

# Databases

SQLDB



- 1. Introduction
- 2. SQLDB
- 3. Replication
- 4. Recovery from failure
- 5. Hardware



# Introduction



- Operating in SA and UK
- Been using Erlang since 2001
- Products for Financial and Telecoms purposes
- Emerald, Amber and Crimson

# SQLDB





**databases**  
the PMT way

# erlang database drivers

---

- Only one “single API” to multiple db’s
- OTP **odbc** application



Powered by  
**Pattern Matched**



but many incompatible and custom drivers:

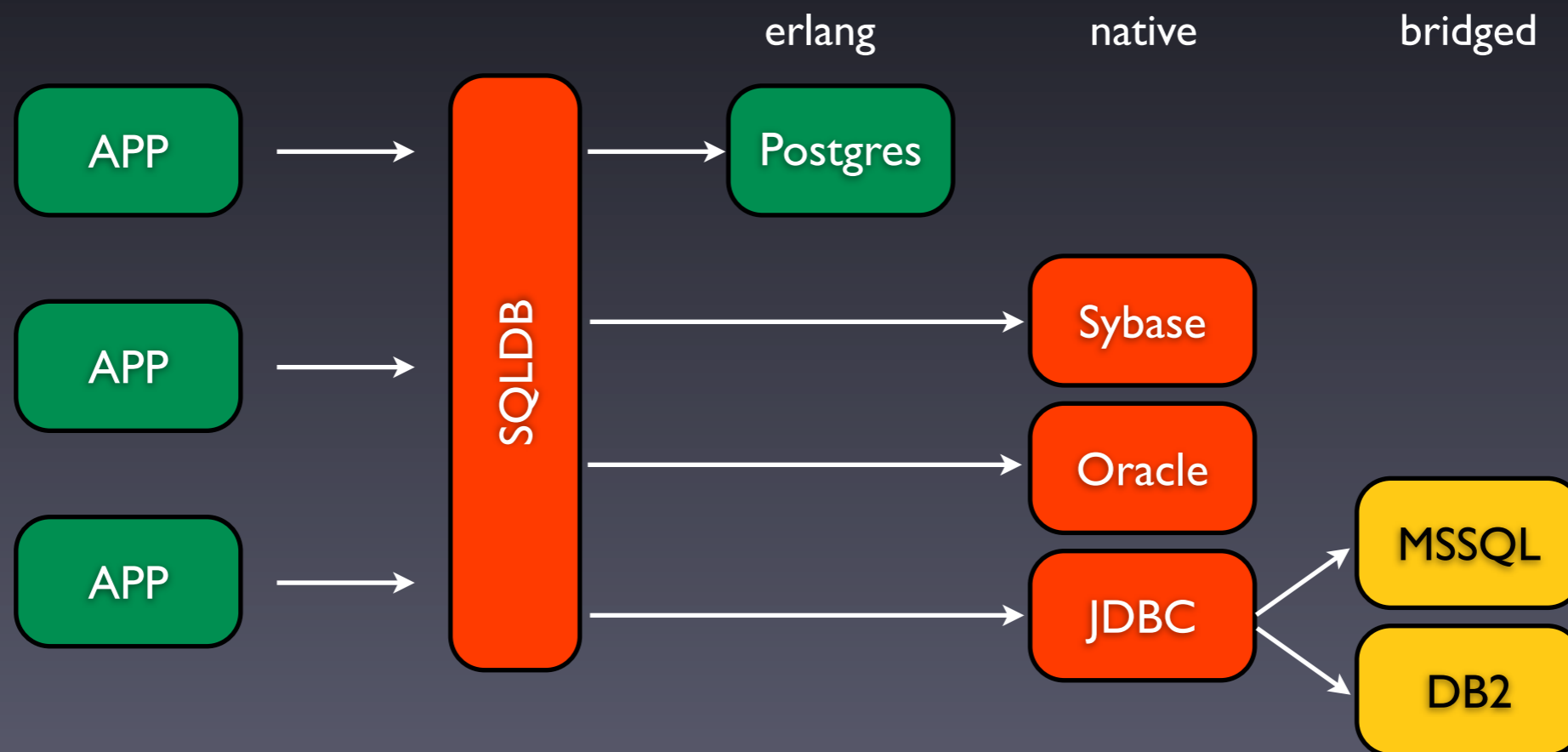
- erlang-mysql-driver
- Yxa
- pgsql
- etc



- single, database-agnostic API
- modular design
- support for connection pools

- we picked **PostgreSQL** as our preferred db
- It is **robust**
- many **features**
- **VERY stable**

- native erlang code
- in use from 2004
- minimum performance 600 TPS/connection
- Ideal for connection pool model





```
ok = sqldb_pool:acquire_connection(1"ThePool"),  
try  
  ...2  
  after  
    ok = sqldb_pool:release_connection(3)  
end.
```

- <sup>1</sup> Acquire a connection - potential to take a long time
- <sup>2</sup> Do your work
- <sup>3</sup> Release the connection





```
ok = sqldb_pool:acquire_connection("ThePool"),  
try  
  ok = sqldb_api:begin_transaction(),  
  ...  
  ok = sqldb_api:commit_transaction(),  
after  
  ok = sqldb_pool:release_connection()  
end.
```

- ① start a transaction
- ② commit or rollback the transaction





```
ok = sqldb_pool:acquire_connection("ThePool"),
try
  Fun =
    fun()1->
      sqldb_api:simple_query("select count(*) from users;"2,[])
    end,
    {ok,
      [{sqldb_result,3[{column,"count",int,8}4],
        [{row,UserCount}]5]}] =
      sqldb_api:transaction(Fun)6
    after
      ok = sqldb_pool:release_connection()
    end.
```

- <sup>1</sup> Wrap the statements in a Fun
- <sup>2</sup> Simple SQL
- <sup>3</sup> Multiple result sets
- <sup>4</sup> Column Information
- <sup>5</sup> Data returned from the query
- <sup>6</sup> Everything executed as a single transaction





```
SQL = "select * from users;",  
{ok, [{sqldb_result, Columns, Rows}] = sqldb_api:simple_query(SQL, []),  
...
```



```
Crash dump was written to: erl_crash.dump  
eheap_alloc: Cannot allocate 1824525600 bytes of memory (of type "heap").
```

- Do not accumulate terms in one process!
- Yes you can use a (server side) cursor
- or...





```
SQL = "select * from users;",  
CallbackModule = user_callback,  
CallbackState = [InitialState],  
ok = sqldb_connection:callback_query(Query, Params, CallbackModule, CallbackState),
```

- 1 The same query
- 2 Name of the callback module
- 3 Initial state
- 4 callback\_query function





# databases

the PMT way

# callback queries

```
sqldb_init(InitialState)1 ->  
{ok, InitialState}.
```

```
sqldb_result_set_begin(ColumnInfo, ModuleState)2 ->  
{ok, ModuleState}.
```

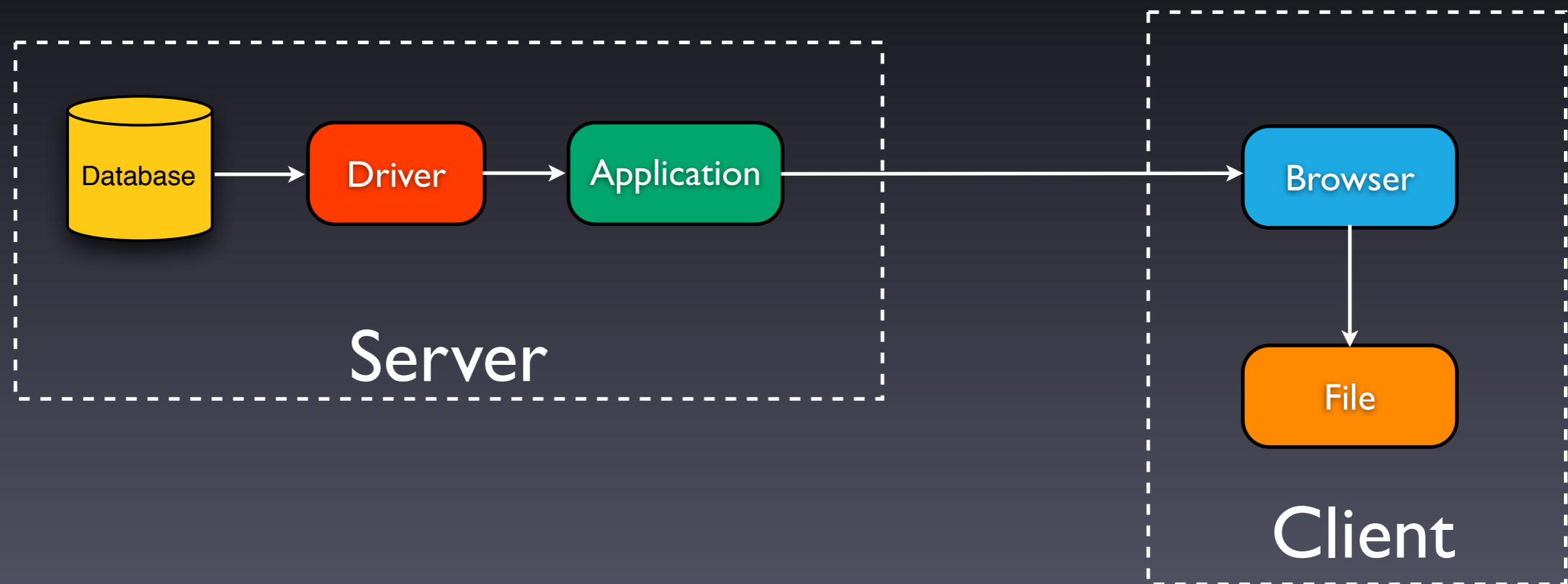
```
sqldb_row(Row, ModuleState)3 ->  
{ok, ModuleState}.
```

```
sqldb_result_set_complete(_Command, ModuleState)4 ->  
{ok, ModuleState}.
```

```
sqldb_results_complete(ModuleState)5 ->  
{ok, ok, State}.
```

- <sup>1</sup> Initialisation
- <sup>2</sup> Function called when a new result set is being processed
- <sup>3</sup> Called for every row in the result set
- <sup>4</sup> Called when a result set finishes
- <sup>5</sup> Called when the query is done





- 1 No accumulation of data anywhere
- 2 Flow control all the way from the client's file to the database
- 3 Very low latency - data is immediately available
- 4 Very efficient - because we have concurrency for free



- 1 Our goal is to make SQLDB the single Erlang database API
- 2 We will soon make the source available... but
- 3 We still need drivers for the native erlang odbc application
- 4 And help to make it perfect...
- 5 Contact [rvg@patternmatched.com](mailto:rvg@patternmatched.com) if you are interested to help!

# Replication





## Mnesia

- ① Fast Reads and writes
- ② Distributed Transactions
- ③ Easy!

- ① Indexes are terrible
- ② Never store “bulk data”
- ③ Practical size limits
- ④ No real data integrity (hacks & cheats)

Yes

## SQL

- ① Huge capacity
- ② Fast queries
- ③ Complex joins
- ④ Strict Integrity if you use it

- ① Hard to make redundant
- ② Very hard to scale - you use hardware
- ③ Hard to replicate

No



- we use both
- sometimes Mnesia leads, sometimes SQL does
- we maintain a single, abstract model for both databases



```
sql_model() ->
[#sqldb_table{tablename=vcr_voucher_pin,
  fields=[#sqldb_field{fieldname=voucher_id,      type=int8,
    #sqldb_field{fieldname=seed,                  type=int4},
    #sqldb_field{fieldname=hash,                  type=bytea},
    #sqldb_field{fieldname=crypto_format,          type=int4},
    #sqldb_field{fieldname=key_index,              type=int4},
    #sqldb_field{fieldname=kcv,                   type=char,length=6}],
  pk = [voucher_id],
  unique = [ #sqldb_unique{ unique_id = 1,
    fields=[hash] }],
  fk
    = [#sqldb_fk{fk_id = 1,
      on_delete_cascade = true,
      fields = [voucher_id],
      ftable = vcr_voucher_base,
      fkeys = [voucher_id]},
      #sqldb_fk{fk_id = 2,
        on_delete_cascade = true,
        ftable = cem_crypto_format,
        fields = [crypto_format],
        fkeys = [crypto_format]
      }]]
].
```

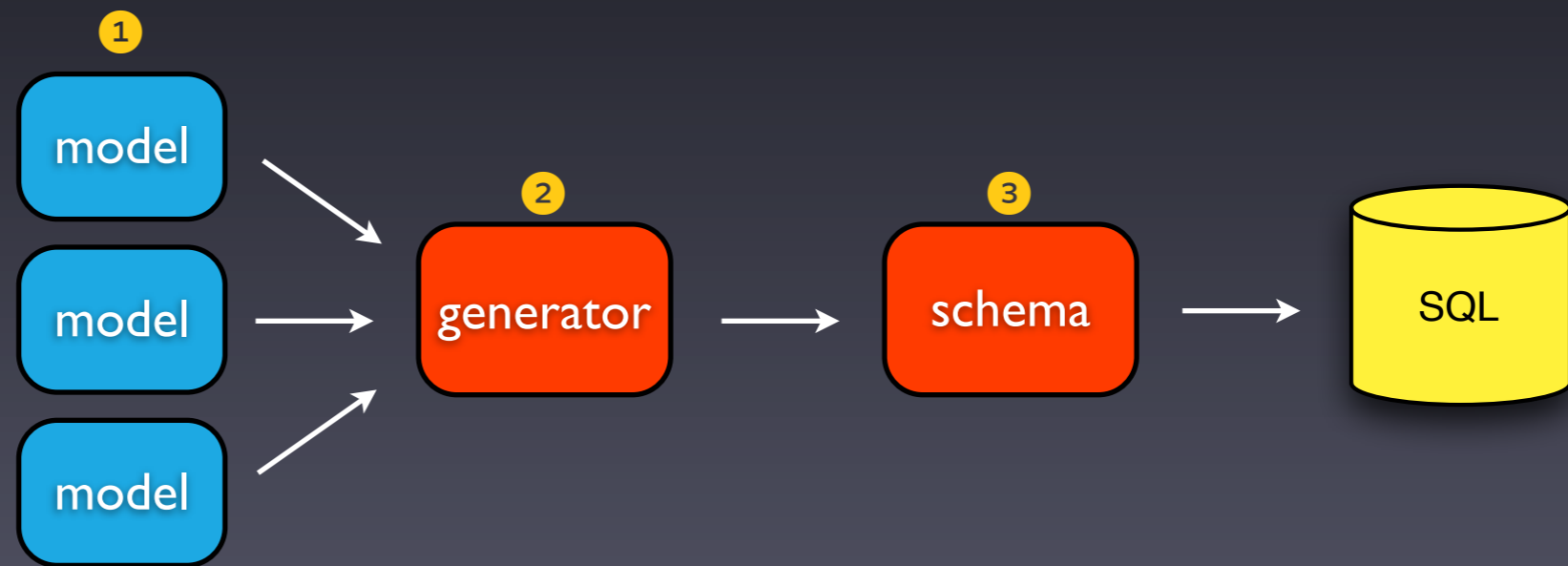




```
-record(user,{username,
               company_id      = 1,
               accounting_entity_id,
               name             = "",
               surname          = "",
               id_nr            = "",
               contact_nr       = "",
               creation_date    = "",
               msisdn           = undefined,
               address1         = "",
               address2         = "",
               address3         = "",
               city             = "",
               code             = "",
               email_address    = "",
               credential_type  = 1
            }).

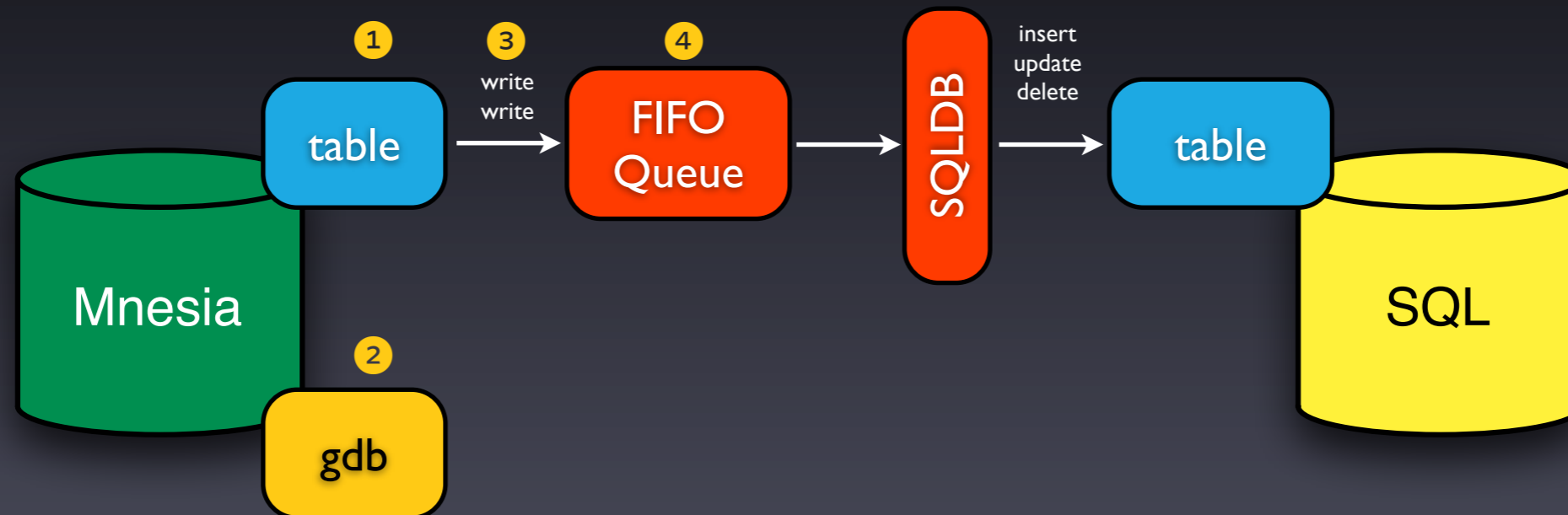
...
table_def() ->
    {emr_user,fields(),[replication,{replication_callback,asr_replication}]}.
...
```





- ① Use the abstract model for all the objects to
- ② generate a database specific
- ③ schema

```
CREATE TABLE vcr_voucher_pin (  
    voucher_id bigint NOT NULL,  
    seed integer,  
    encrypted_pin bytea,  
    hash bytea,  
    crypto_format integer,  
    key_index integer,  
    kcv character(6)  
);
```



- 1 table configured for replication
- 2 gdb hooks into the mnesia activity mechanism
- 3 fast module that intercepts writes and queue events in same transaction
- 4 strict fifo queue that converts events into SQL and push them out via SQLDB





- Both SQL and mnesia **schemas** are GENERATED
- it is simple to **replicate** to SQL
- The model ensures **consistency**



# Dealing with Failure and Recovery



- complete loss of all mnesia nodes
- network partitioning
- Point-in-time recovery (For testing etc)

we solve this by design...

- all **critical data** is also in SQL
- in most cases **SQL leads...** because PostgreSQL is fast
- a **full recovery** is one command:

```
execute() ->
  ok = lock_out_users(),
  ok = stop_services(),
  ok = clear_queues(),
  ok = clear_tables(),
  ok = recover_records(),
  ok = set_counters(),
  ok = clear_queues(),
  ok = start_services(),
  ok = open_for_users().
```

copying back data is trivial...

- we use the sqldb **callback query** feature
- combined with **dirty writes**
- finish off with **full integrity test**
- the model makes this process **automatic**
- we use the **same code** for all tables

# A Quick look at hardware





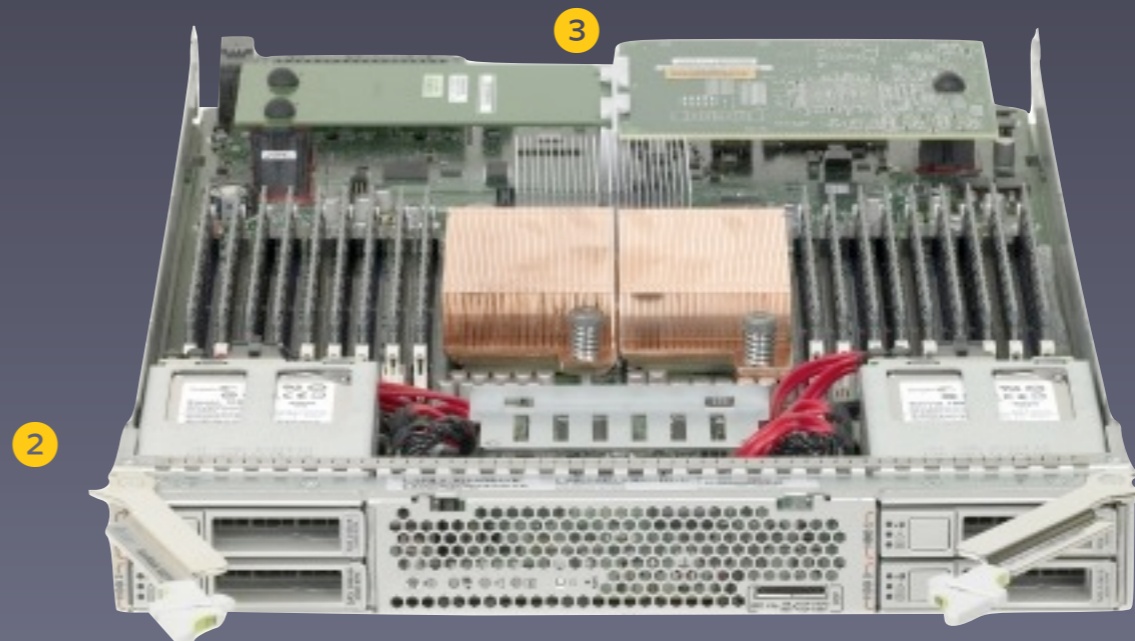
- 1 Sun Blade 6000 Chassis
- 2 Sun 6270 Blades
- 3 Fibre Channel HA

Host

Zone

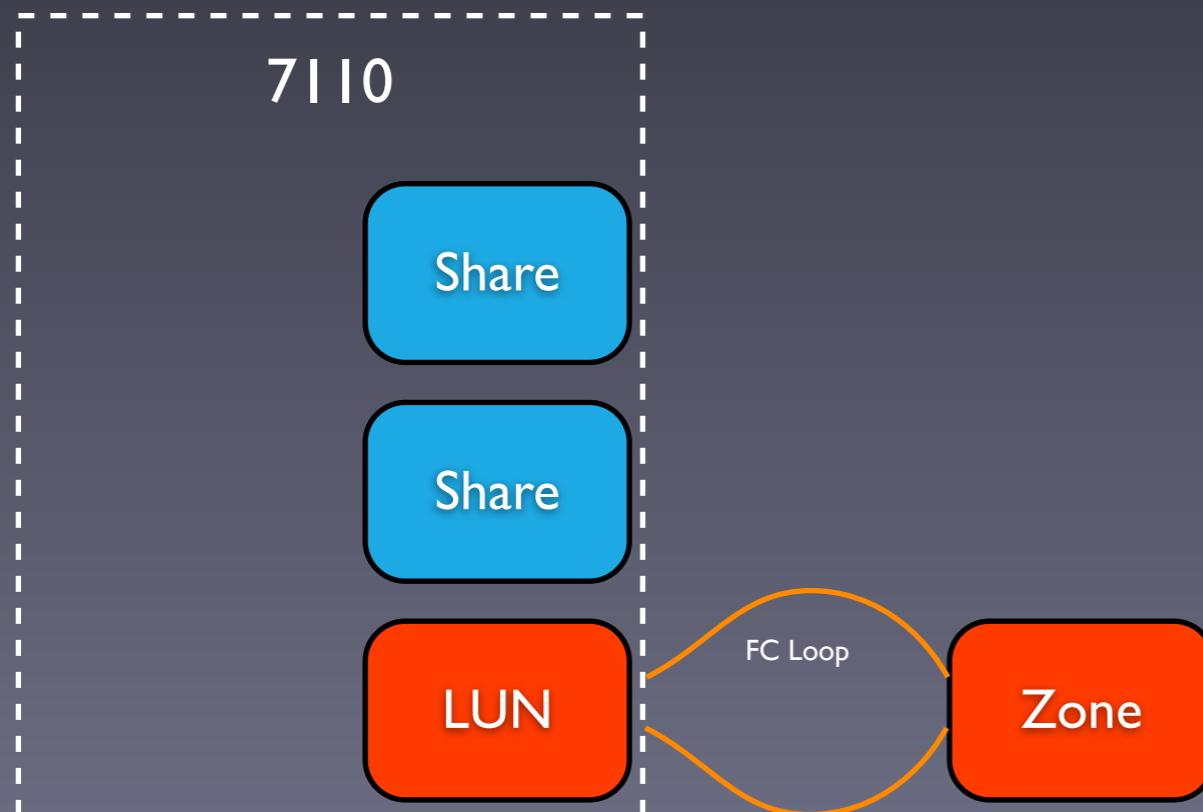
Zone

Zone





- 1 Oracle OpenStorage 7110/7310
- 2 Fibre Channel NAS
- 3 4mbps

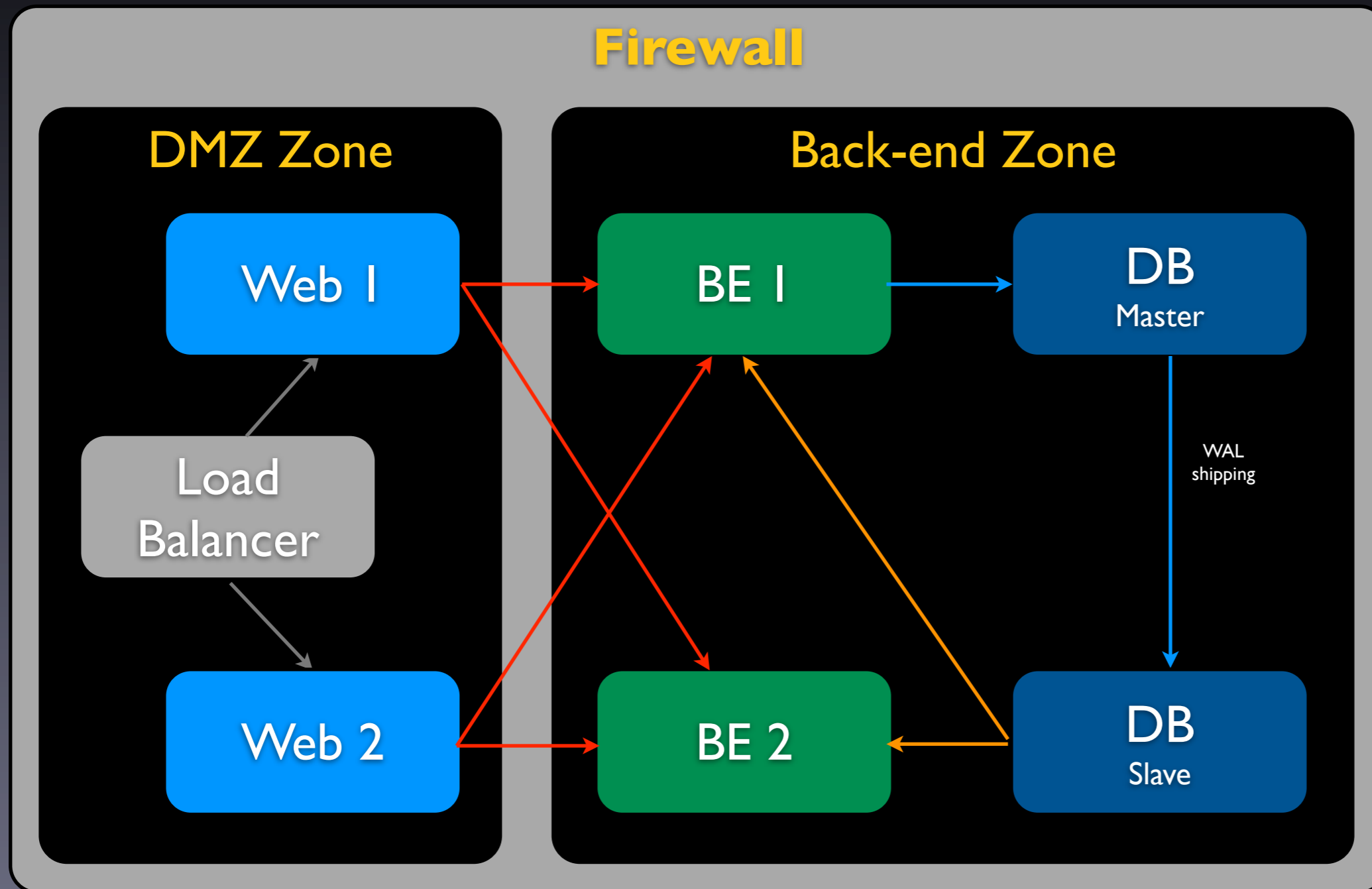


Open Storage 7110



Brocade FC Switch







- **Questions?**