



API. Callers can continue to use regular configuration strings, or can pre-compile any configuration string that would be used in a performance path.

A second compatibility goal is with respect to the internal implementation of the API within WiredTiger. These API(s) need to be changed to allow the compiled strings to be used, as well as allow non-compiled strings to be used. These changes should be easy to do and not prone to error.

## The WT\_CONF structure

The primary structure used by *conf* is WT\_CONF. This provides information about both a single configuration string or a set of configuration strings. That is, in WiredTiger code, it replaces variables or arguments such as `const char *config` as well as arguments like `const char *cfg[]`.

To allow for fast checking of keys, we give an integer value to each key used in the entire set of configurations. This is done automatically by the `api_config.py` script in the `dist` directory. When WT code asks to get a configuration value for a key, it is now able to use an integer key rather than a string. Previously, code used "read\_timestamp" to look up a key, now it uses the integer `WT_CONF_ID_checkpoint_read_timestamp`. Some of this is hidden from view via macros, so that code that previously did:

```
WT_ERR(__wt_config_gets_def(session, cfg, "read_timestamp", 0, &cval));
```

can now do:

```
WT_ERR(__wt_conf_gets_def(session, conf, read_timestamp, 0, &cval));
```

Note that both of these code snippets return `/cval`, which is a **WT\_CONFIG\_ITEM**. So any code following this that uses the `cval` is unchanged. This has allowed us to make mostly mechanical changes on the WiredTiger side, reducing the possibility of errors.

Having one of these `WT_CONF_ID_*` (or "conf id") integer keys allows the conf functions to access needed information quickly. The WT\_CONF struct (which corresponds to `conf` in the second code snippet above) contains a table that is indexed by this `WT_CONF_ID_*`. That is here:

```
struct __wt_conf {
    uint8_t key_map[WT_CONF_ID_COUNT]; /* For each key, a 1-based index into conf_key */
    ...
};
```

The value of `WT_CONF->key_map[conf_id]` for a given `conf_id` gives us another index, which is called `conf_key_index` in the code. This value is either 0, meaning the key is not in this WT\_CONF, or it is an offset into an array of WT\_CONF\_KEY entries. Typically, a WT\_CONF\_KEY has an embedded **WT\_CONFIG\_ITEM**, which is what we want. Thus, returning a value from a WT\_CONF usually means two indirect references, and then we have a populated **WT\_CONFIG\_ITEM** to return.

If we were looking for greater performance, we could have modified this to use one table access by skipping having a `key_map`. Since there are currently over 300 possible keys in the system, this would have meant having a larger chunk of memory for each compiled string, 12K for starters (300 \* ~40 bytes for each WT\_CONF\_KEY). But supporting composite keys means that most compiled strings would need some multiple of that value. For now, we use a `key_map`, this could be reassessed.

## The WT\_CONF\_KEY structure

A WT\_CONF\_KEY is a simple struct that has a union along with a discriminant tag. The tag indicates how the union should be interpreted. The union is one of:

- a WT\_CONFIG\_ITEM. This is the typical case, it contains a value ready to return
- a sub-configuration index. This will include which WT\_CONF in the WT\_CONF array to use for subordinate keys.
- a WT\_CONF\_BIND\_DESC. This struct contains an offset into WT\_SESSION::conf\_bindings.values table, that's where the bound value will be found. It also contains info to perform runtime checks on the bound value when we get it.

## Position Independent Structures

The WT\_CONF struct is meant to be position independent. A position independent struct can be copied or moved to another memory location without "fixing up" any internal pointers, or allocating new structs. The only pointers in WT\_CONF are to items that are immutable and will never change locations, like strings. For example, when we do a compilation, we make a copy of the source string. That copy will always exist for as long as we have the connection open. The WT\_CONFIG\_ITEM entries that are referenced by WT\_CONF have pointers to parts of that immutable string. But the references in WT\_CONF to WT\_CONFIG\_ITEM entries, though, are not pointers. The reference is achieved by using an offset to an adjacent piece of memory.

What if the the WT\_CONF\_KEY array was not adjacent to the WT\_CONF? They could be allocated separately, but then the WT\_CONF would need a pointer to where the array was allocated. Then, copying the WT\_CONF would probably need a new WT\_CONF\_KEY array to be allocated and copied, and the pointer adjusted in the copy. And copies are needed to make changes to a base configuration. When we compile a configuration string, we want to take an (already compiled) WT\_CONF for all the default values for that configuration and copy it, and then overwrite it with the values in the configuration string. To do that, we need a fast copy without extra allocations and fix-ups.

Why do we need compilation to be fast anyway? After all, an efficient program will do pre-compiling in advance. The reason is to also efficiently handle any uses of the API that don't use pre-compiled strings. For those uses, we go ahead and pre-compile the given configuration string. This has several advantages:

- after pre-compiling, we will always have a single way to get configuration values. Imagine an API that supports pre-compiled strings. Suppose, if we are called with a (non-pre-compiled) configuration string, we leave it as is. Then we would need two mechanisms to access any key in the input configuration. We could pull this off, but it's a bit clumsy.
- efficiency. By pre-compiling, we've walked through the config string exactly once, and we never need to walk through it again. This saves us a bunch of time - that is, we get speedups even if callers are not doing explicit pre-compilation.
- testing. All of our existing tests are testing the code paths to pre-compile and access the pre-compiled struct. We don't need two modes of testing, or lots of extra test cases.

The total size of the WT\_CONF struct and its associated WT\_CONF\_KEY structs can be totally determined at compile time. This is also true for WT\_CONF structs that have subordinate WT\_CONF structs, which is used to support composite keys. Subordinate WT\_CONF structs appear after the parent WT\_CONF struct and before the WT\_CONF\_KEY array. Since they each have their own fixed size, but each size may be different, we have, within the WT\_CONF struct, an offset of the WT\_CONF\_KEY array. The WT\_CONF\_KEY\_TABLE\_ENTRY macro is used to access the WT\_CONF\_KEY entries, using the offset.

## Syntactic Sugar and Composite Keys

Suppose we have a fictional key "blah", and we want a slick way to get its associated ID (WT\_CONF\_ID\_blah) while using the syntax:

```
WT_ERR(__wt_conf_gets_def(session, blah, read_timestamp, 0, &cval));
```

Well, we can certainly use identifier pasting (the ## operator) via the C pre-processor to accomplish this. That's fine, but...

WiredTiger also makes heavy use of composite keys to reach down into sub-configurations (or "categories"). For example, suppose our fictional configuration string looks like "blah=7,foo=(bar=123,hello=456)" and we want to get the value for foo.bar. How do we identify these multiple parts to our access functions? And note that we have depths of sub-configuration greater than two. Ideally, we want to take code that currently looks like this:

```
WT_ERR(__wt_config_gets_def(session, cfg, "foo.bar", 0, &cval));
```

and change it to:

```
WT_ERR(__wt_conf_gets_def(session, conf, foo.bar, 0, &cval));
```

Our solution is to create a const struct in conf.h that contains something like this:

```
static const struct {
    uint64_t blah;
    ...
    struct {
        uint64_t bar;
    } Foo;
    ...
} WT_CONF_ID_STRUCTURE = {
    {
        WT_CONF_ID_blah,
        ...
        {
            WT_CONF_ID_foo | (WT_CONF_ID_bar << 16),
            ...
        },
        ...
    },
    ...
};
```

This is pretty close to our ideal, we need to use capitalization for walking through subordinate structures, like so:

```
WT_ERR(__wt_conf_gets_def(session, conf, Foo.bar, 0, &cval));
```

The capitalization was an unfortunate side effect of having an unregulated, but flexible namespace. For example, "checkpoint" is both a key and a category name.

As you might guess \_\_wt\_conf\_gets\_def is a macro that changes its third argument Foo.bar to WT\_CONF\_ID\_STRUCTURE.Foo.bar. And that in turn is a const value, it takes no space and can be fully evaluated at compile time to be:

```
WT_CONF_ID_foo | (WT_CONF_ID_bar << 16)
```

That allows the function that will be called (by the macro) to take a uint64\_t arg, whose layers can be peeled off as we are evaluating composite queries. It also means that we have a depth limit of four for subcategories.

## WT\_CONF storage layout

A WT\_CONF and all its sub-configurations fit contiguously in memory. The struct WT\_CONF for the configuration in question appears first, followed by the WT\_CONF for each sub-configurations, and additional WT\_CONF structs for any other subordinate configuration. These effectively form a WT\_CONF array. Following this array, all the WT\_CONF\_KEY structs that are used appear in a single array. To avoid extra calculations during processing, each WT\_CONF contains an offset to the slice of the WT\_CONF\_KEY array that it uses.

As an example, consider compiling the default configuration string for **WT\_SESSION::begin\_transaction**. Default configuration strings for every API (compiled or not) are listed in `config_def.c`. The string for `begin_transaction` lists all the legal keys and their default value:

```
"ignore_prepare=false,isolation=,name=,no_timestamp=false,"  
"operation_timeout_ms=0,priority=0,read_timestamp=,"  
"roundup_timestamps=(prepared=false,read=false),sync=,"
```

Notice that this has a sub-configuration:

```
"roundup_timestamps=(prepared=false,read=false)
```

that will have its own WT\_CONF struct. In total, we have two WT\_CONF structs. There are a total of eight top level keys, and two keys associated with `roundup_timestamps`. So we'll have an array of ten WT\_CONF\_KEY structs. This looks like:

WT\_CONF

– key\_map array

0
zeroes....
7
zeroes....
3
zeroes....
11
zeroes....
1
2
3
4
5
6
8
zeroes....

– offset to first  
WT\_CONF\_KEY

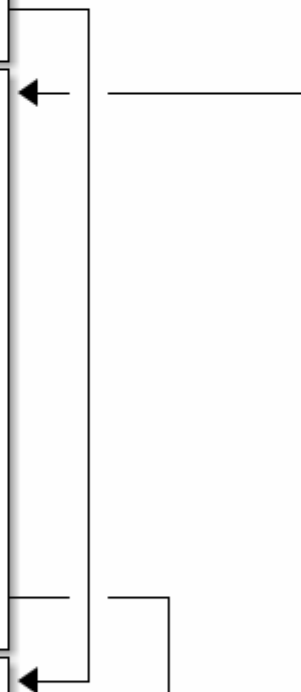
WT\_CONF  
for roundup\_timestamps

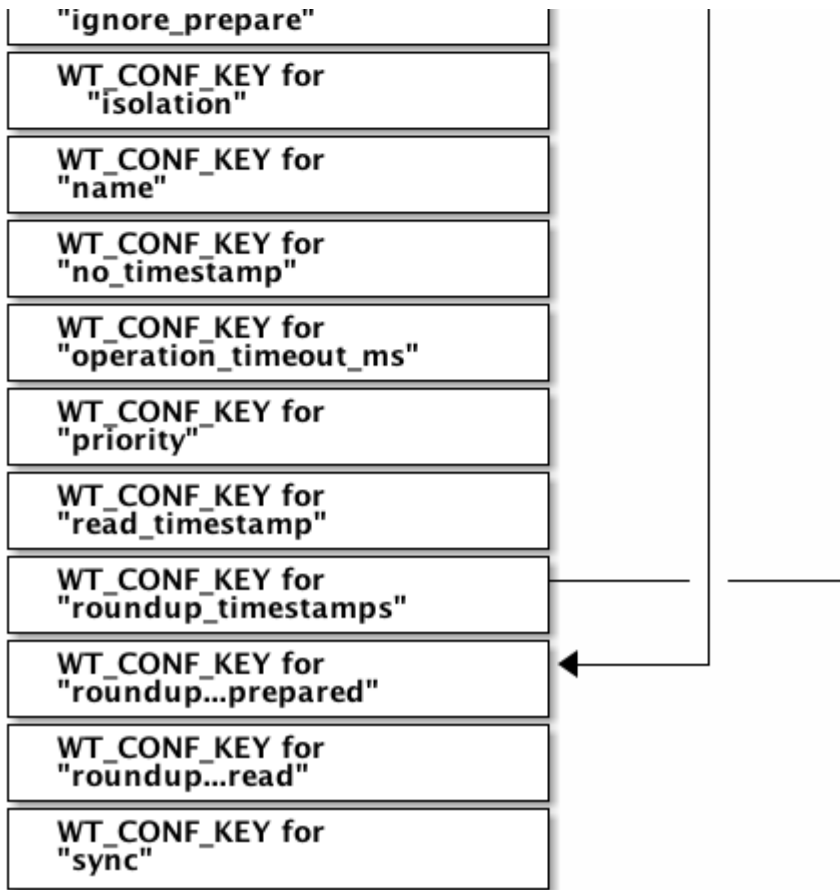
– key\_map array

0
zeroes....
1
2
zeroes....

– offset to first  
WT\_CONF\_KEY

WT\_CONF\_KEY for





You can see that `key_map` is a sparse array; only a small number of its 300 some entries are non-zero. Both `key_map` arrays are actually the same size, the second one is much more sparse, so the zeroed gaps are larger than they appear in this picture. Each entry in the map is a byte, which currently limits the number of `WT_CONF_KEY` array entries to 256.

Let's look at how we would look up `read_timestamp` in this `WT_CONF`. We'll always start at the top `WT_CONF`. `WT_CONF_ID_read_timestamp` is defined as 4, which means we'll look at `key_map[4]`. That is the entry that has 7. We then use 7 along with the offset to the associated `WT_CONF_KEY` array using a 1-based lookup. And that leads us to the `WT_CONF_KEY` for `read_timestamp`. A `WT_CONF_KEY` is a union, and its tag indicates that the `WT_CONFIG_ITEM` part of the union is used. That `WT_CONFIG_ITEM` is filled in according to the value indicated in the default string, that is, `"read_timestamp="`. That `WT_CONFIG_ITEM` is what is returned by the look up.

For a composite key, let's examine looking up `roundup_timestamps.read`. The definition of `WT_CONF_ID_roundup_timestamps` is 157, the top level `key_map[157]` has 8, that takes us to the eighth `WT_CONF_KEY`. Its union tag indicates it has sub-configuration at offset 1. That is, the second `WT_CONF` in the conf array. `WT_CONF_ID_read` is 159, so now we use the `key_map` in the second `WT_CONF`, that takes us to entry 2. We use the offset from the second `WT_CONF` to finally find the matching `WT_CONF_KEY` for `roundup_timestamps.read`. This is tagged as a `WT_CONFIG_ITEM`, so that's what we'll return.

## Binding Parameters

Up to now, we've been talking about fixed configuration strings. A more typical case is a configuration string where one or more parts of it may vary from call to call. For example, with `begin_transaction`, we might want to give each transaction a `read_timestamp`, and we only know that when we are ready to do the call. The compilation system allows unbound parameters that can later be bound. Currently, there are two kinds of bindings:

- strings, represented by "%s"
- integers and boolean, represented by "%d"

These special values indicate a deferred binding. For example, if we compile the string:

```
"ignore_prepare=%d,name=%s,read_timestamp=%d"
```

we'll need a way to mark these unbound parameters in the WT\_CONF struct (discussed below). Before using the compiled string, the API user will call **WT\_SESSION::bind\_configuration** with the actual values. Using the above compiled string, the caller might say:

```
session->bind_configuration(session, compiled_config_string,
    1, /* ignore_prepare, boolean 'true' */
    "txn_name123", /* name */
    0x1234); /* read_timestamp */
```

The binding values are first checked for validity if needed, just as a configuration string checks its values. The values are then put into the session in an array of **WT\_CONFIG\_ITEM**. The ordering of values given in the `bind_configuration` must match the order of the ' ' entries given in the configuration string.

We said we need a way to mark the unbound parameters in the WT\_CONF during compilation. That is done in the associated WT\_CONF\_KEY struct. It has a tag indicating a binding, and thus contains a WT\_CONF\_BIND\_DESC struct. This struct has the offset into the session's array of bound values. Also, an ordered list of pointers to these WT\_CONF\_BIND\_DESC structs makes it easy to find the right binding during the `bind_configuration` call.

This is one of the few times we allow pointers in WT\_CONF. It is safe because unbound parameters are not used internally by WiredTiger. That is, a WT\_CONF built for a default configuration has no unbound parameters, and thus has no pointers. So it can be freely copied when we create a WT\_CONF as a basis for a user's request to compile a configuration string.

The **WT\_SESSION::bind\_configuration** call will use the fill entries in the session's **WT\_CONFIG\_ITEM** array. Then, if we were to look up the key in question, we would find a WT\_CONF\_BIND\_DESC, and know what offset in the session array to find the needed **WT\_CONFIG\_ITEM**.

## Initial Pre-Compilation

TODO: flesh this out Why we pre-compile our default values:

- fast access when no configuration supplied.
- easy way to "pre-populate" a WT\_CONF when we compile a configuration string.

## Connection changes

TODO: flesh this out Where we store all the pre-compiled WT\_CONF structs, and the fake strings.

## Session changes

TODO: flesh this out Where we store any bound parameters.

## Fake Strings (references)

TODO: flesh this out When we compile, what we return to the caller as the "compiled" string, and how we use it to get to the WT\_CONF we need.



TODO: why we need position independent

TODO: why we can't use scoped identifiers.