columns	table_name	TEXT	2	NO
8	davisbase_columns	column_name	TEXT	3	NO
9	davisbase_columns	data_type	TEXT	4	NO
10	davisbase_columns	ordinal_position	TINYINT	5	NO
11	davisbase_columns	is_nullable	TEXT	6	NO