

Oracle Berkeley DB

***Getting Started with
the
SQL APIs***

11g Release 2
(Library Version 11.2.5.1)



Legal Notice

This documentation is distributed under an open source license. You may review the terms of this license at:

Oracle, Berkeley DB, and Sleepycat are trademarks or registered trademarks of Oracle. All rights to these marks are reserved. No third-party use is permitted without the express prior written consent of Oracle.

Other names may be trademarks of their respective owners.

To obtain a copy of this document's original source code, please submit a request to the Oracle Technology Network forum at:

Published 10/25/2011

Table of Contents

Preface	v
Conventions Used in this Book	v
For More Information	v
Contact Us	vi
1. Berkeley DB SQL: The Absolute Basics	1
BDB SQL Is Nearly Identical to SQLite	1
Getting and Installing BDB SQL	1
On Windows Systems	1
On Unix	2
The Journal Directory	2
Unsupported PRAGMAs	2
Changed PRAGMAs	3
PRAGMA journal_size_limit	3
PRAGMA max_page_count	3
Added PRAGMAs	3
PRAGMA TXN_BULK	3
Miscellaneous Differences	4
Berkeley DB Concepts	5
Encryption	5
2. Locking Notes	7
Internal Database Usage	7
Lock Handling	8
SQLite Lock Usage	8
Lock Usage with the DB SQL Interface	9
3. Configuring the Berkeley DB SQL interface	11
Introduction to Environments	11
The DB_CONFIG File	11
Creating the DB_CONFIG File Before Environment Creation	12
Re-creating the Environment	12
Configuring the Database Page Size	12
Selecting the Page Size	13
Selecting the Database File Size	13
Configuring the In-Memory Cache	13
Administering Log Files	14
Setting the Log File Size	14
Configuring the Logging Region Size	15
Setting the In-Memory Log Buffer Size	15
Managing the Locking Subsystem	16
4. Administrating Berkeley DB SQL Databases	18
Backing Up Berkeley DB SQL Databases	18
Offline Backups	18
Hot Backup	18
Incremental Backups	19
About Unix Copy Utilities	20
Recovering from a Backup	20
Catastrophic Recovery	21

Syncing with Oracle Databases	21
Syncing on Unix Platforms	21
Syncing on Windows Platforms	22
Syncing on Windows Mobile Platforms	22
Data Migration	22
Migration Using the Shells	23
Supported Data and Schema	23
Replicating Berkeley DB SQL Databases	24
Preparing to use Replication with the Berkeley DB SQL API	24
Using Replication with the Berkeley DB SQL API	25

Preface

Welcome to the Berkeley DB SQL interface. This manual describes how to configure and use the SQL interface to Berkeley DB 11g Release 2. This manual also describes common administrative tasks, such as backup and restore, database dump and load, and data migration when using the BDB SQL interface.

This manual is intended for anyone who wants to use the BDB SQL interface. Because usage of the BDB SQL interface is very nearly identical to SQLite, prior knowledge of SQLite is assumed by this manual. No prior knowledge of Berkeley DB is necessary, but it is helpful.

To learn about SQLite, see the official SQLite website at: <http://www.sqlite.org>

Conventions Used in this Book

The following typographical conventions are used within in this manual:

Keywords or literal text that you are expected to type is presented in a monospaced font. For example: "Use the `DB_HOME` environment variable to identify the location of your environment directory."

Variable or non-literal text is presented in *italics*. For example: "Go to your *DB_INSTALL* directory."

Program examples and literal text that you might type are displayed in a monospaced font on a shaded background. For example:

```
/* File: gettingstarted_common.h */
typedef struct stock_dbs {
    DB *inventory_dbp; /* Database containing inventory information */
    DB *vendor_dbp;    /* Database containing vendor information */

    char *db_home_dir; /* Directory containing the database files */
    char *inventory_db_name; /* Name of the inventory database */
    char *vendor_db_name; /* Name of the vendor database */
} STOCK_DBS;
```

Note

Finally, notes of interest are represented using a note block such as this.

For More Information

Beyond this manual, you may also find the following sources of information useful when using the Berkeley DB SQL interface:

- [Berkeley DB Installation and Build Guide](#)
- [Berkeley DB Programmer's Reference Guide](#)

To download the latest documentation along with white papers and other collateral, visit <http://www.oracle.com/technetwork/indexes/documentation/index.html>.

For the latest version of the Oracle downloads, visit <http://www.oracle.com/technetwork/database/berkeleydb/downloads/index.html>.

Contact Us

You can post your comments and questions at the Oracle Technology (OTN) forum for Oracle Berkeley DB at: <http://forums.oracle.com/forums/forum.jspa?forumID=271>, or for Oracle Berkeley DB High Availability at: <http://forums.oracle.com/forums/forum.jspa?forumID=272>.

For sales or support information, email to: berkeleydb-info_us@oracle.com You can subscribe to a low-volume email announcement list for the Berkeley DB product family by sending email to: bdb-join@oss.oracle.com

Chapter 1. Berkeley DB SQL: The Absolute Basics

Welcome to the Berkeley DB SQL interface. If you are a SQLite user who is using the BDB SQL interface for reasons other than performance enhancements, this chapter tells you the minimum things you need to know about the interface. You should simply read this chapter and then skip the rest of this book.

If, however, you are using the BDB SQL interface for performance reasons, then you need to read this chapter, plus most of the rest of the chapters in this book (although you can probably skip most of [Administrating Berkeley DB SQL Databases \(page 18\)](#), unless you want to administer your database "the Berkeley DB way").

Also, if you are an existing Berkeley DB user who is interested in the BDB SQL interface, read this chapter plus the rest of this book.

BDB SQL Is Nearly Identical to SQLite

Your interaction with the BDB SQL interface is almost identical to SQLite. You use the same APIs, the same command shell environment, the same SQL statements, and the same PRAGMAs to work with the database created by the BDB SQL interface as you would if you were using SQLite.

To learn how to use SQLite, see the official [SQLite Documentation Page](#).

That said, there are a few small differences between the two interfaces. These are described in the remainder of this chapter.

Getting and Installing BDB SQL

The BDB SQL interface comes as a part of the Oracle Berkeley DB download. This can be downloaded from the [Oracle Berkeley DB download page](#).

How you install the BDB SQL interface differs depending on whether you are using a Unix or a Windows system.

On Windows Systems

The BDB SQL interface is automatically built and installed whenever you build or install Berkeley DB for a Windows system. The BDB SQL interface dlls and the command line interpreter have names that differ from a standard SQLite distribution as follows:

- `dbsql.exe`

This is the command line shell. It operates identically to the SQLite `sqlite3.exe` shell.

- `libdb_sql150.dll`

This is the library that provides the BDB SQL interface. It is the equivalent of the SQLite `sqlite3.dll` library.

On Unix

In order to build the BDB SQL interface, you download and build Berkeley DB, configuring it so that the BDB SQL interface is also built. Be aware that it is not built by default. Instead, you need to tell the Berkeley DB configure script to also build the BDB SQL interface. For instructions on building the BDB SQL interface, see Building the DB SQL Interface in the *Berkeley DB Installation and Build Guide*.

The library and application names used when building the BDB SQL interface are different than those used by SQLite. If you want library and command shell names that are consistent with the names used by SQLite, configure the BDB SQL interface build using the compatibility (`--enable-sql_compat`) option.

Warning

The compatibility option can break other applications on your platform that rely on standard SQLite. This is especially true of Mac OS X, which uses standard SQLite for a number of default applications.

Use the compatibility option only if you know exactly what you are doing.

Unless you built the BDB SQL interface with the compatibility option, libraries and a command line shell are built with the following names:

- `dbsql`

This is the command line shell. It operates identically to the SQLite `sqlite3` shell.

- `libdb_sql`

This is the library that provides the BDB SQL interface. It is the equivalent of the SQLite `libsqlite3` library.

The Journal Directory

When you create a database using the BDB SQL interface, a directory is created alongside of it. This directory has the same name as your database file, but with a `-journal` suffix.

That is, if you create a database called "mydb" then the BDB SQL interface also creates a directory alongside of the "mydb" file called "mydb-journal".

This directory contains files that are very important for the proper functioning of the BDB SQL interface. Do not delete this directory or any of its files unless you know what you are doing.

For more information on the journal directory, see [Introduction to Environments \(page 11\)](#).

Unsupported PRAGMAs

The following PRAGMAs are not supported by the BDB SQL interface.

[PRAGMA journal_mode](#)

[PRAGMA legacy_file_format](#)

Also, [PRAGMA fullsync](#) is always on for the BDB SQL interface. (This is an issue only for Mac OS X platforms.)

Changed PRAGMAs

The following PRAGMAs are available in the BDB SQL interface, but they behave differently in some way.

PRAGMA journal_size_limit

For standard SQLite, this pragma identifies the maximum size that the journal file is allowed to be.

Berkeley DB does not have a journal file, but it does write and use *log files*. Over the course of the database's lifetime, Berkeley DB will probably create multiple log files. A new log file is created when the current log file has reached the defined maximum size for a log file.

You use `PRAGMA journal_size_limit` to define this maximum size for a log file.

For more information, see [Setting the Log File Size \(page 14\)](#).

PRAGMA max_page_count

For standard SQLite, this identifies the maximum number of pages allowed in the database. For the BDB SQL interface, this identifies the maximum size (in bytes) that the database file is allowed to be.

For both interfaces, this pragma performs essentially the same function, but you express the upper bound in a slightly different way depending on which interface you are using.

For more information, see [Configuring the Database Page Size \(page 12\)](#).

Added PRAGMAs

The following PRAGMAs are added in the Berkeley DB SQL interface.

PRAGMA TXN_BULK

`PRAGMA TXN_BULK = 0 | 1;`

Enables transactional bulk loading optimization. The default setting for TXN_BULK is 0. This means the PRAGMA is turned off. When TXN_BULK is set to 1 in the SQL source, it causes the application to enter a mode in which:

- transactional bulk loading optimization is enabled for top-level transactions
- nested per-statement transactions are not used

PRAGMA TXN_BULK enables two optimizations. The transactional bulk loading optimization uses the DB_TXN_BULK flag when starting transactions. Note that there are implications of using the DB_TXN_BULK flag, particularly with regard to its interaction with hot backup. The other optimization omits nested subtransactions for each statement. When this optimization

is enabled, you can not undo any single statement. If any statement must be undone, then the entire encompassing transaction must be aborted. This is a compromise, trading speed of bulk inserts against the standard statement undo guarantees. With these two optimizations, a transaction that inserts a large number of new records can run with much less I/O and transaction management overhead.

Miscellaneous Differences

The following miscellaneous differences also exist between the BDB SQL interface and SQLite:

- The BDB SQL interface does not support the IMMEDIATE keyword (BEGIN IMMEDIATE behaves just like BEGIN).
- When an exclusive transaction is active, it will block any new transactions from beginning (they will be blocked during their first operation until the exclusive transactions commits or aborts). Non-exclusive transactions that are active when the exclusive transaction begins will not be able to execute any more operations without being blocked until the exclusive transactions finishes.
- There are differences in how the two products work in a concurrent application that will cause the BDB SQL interface to deadlock where SQLite would result in a different error. This is because the products use different locking paradigms. See [Locking Notes \(page 7\)](#) for more information.
- The BDB SQL does not call the busy callback when a session attempts to operate the same database page that another session has locked. It blocks instead. That is to say, the functions `sqlite3_busy_handler` and `sqlite3_busy_timeout` are not effective in BDB SQL.
- The BDB SQL does not support two phase commit across databases. Attaching to multiple databases can lead to inconsistency after recovery and undetected deadlocks when accessing multiple databases from concurrent transactions in different order. Hence, applications must ensure that they access databases in the same order in any transaction that spans multiple databases. Else, a deadlock can occur that causes threads to block, and the deadlock will not be detected by Berkeley DB.
- In BDB SQL, when two sessions accessing the same database perform conflicting operations on the same page, one session will be blocked until the conflicting operations are resolved. For example,

Session 1:

```
dbsql> insert into a values (4);
dbsql> begin;
dbsql> insert into a values (5);
```

Session 2:

```
dbsql> select * from a;
```

What happens here is that Session 2 is blocked until Session 1 commits the transaction.

Session 1:

```
dbsql> commit;
```

Session 2:

```
dbsql> select * from a;  
4  
5
```

Under such situations in SQLite, operations poll instead of blocking, and a callback is used to determine whether to continue polling.

- By default, you always only have a single database file when you use BDB SQL interface SQL, just as you do when you use SQLite. However, you can configure BDB SQL interface at compile time to create one BDB SQL interface database file for each SQL table that you create. How to perform this configuration is described in the *Berkeley DB Installation and Build Guide*.

Berkeley DB Concepts

If you are a SQLite user who is migrating to the BDB SQL interface, then there are a few Berkeley DB-specific concepts you might want to know about.

- Environments. The directory that is created alongside your database file, and which ends with the "-journal" suffix, is actually a Berkeley DB environment directory. This might be interesting to you in some administrative situations. For some minimal information on what an environment is, see [Introduction to Environments \(page 11\)](#).

- The Locking Subsystem

You can configure the maximum number of locks that can be in use at any given time when you use the BDB SQL interface. This is probably only interesting to you if you are using the BDB SQL interface in a concurrent application that is running a very large number of transactions.

For information on configuring your locking subsystem, see [Managing the Locking Subsystem \(page 16\)](#).

- The Logging Subsystem

The BDB SQL interface maintains log files in its journal directory, and you can manage various aspects of these. For the overwhelming majority of applications, there is no need to manage this. But for the sake of completeness, this topic is described in this manual.

For more information, see [Administering Log Files \(page 14\)](#).

Encryption

The Berkeley DB SQL interface supports the SQLite Encryption Extension (SEE) to ensure security of your data. The supported encryption algorithm is AES-128 in CBC mode. For more

information on the concepts relating to BDB encryption, see the Programmer's Reference Guide.

To learn how to use the SQLite Encryption Extension (SEE), see the official [SQLite Documentation Page](#).

Note

The Berkeley DB SQL interface does not support the `sqlite3_rekey` method.

Chapter 2. Locking Notes

There are some important performance differences between the BDB SQL interface and SQLite, especially in a concurrent environment. This chapter gives you enough information about how the BDB SQL interface uses its database, as opposed to how SQLite uses its database, in order for you to understand the difference between the two interfaces. It then gives you some advice on how to best approach working with the BDB SQL interface in a multi-threaded environment.

If you are an existing user of SQLite, and you care about improving your application performance when using the BDB SQL interface in a concurrent situation, you should read this chapter. Existing users of Berkeley DB may also find some interesting information in this chapter, although it is mostly geared towards SQLite users.

Internal Database Usage

The BDB SQL interface and SQLite do different things when it comes to locking data in their databases. In order to provide ACID transactions, both products must prevent concurrent access during write operations. Further, both products prevent concurrent access by obtaining software level locks that allow only the current holder of the lock to perform write access to the locked data.

The difference between the two is that when SQLite requires a lock (such as when a transaction is underway), it locks the entire database and all tables. (This is known as *database level locking*.) The BDB SQL interface, on the other hand, only locks the portion of the table being operated on within the current transactional context (this is known as *page level locking*). In most situations, this allows applications using the BDB SQL interface to operate concurrently and so have better read/write throughput than applications using SQLite. This is because there is less lock contention.

By default, one Berkeley DB logical database is created within the single database file for every SQL table that you create. Within each such logical database, each table row is represented as a Berkeley DB key/data pair.

This is important because the BDB SQL interface uses Berkeley DB's Transaction Data Store product. This means that Berkeley DB does not have to lock an entire database (all the tables within a database file) when it acquires a lock. Instead, it locks a single Berkeley DB database page (which usually contains a small sub-set of rows within a single table).

The size of database pages will differ from platform to platform (you can also manually configure this), but usually a database page can hold multiple key/data pairs; that is, multiple rows from a SQL table. Exactly how many table rows fit on a database page depends on the size of your page and the size of your table rows.

If you have an exceptionally small table, it is possible for the entire table to fit on a single database page. In this case, Berkeley DB is in essence forced to serialize access to the entire table when it requires a lock for it.

Note, however, that the case of a single table fitting on a single database page is very rare, and it in fact represents the abnormal case. Normally tables span multiple pages and so

Berkeley DB will lock only portions of your tables. This locking behavior is automatic and transparent to your application.

Lock Handling

There is a difference in how applications written for the BDB SQL interface handle deadlocks as opposed to how deadlocks are handled for SQLite applications. For the SQLite developer, the following information is a necessary review in order to understand how the BDB SQL interface behaves differently.

From a usage point of view, the BDB SQL interface behaves in the same way as SQLite in shared cache mode. The implications of this are explained below.

SQLite Lock Usage

As mentioned previously in this chapter, SQLite locks the entire database while performing a transaction. It also has a locking model that is different from the BDB SQL interface, one that supports multiple readers, but only a single writer. In SQLite, transactions can start as follows:

- `BEGIN`

Begins the transaction, locking the entire database for reading. Use this if you only want to read from the database.

- `BEGIN IMMEDIATE`

Begins the transaction, acquiring a "modify" lock. This is also known as a `RESERVED` lock. Use this if you are modifying the database (that is, performing `INSERT`, `UPDATE`, or `DELETE`). `RESERVED` locks and read locks can co-exist.

- `BEGIN EXCLUSIVE`

Begins the transaction, acquiring a write lock. Transactions begun this way will be written to the disk upon commit. No other lock can co-exist with an exclusive lock.

The last two statements are a kind of a contract. If you can get them to complete (that is, not return `SQLITE_LOCKED`), then you can start modifying the database (that is, change data in the in-memory cache), and you will eventually be able to commit (write) your modifications to the database.

In order to avoid deadlocks in SQLite, programmers who want to modify a SQLite database start the transaction with `BEGIN IMMEDIATE`. If the transaction cannot acquire the necessary locks, it will fail, returning `SQLITE_BUSY`. At that point, the transaction falls back to an unlocked state whereby it holds no locks against the database. This means that any existing transactions in a `RESERVED` state can safely wait for the necessary `EXCLUSIVE` lock in order to finally write their modifications from the in-memory cache to the on-disk database.

The important point here is that so long as the programmer uses these locks correctly, he can assume that he can proceed with his work without encountering a deadlock. (Assuming that all database readers and writers are also using these locks correctly.)

Lock Usage with the DB SQL Interface

When you use the BDB SQL interface, you can begin your transaction with `BEGIN` or `BEGIN EXCLUSIVE`.

Note that the `IMMEDIATE` keyword is ignored in the BDB SQL interface (`BEGIN IMMEDIATE` behaves like `BEGIN`).

When you begin your transaction with `BEGIN`, Berkeley DB decides what kind of a lock you need based on what you are doing to the database. If you perform an action that is read-only, it acquires a read lock. If you perform a write action, it acquires a write lock.

Also, the BDB SQL interface supports multiple readers *and* multiple writers. This means that multiple transactions can acquire locks as long as they are not trying to modify the same page. For example:

Session 1:

```
dbsql> create table a(x int);
dbsql> begin;
dbsql> insert into a values (1);
dbsql> commit;
```

Session 2:

```
dbsql> create table b(x int);
dbsql> begin;
dbsql> insert into b values (1);
dbsql> commit;
```

Because these two sessions are operating on different pages in the Berkeley DB cache, this example will work. If you tried this with SQLite, you could not start the second transaction until the first had completed.

However, if you do this using the BDB SQL interface:

Session 1:

```
dbsql> begin;
dbsql> insert into a values (2);
```

Session 2:

```
dbsql> begin;
dbsql> insert into a values (2);
```

The second session blocks until the first session commits the transaction. Again, this is because both sessions are operating on the same database page(s). However, if you simultaneously attempt to write pages in reverse order, you can deadlock. For example:

Session 1:

```
dbsql> begin;
```

```
dbsql> insert into a values (3);  
dbsql> insert into b values (3);
```

Session 2:

```
dbsql> begin;  
dbsql> insert into b values (3);  
dbsql> insert into a values (3);  
Error: database table is locked
```

What happens here is that Session 1 is blocked waiting for a lock on table b, while Session 2 is blocked waiting for a lock on table a. The application can make no forward progress, and so it is deadlocked.

When such a deadlock is detected one session loses the lock it got when executing its last statement, and that statement is automatically rolled back. The rest of the statements in the session will still be valid, and you can continue to execute statements in that session. The session that does not lose its lock to deadlock detection will continue to execute as if nothing happened.

Assume Session 2 was sacrificed to deadlock detection, no value would be inserted into a and an error will be returned. But the insertion of value 3 into b would still be valid. Session 1 would continue to wait while inserting into table b until Session 2 either commits or aborts, thus freeing the lock it has on table b.

When you begin your transaction with `BEGIN EXCLUSIVE`, the session is never aborted due to deadlock or lock contention with another transaction. Non-exclusive transactions are allowed to execute concurrently with the exclusive transaction, but the non-exclusive transactions will have their locks released if deadlock with the exclusive transaction occurs. If two or more exclusive transactions are running at the same time, they will be forced to execute in serial.

If Session 1 was using an exclusive transaction, then Session 2 would lose its locks when deadlock is detected between the two. If both Session 1 and Session 2 start an exclusive transaction, then the last one to start the exclusive transaction would be blocked after executing `BEGIN EXCLUSIVE` until the first one is committed or aborted.

Chapter 3. Configuring the Berkeley DB SQL interface

In almost all cases, there is no need for you to directly configure Berkeley DB resources; instead, you can use the same configuration techniques that you always use for SQLite. The Berkeley DB SQL interface will take care of the rest.

However, there are a few configuration activities that some unusually large or busy installations might need to make and for which there is no SQLite equivalent. This chapter describes those activities.

Introduction to Environments

Before continuing with this section, it is necessary for you to have a high-level understanding of Berkeley DB's environments.

In order to manage its resources (data, shared cache, locks, and transaction logs), Berkeley DB often uses a directory that is called the *Berkeley DB environment*. As used with the BDB SQL interface, environments contain log files and the information required to implement a shared cache and fine-grained locking. This environment is placed in a directory that appears on the surface to be a SQLite rollback file.

That is, if you use BDB SQL interface to create a database called `mydb.db`, then a directory is created alongside of it called `mydb.db-journal`. Normally, SQLite creates a journal file only when a transaction is underway, and deletes this file once the transaction is committed or rolled back. However, that is not what is happening here. The BDB SQL interface journal directory contains important Berkeley DB environment information that is meant to persist between transactions and even between process runtimes. So it is very important that you do *not* delete the contents of your Berkeley DB journal directory. Doing so will cause improper operation and could lead to data loss.

Note that the environment directory is also where you put your `DB_CONFIG` file. This file can be used to configure additional tuning parameters of Berkeley DB, if its default behavior is not appropriate for your application. For more information on the `DB_CONFIG` file, see the next section.

Note

Experienced users of Berkeley DB should be aware that neither `DB_USE_ENVIRON` nor `DB_USE_ENVIRON_ROOT` are specified to `DB_ENV->open()`. As a result, the `DB_HOME` environment variable is ignored. This means that the BDB SQL interface will always create a database in the location defined by the database name given to the BDB SQL interface.

The DB_CONFIG File

You can configure most aspects of your Berkeley DB environment by using the `DB_CONFIG` file. This file must be placed in your environment directory. When using the BDB SQL interface, this

is the directory created alongside of your database. It has the same name as your database, followed by a `-journal` extension. For example, if your database is named `mydb.db`, then your environment directory is created next to the `mydb.db` file, and it is called `mydb.db-journal`.

If a `DB_CONFIG` file exists in your environment directory, it will be read for lines of the format **NAME VALUE** when your environment is opened. This happens when your application starts up and creates its first connection to the database.

One or more whitespace characters are used to delimit the two parts of the line, and trailing whitespace characters are discarded. All empty lines or lines whose first character is a whitespace or hash (`#`) character are ignored. Each line must specify both the NAME and the VALUE of the pair. The specific NAME VALUE pairs you can use with the BDB SQL interface are documented in the [Berkeley DB C API](#).

In some cases, you must either specify a configuration option before the environment is created, or the environment must be re-created before the configuration option will take effect. The documentation for each configuration option will indicate where this is true.

Creating the DB_CONFIG File Before Environment Creation

In order to provide the `DB_CONFIG` file before the environment is first created, physically make the environment directory in the correct location in your filesystem (this is wherever you want to place your database file), and put the `DB_CONFIG` file there before you create your database.

Re-creating the Environment

Some `DB_CONFIG` parameters require you to re-create your environment before they take effect. The `DB_CONFIG` parameter descriptions indicate where this is the case.

To re-create your environment:

- Make sure the `DB_CONFIG` file contains the following line:

```
add_data_dir ..
```

(This line should already be in the `DB_CONFIG` file.)

- Run the `db_recover` command line utility. If you run it from within your environment (`-journal`) directory, no command line arguments are required. If you run it from outside your environment directory, use the `-h` parameter to identify the location of the environment:

```
db_recover -h /some/path/to/mydb.db-journal
```

Configuring the Database Page Size

When using the BDB SQL interface, you configure your database page size in exactly the same way as you do when using SQLite. That is, use `PRAGMA page_size` to report and set the page

size. This PRAGMA must be called before you create your first SQLite table. See the [PRAGMA page_size](#) documentation for more information.

When you use PRAGMA `cache_size` to size your in-memory cache, you provide the cache size in terms of a number of pages. Therefore, your database page size influences how large your cache is, and so determines how much of your database will fit into memory. If you adjust the database page size, you may also want to adjust the in-memory cache size, as described in [Configuring the In-Memory Cache \(page 13\)](#).

The size of your pages can also affect how efficient your application is at performing disk I/O. It will also determine just how fine-grained the fine-grained locking actually is. This is because Berkeley DB locks database pages when it acquires a lock.

Selecting the Page Size

Note that the default value for your page size is probably correct for the physical hardware that you are using. In almost all situations, the default page size value will give your application the best possible I/O performance. For this reason, tuning the page size should rarely, if ever, be attempted.

That said, when using the BDB SQL interface, the page size affects how much of your tables are locked when read and/or write locks are acquired. (See [Internal Database Usage \(page 7\)](#) for more information.) Increasing your page size will typically improve the bandwidth you get accessing the disk, but it also may increase contention if too many key data pairs are on the same page. Decreasing your page size frequently improves concurrency, but may increase the number of locks you need to acquire and may decrease your disk bandwidth.

When changing your page size, make sure the value you select is a power of 2 that is greater than 512 and less than or equal to 64KB. (Note that the standard SQLite MAX_PAGE_SIZE limit is not examined for this upper bound.)

Beyond that, there are some additional things that you need to consider when selecting your page size. For a thorough treatment of selecting your page size, see the section on Selecting a page size in the *Berkeley DB Programmer's Reference Guide*.

Selecting the Database File Size

Berkeley DB sets an upper bound on how large your database file size is allowed to be. Any attempt to insert data into the database that grows this file beyond this upper bound results in a failure.

You can set the upper bound for your database file size using PRAGMA `max_page_count`. Issue this PRAGMA with no value to see what the current maximum database file is.

Configuring the In-Memory Cache

SQLite provides an in-memory cache which you size according to the maximum number of database pages that you want to hold in memory at any given time.

Berkeley DB also provides an in-memory cache that performs the same function as SQLite. You can configure this cache using the exact same PRAGMAs as you are used to using with

SQLite. See [PRAGMA cache_size](#) and [PRAGMA default_cache_size](#) for details. As is the case with SQLite, you use these PRAGMAs to describe the total number of pages that you want in the cache.

Note that you can change the cache size only if no table operations have been executed on the database. In other words, to change your cache size:

- Open a handle to your database.
- Execute `PRAGMA cache_size`
- Proceed with any table modification operations (CREATE, UPDATE, INSERT, SELECT) that you might want to perform.

Alternatively, you can set your cache size with your `DB_CONFIG` file, and so skip the necessity of executing the PRAGMA. See the [Berkeley DB C API](#) for details.

Administering Log Files

Your environment directory contains log files. Berkeley DB log files are used to record all the transactional activity performed against the Berkeley DB database files. This information is used after an application or system failure to automatically restore the database to an up-to-date consistent point.

Your log files are maintained by Berkeley DB's logging subsystem. There are some aspects of the Berkeley DB logging subsystem that you can configure using `DB_CONFIG` parameters, and (sometimes) by using PRAGMAs.

Note

For most users of the BDB SQL interface, there should not normally be any reason to manage your log files or otherwise worry about them. However, it is important to realize that they can not simply be deleted. Note that when using the Berkeley DB SQL interface, your log files will be automatically deleted by Berkeley DB when they are no longer needed.

The things you can manage for your logging subsystem are:

- Size of the log files. See [Setting the Log File Size \(page 14\)](#).
- Size of the logging subsystem's region. See [Configuring the Logging Region Size \(page 15\)](#).
- Size of the log buffer in memory. [Setting the In-Memory Log Buffer Size \(page 15\)](#).

Setting the Log File Size

Whenever a pre-defined amount of data is written to a log file (10 MB by default), the BDB SQL interface stops using the current log file and starts writing to a new file. You can change the maximum amount of data contained in each log file by using either `PRAGMA journal_size_limit` or the `set_lg_max` `DB_CONFIG` file parameter.

If you use `PRAGMA journal_size_limit`, then using this `PRAGMA` with no value simply returns the current journal size limit. Using:

```
PRAGMA journal_size_limit=N
```

sets the log size to *N* bytes. If the `PRAGMA` is successful, *N* is returned. If it fails, the previous log file size is returned. Failures can occur if you specify a log file size that is less than 4K bytes, or if you specify a log file size larger than the permitted file size on the system.

If you use the `DB_CONFIG` file to manage this value, `set_lg_max` may be changed without re-creating the environment. You will, however, have to restart your application in order for the `DB_CONFIG` file to be re-read.

The `DB_CONFIG` file is described in [The DB_CONFIG File \(page 11\)](#). The `set_lg_max` parameter is described in the [Berkeley DB C API](#).

For a description of how, when and why you should change the size of your log files, see the [Selecting a page size](#) section in the *Berkeley DB Programmer's Reference Guide*.

Configuring the Logging Region Size

The logging subsystem's default region size is 512 KB. The logging region is used to store database and table names, and so you may need to increase its size if you will be using a large number of tables.

You can set the size of your logging region by using the `set_lg_regionmax` `DB_CONFIG` parameter. Note that to manage this value you must set it before you create your environment, or you must re-create your environment. See [The DB_CONFIG File \(page 11\)](#) for more information.

The `set_lg_regionmax` parameter is described in the [Berkeley DB C API](#).

Setting the In-Memory Log Buffer Size

When using named (persistent) databases, log information is stored in-memory until the storage space fills up, or a transaction commit forces the log information to be flushed to disk.

It is possible to increase the amount of memory available to your file log buffer. Doing so improves throughput for long-running transactions, or for transactions that produce a large amount of data. Note that for named (persistent) databases, the default log buffer space is 32 KB.

You can increase your log buffer space by using the `set_lg_bsize` `DB_CONFIG` parameter. For the BDB SQL interface, when the logging subsystem is configured for on-disk logging, the default size of the in-memory log buffer is approximately 64KB. Note that this method can only be called before the environment is first opened, so you will have to set this by creating your `-journal` directory, and then creating your database. See [The DB_CONFIG File \(page 11\)](#) for more information.

The `set_lg_bsize` parameter is described in the [Berkeley DB C API](#).

Note

When working with in-memory databases, the environment is configured to perform logging in-memory. The log buffer is set to 64 * 1024, and the maximum log size is set to 32 * 1024.

Managing the Locking Subsystem

Whenever the BDB SQL interface reads from or writes to the database, the underlying Berkeley DB code must acquire locks. These locks represent a finite resource. For most installations, you should never have to worry about the locking resources available to Berkeley DB because the default values are appropriate for most applications.

However, if your application is using an extremely large number of threads that are all simultaneously accessing your data, then you might have to increase your locking resources. Similarly, if your database contains a very large number of tables that you are accessing using one or more simultaneous threads or processes, then you might also need to increase your locking resources.

On the other hand, if you are using the BDB SQL interface on devices with extremely limited resources, then you might want to reduce your locking resources.

All of these values must be configured before your environment is first created. To change these values after environment creation time, you must re-create the environment. See [The DB_CONFIG File \(page 11\)](#) for more information.

The maximum locking values that you can manage, and the DB_CONFIG parameter that you use to manage that value, are:

- The maximum number of lockers supported by the environment. This value is used by the environment when it is opened to estimate the amount of space that it should allocate for various internal data structures. By default, 2,000 lockers are supported.

The maximum number of lockers corresponds roughly to the maximum number of concurrent transactions in the system.

To configure this value, use the `set_lk_max_lockers` DB_CONFIG parameter. See the [Berkeley DB C API](#) for details.

- The maximum number of locks supported by the environment. By default, 10,000 locks are supported.

To configure this value, use the `set_lk_max_locks` DB_CONFIG parameter. See the [Berkeley DB C API](#) for details.

- The maximum number of locked objects supported by the environment. By default, 10,000 objects can be locked.

To configure this value, use the `set_lk_max_objects` DB_CONFIG parameter. See the [Berkeley DB C API](#) for details.

Note that when you are using the BDB SQL interface, the default values provided in the previous list are different from the default values used by Berkeley DB in general. For Berkeley DB in general, the defaults for all these values are set to 1,000.

Chapter 4. Administrating Berkeley DB SQL Databases

This chapter provides administrative procedures that are unique to the Berkeley DB SQL interface.

Backing Up Berkeley DB SQL Databases

You can use the standard SQLite `.dump` command to backup the data managed by the BDB SQL interface. You can also use the standard Berkeley DB backup mechanisms on the database.

The BDB SQL interface supports the standard SQLite Online Backup API. However, there is a small difference between the two interfaces. In the BDB SQL interface, the value returned by the `sqlite3_backup_remaining` method and the number of pages passed to the `sqlite3_backup_step` method, are estimates of the number of pages to be copied and not exact values. To be certain that the backup process is complete, check if the `sqlite3_backup_step` method has returned `SQLITE_DONE`. To learn how to use SQLite Online Backup API, see the official [SQLite Documentation Page](#).

This section describes the mechanisms that can be performed from the command line.

Offline Backups

To create an offline backup:

1. Commit or abort all on-going transactions.
2. Pause all database writes.
3. Force a checkpoint. See the `db_checkpoint` command line utility.
4. Copy your database file to the backup location. Note that in order to perform recovery from this backup, do not change the name of the database file.
5. Copy the *last* log file to your backup location. Your log files are named `log.xxxxxxxxxx`, where `xxxxxxxxxx` is a sequential number. The last log file is the file with the highest number.

Remember that your log files are placed in the environment directory, which is created on-disk next to your database file. It has the same name as your database file, but adds a `-journal` extension. For example, if your database is named `mydb.db`, then your environment directory is named `mydb.db-journal`

Hot Backup

To create a hot backup, you do not have to stop database operations. Transactions may be on-going and you can be writing to your database at the time of the backup. However, this means that you do not know exactly what the state of your database is at the time of the backup.

You can use the `db_hotbackup` command line utility to create a hot backup for you. This utility will (optionally) run a checkpoint and then copy all necessary files to a target directory. To do this when you are using the BDB SQL interface:

1. Create a `DB_CONFIG` file in your environment directory.
2. Add a `set_data_dir` parameter to the `DB_CONFIG` file. This parameter indicates what directory contains the actual Berkeley DB database managed by this environment. That directory is one level up from your environment, so you want this parameter to be:

```
set_data_dir ..
```

3. Add a `set_lg_dir` parameter to the `DB_CONFIG` file. This parameter identifies the directory that contains the environment's log files. This parameter should be:

```
set_lg_dir .
```

4. Run the `db_hotbackup` command:

```
db_hotbackup -h [environment directory] -b [target directory] -D
```

The `-D` option tells the utility to read the `DB_CONFIG` file before running the backup.

Alternatively, you can manually create a hot backup as follows:

1. Copy your database file to the backup location. Note that in order to perform recovery from this backup, do not change the name of the database file.
2. Copy all logs to your backup location.

Remember that your log files are placed in the environment directory.

Note

It is important to copy your database file *and then* your logs. In this way, you can complete or roll back any database operations that were only partially completed when you copied the database.

Incremental Backups

Once you have created a full backup (that is, either a offline or hot backup), you can create incremental backups. To do this, simply copy all of your currently existing log files to your backup location.

Incremental backups do not require you to run a checkpoint or to cease database write operations.

When you are working with incremental backups, remember that the greater the number of log files contained in your backup, the longer recovery will take. You should run full backups on some interval, and then do incremental backups on a shorter interval. How frequently you need to run a full backup is determined by the rate at which your database changes and how sensitive your application is to lengthy recoveries (should one be required).

You can also shorten recovery time by running recovery against the backup as you take each incremental backup. Running recovery as you go means that there will be less work for the BDB SQL interface to do if you should ever need to restore your environment from the backup.

About Unix Copy Utilities

If you are copying database files you must copy databases atomically, in multiples of the database page size. In other words, the reads made by the copy program must not be interleaved with writes by other threads of control, and the copy program must read the databases in multiples of the underlying database page size. Generally, this is not a problem because operating systems already make this guarantee and system utilities normally read in power-of-2 sized chunks, which are larger than the largest possible Berkeley DB database page size.

On some platforms (most notably, some releases of Solaris), the copy utility (`cp`) was implemented using the `mmap()` system call rather than the `read()` system call. Because `mmap()` did not make the same guarantee of read atomicity as did `read()`, the `cp` utility could create corrupted copies of the databases.

Also, some platforms have implementations of the `tar` utility that performs 10KB block reads by default. Even when an output block size is specified, the utility will still not read the underlying database in multiples of the specified block size. Again, the result can be a corrupted backup.

To fix these problems, use the `dd` utility instead of `cp` or `tar`. When you use `dd`, make sure you specify a block size that is equal to, or an even multiple of, your database page size. Finally, if you plan to use a system utility to copy database files, you may want to use a system call trace utility (for example, `ktrace` or `truss`) to make sure you are not using a I/O size that is smaller than your database page size. You can also use these utilities to make sure the system utility is not using a system call other than `read()`.

Recovering from a Backup

If you used standard Berkeley DB backup procedures to backup your database, then you can restore your database using the procedures described in this section.

Note that Berkeley DB supports two types of recovery:

- Normal recovery, which examines only those log records needed to bring the database to a consistent state since the last checkpoint. Normal recovery starts with any logs used by any transactions active at the time of the last checkpoint, and examines all logs from then to the current logs.

Normal recovery is automatically run (if necessary) when you open your database. It is necessary to run recovery if a thread or process shuts down without properly closing the database.

- Catastrophic recovery examines all available log files. You use catastrophic recovery to restore your database from a previously created backup.

Catastrophic Recovery

Use catastrophic recovery when you are recovering your database from a previously created backup. Note that to restore your database from a previous backup, you should copy the backup to a new environment directory, and then run catastrophic recovery. Failure to do so can lead to the internal database structures being out of sync with your log files.

To run catastrophic recovery:

- Shutdown all database operations.
- Restore the backup to an empty directory. This means you need your database file, as well as the `-journal` directory, and any available log files that the backup contains.

Note that the backup database file and the journal directory must have the same name as the database and journal directory that you are restoring. You can put the backup in a different location on disk, but the name of the file and directory must remain the same.

- Make sure that a `DB_CONFIG` file exists in the journal directory that you are using to restore your database. This file must contain a the following line:

```
set_data_dir ..
```

- Run the `db_recover` command line utility with the `-c` option.

Note that catastrophic recovery examines every available log file — not just those log files created since the last checkpoint as is the case for normal recovery. For this reason, catastrophic recovery is likely to take longer than does normal recovery.

Syncing with Oracle Databases

Oracle's SQLite Mobile Client product allows you to synchronize a SQLite database with a back-end Oracle database. Because the BDB SQL interface is a drop-in replacement for SQLite, this means you can synchronize a Berkeley DB database with an Oracle back-end as well.

Note

Berkeley DB SQL databases are not compatible with SQLite databases. In order for sync to work, you must remove any currently existing SQLite databases.

Syncing on Unix Platforms

For Unix platforms, the easiest way to use Oracle's SQLite Mobile Client is to build the BDB SQL interface with the compatibility option. That is, specify both `--enable-sql` and `--enable-sql-compat` when you configure your Berkeley DB installation. This causes libraries with the exact same name as the SQLite libraries to be created when you build Berkeley DB.

Having done that, you must then change your platform's library search path so that it finds the Berkeley DB libraries *before* any installed SQLite libraries. On many (but not all) Unix platforms, you do this by modifying the `LD_LIBRARY_PATH` environment variable. See your operating system documentation for information on how to change your search path for dynamically linked libraries.

Once you have properly configured and built your Berkeley DB installation, and you have properly configured your operating system, you can use the Oracle SQLite Mobile Client in exactly the same way as you would if you were using standard SQLite libraries and databases with it. See the [Oracle Database Lite](#) documentation for information on using SQLite Mobile Client.

For information on building the BDB SQL interface, see the Configuring the SQL Interface section in the *Berkeley DB Installation and Build Guide*.

Syncing on Windows Platforms

For Windows platforms, you use Oracle's SQLite Mobile Client by building the BDB SQL interface in the same way as you normally do. See the Building Berkeley DB for Windows chapter in the *Berkeley DB Installation and Build Guide* for more information.

Once you have built the product, rename the Berkeley DB SQL dlls so that they are named identically to the standard SQLite dlls (sqlite3.dll). Install the renamed Berkeley DB SQL dll along with the main Berkeley DB dll (libdb5x.dll) in the same directory as the SQLite dlls. See the Building the SQL API section for details.

Finally, configure your Windows PATH environment variable so that it finds your Berkeley DB dlls before it finds any standard SQLite dlls that might be installed on your system.

Once you have built your Berkeley DB installation and renamed your dlls, and you have properly configured your operating system, you can use the Oracle SQLite Mobile Client in exactly the same way as you would if you were using standard SQLite libraries and databases with it. See the [Oracle Database Lite](#) documentation for information on using SQLite Mobile Client.

Syncing on Windows Mobile Platforms

For Windows Mobile platforms, you use Oracle's SQLite Mobile Client by building the BDB SQL interface in the same way as you normally do. See the Building Berkeley DB for Windows Mobile chapter in the *Berkeley DB Installation and Build Guide* for more information.

Once you have built the product, rename the Berkeley DB SQL dll to sqlite3.dll. Then, copy the dll to the \Windows path on the phone. Note that you only need the new sqlite3.dll; you do not need any of the other Berkeley DB dlls.

Once you have built your Berkeley DB installation and renamed your dlls, and you have properly configured your operating system, you can use the Oracle SQLite Mobile Client in exactly the same way as you would if you were using standard SQLite libraries and databases with it. See the [Oracle Database Lite](#) documentation for information on using SQLite Mobile Client.

Data Migration

If you have a database created by SQLite, you can migrate it to a Berkeley DB database for use with the BDB SQL interface. For production applications, you should do this only when your application is shutdown.

Migration Using the Shells

To migrate your data from SQLite to a Berkeley DB database:

1. Make sure your application is shutdown.
2. Open the SQLite database within the **sqlite3** shell.
3. Execute the `.output` command to specify the location where you want to dump data.
4. Dump the database using the SQLite `.dump` command.
5. Close the **sqlite3** shell and open the Berkeley DB `dbsql` shell. Note that if you build the BDB SQL interface with the compatibility option, you can alternatively use Berkeley DB's `sqlite3` utility.
6. Load the dumped data using the `.read` command.

Note that you can migrate in the reverse direction as well. Dump the Berkeley DB database by calling `.dump` from within the `dbsql` shell, and load it into SQLite by calling `.read` from within SQLite's **sqlite3** shell.

Supported Data and Schema

You can migrate data between SQLite and Berkeley DB that uses the UTF-8 character encoding.

The following data types can be migrated between SQLite and Berkeley DB:

- CHAR, TEXT , VARCHAR, NVARCHAR, STRING
- REAL, DOUBLE, FLOAT
- INTEGER, BOOLEAN, BIG INTEGER, NUMBER
- NUMERIC
- BLOB, CLOB
- NULL, NOT NULL
- COLLATE BINARY, COLLATE RTRIM, COLLATE NOCASE
- DATETIME, CURRENT_TIME, CURRENT_DATE, CURRENT_TIMESTAMP

The following schema can be migrated between SQLite and Berkeley DB:

- PRAGMA writable_schema=ON/OFF
- PRAGMA foreign_keys=ON/OFF
- PRAGMA cache_size
- CREATE TABLE
 - PRIMARY KEY

- UNIQUE
- CONFLICT IGNORE, FAIL, REPLACE, ABORT, ROLLBACK
- REFERENCE ON ... CASCADE, ON ... NO ACTION, DEFERRABLE INITIALLY DEFERRED, and so forth.
- AUTOINCREMENT
- Static DEFAULT value, dynamic DEFAULT value
- Functions such as datetime, typeof, and so forth.
- ASC, DESC
- HIDDEN
- CHECK
- CREATE INDEX, UNIQUE INDEX
- CREATE VIEW
 - SELECT statement, ANALYZE
 - JOIN
 - UNION
- CREATE TRIGGER AFTER/BEFORE BEGIN
- CREATE VIRTUAL TABLE USING
- INSERT

Replicating Berkeley DB SQL Databases

It is possible to replicate a Berkeley DB SQL database using the `db_replicate` utility. This section outlines the configuration steps required to enable replication, and provides an example of how to use the utility with a Berkeley DB SQL database.

Note

The following section assumes that you are familiar with Berkeley DB replication concepts. If not, you should read the Running Replication using the `db_replicate` Utility section in the *Berkeley DB Programmer's Reference Guide*.

Preparing to use Replication with the Berkeley DB SQL API

In order to use replication with a Berkeley DB SQL application, you must do a few things prior to creating a database.

1. Add the following lines to the DB_CONFIG file:

```
set_open_flags db_init_rep
add_data_dir ../
set_open_flags db_thread
set_open_flags db_register
```

This causes the Berkeley DB SQL engine to create a database that is compatible with replication.

See [The DB_CONFIG File \(page 11\)](#) for more information about that file.

2. Tell Berkeley DB how to startup the local (current) site. Do this by adding a line to the DB_CONFIG file of the form:

```
repmgr_set_local_site <address> <port>
```

where <address> is the URL for the local machine and <port> is the port used on this machine for replication communications.

For more information, on configuring the local site, see the DB_ENV->repmgr_set_local_site() method.

3. Tell Berkeley DB about the other machines that are participating in your replication group. Do this by adding a line to the DB_CONFIG file of the form:

```
repmgr_add_remote_site <address> <port>
```

for each machine participating in your replication group, other than the local machine. Here, <address> is the URL for the remote machine that you are identifying, and <port> is the port used on that machine for replication communications.

For more information, see the DB_ENV->repmgr_add_remote_site() method.

Using Replication with the Berkeley DB SQL API

Once you have performed the configuration steps described in the previous section, you can start using (and populating) your Berkeley DB SQL database as normal.

When you are ready to start replicating your Berkeley DB SQL database, do the following:

1. Each site in the replication group needs to have a starting point. The best way to create that is to use a backup of the database that will be replicated. Install a copy of the backup at each of the sites in the replication group. For detailed instructions on creating a backup, see [Backing Up Berkeley DB SQL Databases \(page 18\)](#).
2. Create or update the DB_CONFIG file at each site in the replication group as discussed in the previous section. Make sure that each site has the correct settings for DB_ENV->repmgr_set_local_site() and DB_ENV->repmgr_add_remote_site().
3. On the site that you want to start as the master of the replication group, run the db_replicate utility in the following way:

```
db_replicate -M -h <journal-path>
```

where <journal-path> is the path to the journal directory of the database that you want to replicate.

4. For each of the other sites in the group, run the db_replicate utility in the following way:

```
db_replicate -h <journal-path>
```

where <journal-path> is the path to the journal directory of the database that you want to replicate.

5. Verify that replication has started successfully. You can do this by issuing an update operation on the site you selected as the master, and then running a query on a remote site to verify that the updated contents are visible on the remote site.