

mod\_ndb:

## A REST Web Services API for MySQL Cluster

### Design goals

- Build a database server that conforms to HTTP 1.1.
- Have a lock-free design, with no mutexes in the mod\_ndb code.
- Build mod\_ndb for multiple versions of Apache, MySQL, and NDB from a single source tree.
- Do as much work as possible when processing the configuration file, and as little as possible when servicing a request.
- Be able to process configuration files without connecting to a cluster or using the NDB Data Dictionary.

# Apache processes and threads in mod\_ndb

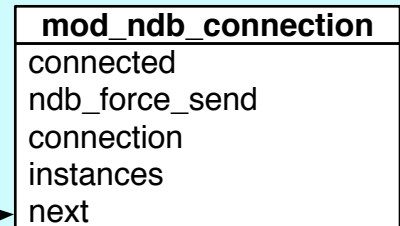
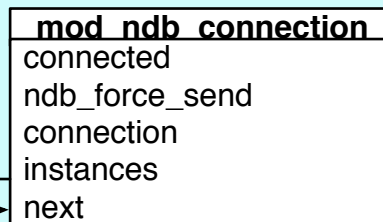
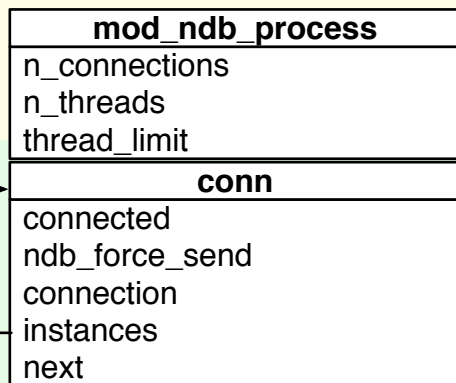
*mod\_ndb.h*

```
struct mod_ndb_process {
    int n_connections;
    int n_threads;
    int thread_limit;
    struct mod_ndb_connection conn; // not a pointer
};
```

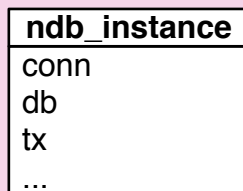
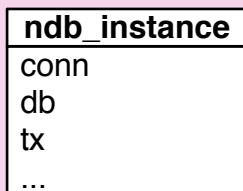
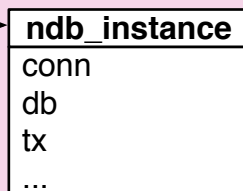
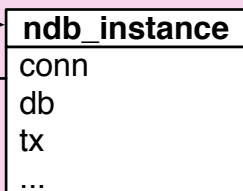
*mod\_ndb.h*

```
struct mod_ndb_connection {
    unsigned int connected;
    int ndb_force_send;
    Ndb_cluster_connection *connection;
    ndb_instance **instances;
    struct mod_ndb_connection *next;
};
typedef struct mod_ndb_connection ndb_connection;
```

One mod\_ndb\_process  
per Apache process



One mod\_ndb\_connection per NDB connect string



n\_threads

One ndb\_instance per Apache thread,  
per NDB connect string

*mod\_ndb.h*

```
struct mod_ndb_instance {
    struct mod_ndb_connection *conn;
    Ndb *db;
    NdbTransaction *tx;
    int n_read_ops;
    int max_read_ops;
    struct data_operation *data;
    struct {
        unsigned int has_blob : 1;
        unsigned int aborted : 1;
        unsigned int use_etag : 1;
    } flag;
    unsigned int requests;
    unsigned int errors;
};
```

```
typedef struct mod_ndb_instance
    ndb_instance;
```

## Some basics of query execution

- In the configuration for an endpoint, all of the "key columns" -- parameters like "id=4" and "year=2000" that may appear in the query string -- are stored in a sorted list. When the parameters are read from *r->args*, we use a binary search to find each parameter in the key columns.
- Besides named parameters, key columns can also be passed in *r->path\_info*, as in the example *http://server/ndb/mytable/2000/4*. Pathinfo configuration is stored as a mapping from the position in the *path\_info* string to the key column's index number in the sorted list – so the value gets associated with a named key column *without* having to use the binary sort.
- Once a *key\_column* is found, *set\_key()* in *Query.cc* determines how to use it. Either it to a filter, or it belongs to an index and therefore implies an access plan. If the implied plan is better than the current plan, then use it:

```
if(keycol.implied_plan > q->plan) {  
    q->plan = keycol.implied_plan;  
    q->active_index = keycol.index_id;  
}
```

- The request body – *i.e.* the data sent with a POST request – is handled differently. When the body is read (in *read\_http\_post.cc*), the names and values are stored in an apache table, *q->form\_data*. Later, in *set\_up\_write()*, we iterate over the list of updatable columns *dir->updatable* and retrieve each column's new value (if any) from *q->form\_data* using *ap\_table\_get()*.

When multipart/form-data is supported, this might change.

## Per-server (i.e. per-VHOST) config structure

config::srv
connect_string
max_read_operations

```
struct srv {
    char *connect_string;
    int max_read_operations;
};
```

## Apache per-directory config structure

config::dir
database
table
pathinfo_size
pathinfo
allow_delete
use_etags
results
sub_results
format_param[]
incr_prefetch
flag.pathinfo_always
flag.has_filters
visible
updatable
indexes
key_columns

```
/* Apache per-directory configuration */
struct dir {
    char *database;
    char *table;
    int pathinfo_size;
    short *pathinfo;
    int allow_delete;
    int use_etags;
    result_format_type results;
    result_format_type sub_results;
    char *format_param[2];
    int incr_prefetch;
    struct {
        unsigned pathinfo_always : 1;
        unsigned has_filters : 1;
    } flag;
    apache_array<char*> *visible;
    apache_array<char*> *updatable;
    apache_array<config::index> *indexes;
    apache_array<config::key_col> *key_columns;
};
```

## Configuration Directives

Directive	Function	Data Structure	Inheritable
ndb-connectstring	connectstring()	srv->connect_string	Yes
ndb-max-read-subrequests	maxreadsubrequests()	srv->max_read_operations	Yes
Database	ap_set_string_slot()	dir->database	Yes
Table	ap_set_string_slot()	dir->table	Yes
Deletes	ap_set_flag_slot()	dir->allow_delete	Yes
Format	result_format()	dir->results	Yes
Columns	non_key_column()	dir->visible	No
AllowUpdate	non_key_column()	dir->updatable	No
PrimaryKey	primary_key()	dir->key_columns	No
UniqueIndex	named_index()	dir->key_columns	No
OrderedIndex	named_index()	dir->key_columns	No
PathInfo	pathinfo()	dir->pathinfo	No
Filter	filter()	dir->key_columns	No

## Configuration: Indexes and key columns

config::index
name
type
n_columns
first_col_serial
first_col_idx

```
struct index {
    char *name;
    char type;
    unsigned short n_columns;
    short first_col_serial;
    short first_col;
};
```

config::key_col
name
index_id
serial_no
idx_map_bucket
next_in_key_serial
next_in_key
is.in_pk
is.filter
is.alias
is.in_ord_idx
is.in_hash_idx
is.in_pathinfo
filter_op
implied_plan

```
struct key_col {
    char *name;
    short index_id;
    short serial_no;
    short idx_map_bucket;
    short next_in_key_serial;
    short next_in_key;
    struct {
        unsigned int in_pk      : 1;
        unsigned int filter    : 1;
        unsigned int alias      : 1;
        unsigned int in_ord_idx : 1;
        unsigned int in_hash_idx : 1;
        unsigned int in_pathinfo : 1;
    } is;
    int filter_op;
    AccessPlan implied_plan;
};
```

```
/*
```

Every time a new column is added, the columns get reshuffled some, so we have to fix all the mappings between serial numbers and actual column id numbers.

The configuration API in Apache never gives the module a chance to "finalize" a configuration structure. You never know when you're finished with a particular directory. So, we run `fix_all_columns()` every time we create a new column, which, alas, does not scale too well.

While processing the config file, the CPU time spent fixing columns grows with  $n^2$ , the square of the number of columns. This could be improved using config handling that was more complex (a container directive) or less user-friendly (an explicit "end" token).

On the other hand, the design is optimized for handling queries at runtime, where some operations (e.g. following the list of columns that belong to an index) are constant, and the worst (looking up a column name in the columns table) grows at  $\log n$ .

```
*/
```

## N-SQL

The N-SQL language is built using the Coco/R C++ compiler generator from <http://www.ssw.uni-linz.ac.at/coco/> -- all basic configuration in the parser is implemented by calls in to the older configuration routines in *config.cc*

## Using C++ class templates above the Apache API

Apache's C-language API relies heavily on void pointers that you can cast to different data types. In C++, though, casting is no fun – the compiler requires you to make every cast explicitly, and casting defeats the type-safe design of the language.

Here are some examples from the array API: `array_header->elts` is a `char *` which you cast to an array pointer, and `ap_push_array()` returns a void pointer to a new element.

*httpd/ap\_alloc.h*

```
typedef struct {
    ap_pool *pool;
    int elt_size;
    int nelts;
    int nalloc;
    char *elts;
} array_header;

array_header * ap_make_array(pool *p, int nelts, int elt_size);
void * ap_push_array(array_header *);
```

---

```
template <class T>
class apache_array: public array_header {
public:
    int size() { return this->nelts; }
    T **handle() { return (T**) &(this->elts); }
    T *items() { return (T*) this->elts; }
    T &item(int n){ return ((T*) this->elts)[n]; }
    T *new_item() { return (T*) ap_push_array(this); }
    void * operator new(size_t, ap_pool *p, int n) {
        return ap_make_array(p, n, sizeof(T));
    };
};
```

*mod\_ndb.h*

In `mod_ndb`, the template `apache_array<T>` builds a subclass of `array_header` to manage an array of any type. All of the casting is done here in the template definition, so the code in the actual source files is cleaner:

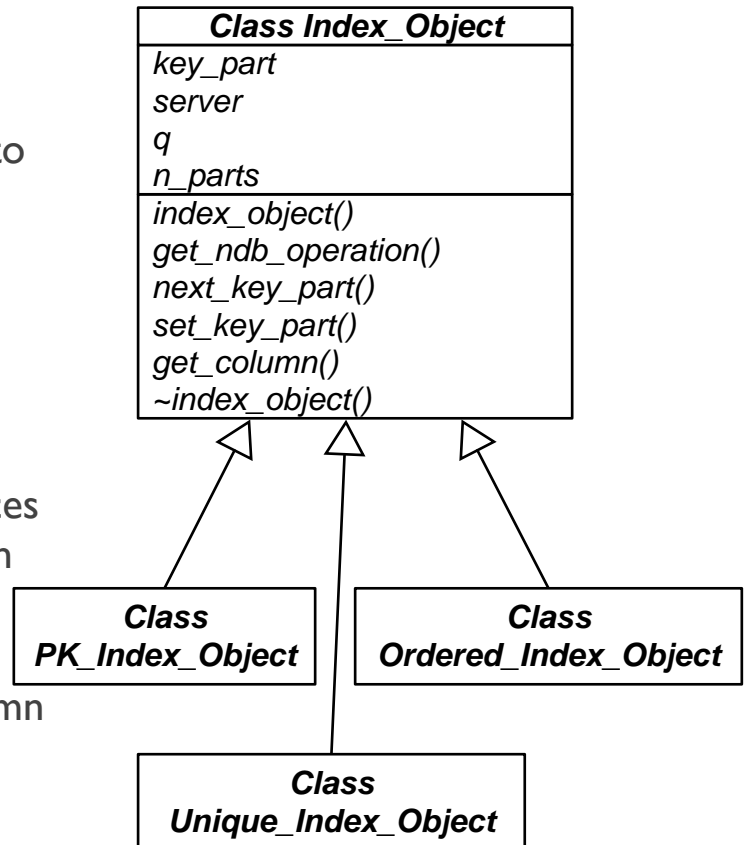
```
dir->visible      = new(p, 4) apache_array<char *>;
dir->updatable    = new(p, 4) apache_array<char *>;
dir->indexes       = new(p, 2) apache_array<config::index>;

*dir->visible->new_item() = ap_pstrdup(cmd->pool, arg);
```

## Class index\_object: Standardizing index access in mod\_ndb

The index\_object class hierarchy is defined and implemented entirely in the file "index\_object.h"

- get\_ndb\_operation() is a single interface to getNdbOperation, getNdbIndexOperation, and getNdbIndexScanOperation.
- set\_key\_part() is a single interface for op->equal() and scanop->setBound().
- next\_key\_part() is an iterator that advances the counter key\_part and returns false when you reach the end of the key
- get\_column() maps a key part to its Column in the dictionary



```

class index_object {
public:
    int key_part;
    server_rec *server;
    struct QueryItems *q;
    int n_parts;

    index_object(struct QueryItems *queryitems, request_rec *r) {
        q = queryitems;
        server = r->server;
        key_part = 0;
    };
    virtual ~index_object() {};

    virtual NdbOperation *get_ndb_operation(NdbTransaction *) = 0;
    bool next_key_part() { return (key_part++ < n_parts); };
    virtual int set_key_part(config::key_col &, mvalue &) = 0;
    virtual const NdbDictionary::Column *get_column() {
        return q->idx->getColumn(key_part);
    };
};

```



# Transactions and Operations

*mod\_ndb.h*

```
struct mod_ndb_instance {
    struct mod_ndb_connection *conn;
    Ndb *db;
    NdbTransaction *tx;
    int n_read_ops;
    int max_read_ops;
    struct data_operation *data;
    struct {
        unsigned int has_blob : 1 ;
        unsigned int aborted : 1 ;
        unsigned int use_etag : 1 ;
    } flag;
    unsigned int requests;
    unsigned int errors;
};
```

```
typedef struct mod_ndb_instance
    ndb_instance;
```

```
/* An operation */
struct data_operation {
    NdbOperation *op;
    NdbIndexScanOperation *scanop;
    NdbBlob *blob;
    unsigned int n_result_cols;
    const NdbRecAttr **result_cols;
    result_format_type result_format;
};
```

At startup time, an array of *max\_read\_ops* *data\_operation* structures is allocated for each *ndb\_instance*.

ndb_instance
conn
db
tx
n_read_ops
max_read_ops
data
flag
requests
errors

data_operation
op
scanop
blob
n_result_cols
result_cols
result_format

0

...

data_operation
op
scanop
blob
n_result_cols
result_cols
result_format

↓

*max\_read\_ops*

*Query.cc*

Individual operations are processed in *Query.cc*. The *Query()* function uses the configuration and the query string to determine an "access plan" and create an appropriate *NdbOperation*.

In a subrequest, processing ends after *Query()*, but in a complete request it passes immediately into *ExecuteAll()*.

*Execute.cc*

In *ExecuteAll()* (*Execute.cc*), we execute the transaction and then collect and format the results. In an ordinary request, a single result page is sent to the client. In a subrequest, though, the final call into *"/ndb-exec-batch"* (the *execute handler*) calls directly into *Execute.cc*, executes the transaction, and iterates over the all the operations (from 0 to *n\_read\_ops*), storing the results in the Apache notes table.

# Encoding and decoding NDB & MySQL data types

```
namespace MySQL {
    void result(result_buffer &, const NdbRecAttr &);
    void value(mvalue &, ap_pool *,
               const NdbDictionary::Column *,
               const char *);
};
```

*MySQL\_Field.h*

MySQL
result()
value()

---

## Decoding

- result() is a generic "decode" function; it converts an NdbRecAttr to a printable ASCII value
- Decoding is handled by some private functions inside of MySQL\_Field.cc, including String(), Time(), Date(), and Datetime()...

- String() can unpack three different sorts of strings packed into NDB character arrays.

```
enum ndb_string_packing {
    char_fixed,
    char_var,
    char_longvar
};
```

- Time(), Date() and Datetime() decode specially packed mysql data types.

---

## Encoding

- value() is a generic "encode" function; given an ASCII value (from HTTP) and an NdbDictionary::Column (which specifies how to encode the value), it will return an *mvalue* properly encoded for the database.

```
enum mvalue_use {
    can_not_use, use_char,
    use_signed, use_unsigned,
    use_64, use_unsigned_64,
    use_float, use_double,
    use_interpreted, use_null,
    use_autoinc
};
```

```
enum mvalue_interpreted {
    not_interpreted = 0,
    is_increment, is_decrement
};
```

### mvalues

```
struct mvalue {
    const NdbDictionary::Column *ndb_column;
    union {
        const char *      val_const_char;
        char *             val_char;
        int                val_signed;
        unsigned int       val_unsigned;
        time_t             val_time;
        long long           val_64;
        unsigned long long val_unsigned_64;
        float               val_float;
        double              val_double;
        const NdbDictionary::Column * err_col;
    } u;
    size_t len;
    mvalue_use use_value;
    mvalue_interpreted interpreted;
};
typedef struct mvalue mvalue;
```

# Output Formats and Result Buffers

Output formats are compiled using a hand-written scanner and parser into a tree structure, with Cells at the base.

<b>result_buffer</b>
<i>size_t</i> alloc_sz
<i>char</i> * buff
<i>size_t</i> sz
<i>char</i> * init()
<i>bool</i> prepare()
<i>void</i> putc()
<i>void</i> out()

*result\_buffer.h*

<b>len_string</b>
<i>size_t</i> len
<i>const char</i> * string

*output\_format.h*

```
enum re_type { const_string, item_name, item_value };
enum re_esc { no_esc, esc_xml, esc_json };
enum re_quot { no_quot, quote_char, quote_all };
```

<b>output_format</b>
<i>name</i>
<i>flags</i>
<i>symbol_table</i> []
<i>Node</i> *top_node
<i>Node</i> *symbol()
<i>char</i> *compile()
<i>void</i> dump()

<b>Node</b>
<i>char</i> *Name
<i>char</i> *unresolved
<i>Cell</i> *cell
<i>Node</i> *next_node
<b>virtual</b> <i>void</i> compile()
<b>virtual</b> <i>int</i> Run()
<b>virtual</b> <i>void</i> out()
<b>virtual</b> <i>void</i> dump()

<b>Cell : public len_string</b>
<i>re_type</i> elem_type
<i>re_quot</i> elem_quote
<i>const char</i> **escapes
<i>unsigned int</i> i
<i>Cell</i> *next
<i>void</i> out()
<i>void</i> chain_out()
<i>void</i> dump()

<b>Loop : public Node</b>
<i>Cell</i> *begin
<i>Node</i> *core
<i>len_string</i> *sep
<i>Cell</i> *end

<b>ScanLoop : public Loop</b>
<i>Cell</i> *begin
<i>Node</i> *core
<i>len_string</i> *sep
<i>Cell</i> *end

<b>RowLoop : public Loop</b>
<i>Cell</i> *begin
<i>Node</i> *core
<i>len_string</i> *sep
<i>Cell</i> *end

<b>RecAttr : public Node</b>
<i>char</i> *unresolved2
<i>Cell</i> *fmt
<i>Cell</i> *null_fmt

*output\_format.cc*

```
int build_results(request_rec *r, data_operation *data, result_buffer &res) {
    output_format *fmt = data->fmt;
    int result_code;

    if(fmt->flag.is_raw) return Results_raw(r, data, res);
    res.init(r, 8192);
    for(Node *N = fmt->top_node; N != 0 ; N=N->next_node) {
        result_code = N->Run(data, res);
        if(result_code != OK) return result_code;
    }
    return OK;
}
```

In *build\_results()*, a query result is built by running the nodes of the output format.