

CHAPTER 9

pmemkv: A Persistent In-Memory Key-Value Store

Programming persistent memory is not easy. In several chapters we have described that applications that take advantage of persistent memory must take responsibility for atomicity of operations and consistency of data structures. PMDK libraries like `libpmemobj` are designed with flexibility and simplicity in mind. Usually, these are conflicting requirements, and one has to be sacrificed for the sake of the other. The truth is that in most cases, an API's flexibility increases its complexity.

In the current cloud computing ecosystem, there is an unpredictable demand for data. Consumers expect web services to provide data with predictable low-latency reliability. Persistent memory's byte addressability and huge capacity characteristics make this technology a perfect fit for the broadly defined cloud environment.

Today, as greater numbers of devices with greater levels of intelligence are connected to various networks, businesses and consumers are finding the cloud to be an increasingly attractive option that enables fast, ubiquitous access to their data. Increasingly, consumers are fine with lower storage capacity on endpoint devices in favor of using the cloud. By 2020, IDC predicts that more bytes will be stored in the public cloud than in consumer devices (Figure 9-1).

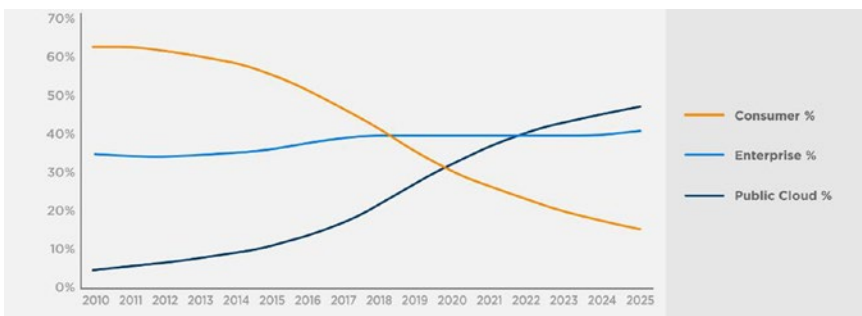


Figure 9-1. *Where is data stored?* Source: IDC White Paper – #US44413318

The cloud ecosystem, its modularity, and variety of service modes define programming and application deployment as we know it. We call it cloud-native computing, and its popularity results in a growing number of high-level languages, frameworks, and abstraction layers. Figure 9-2 shows the 15 most popular languages on GitHub based on pull requests.

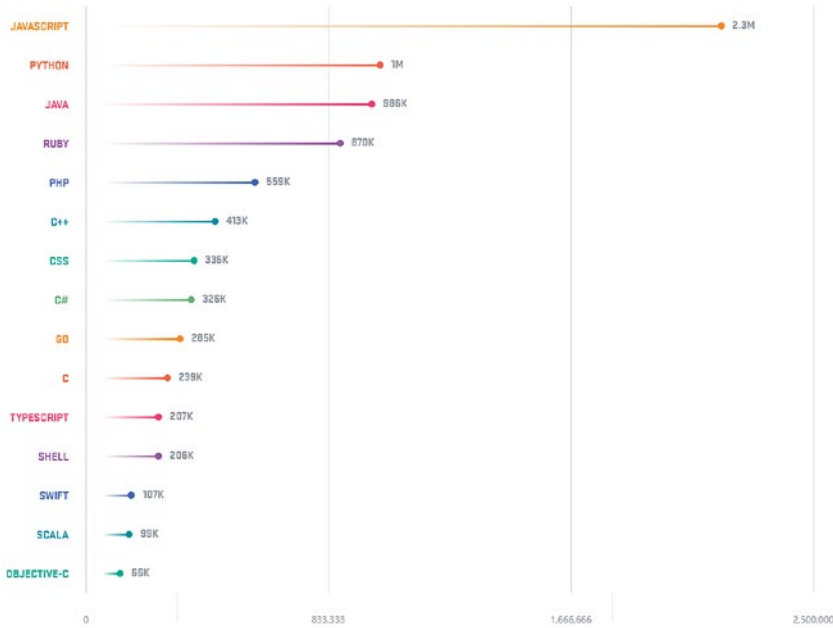


Figure 9-2. The 15 most popular languages on GitHub by opened pull request (2017). Source: <https://octoverse.github.com/2017/>

In cloud environments, the platform is typically virtualized, and applications are heavily abstracted as to not make explicit assumptions about low-level hardware details. The question is: how to make programming persistent memory easier in cloud-native environment given the physical devices are local only to a specific server?

One of the answers is a key-value store. This data storage paradigm designed for storing, retrieving, and managing associative arrays with straightforward API can easily utilize the advantages of persistent memory. This is why pmemkv was created.

pmemkv Architecture

There are many key-value data stores available on the market. They have different features and licenses and their APIs are targeting different use cases. However, their core API remains the same. All of them provide methods like put, get, remove, exists, open, and close. At the time we published this book, the most popular key-value data store is Redis. It is available in open source (<https://redis.io/>) and enterprise (<https://redislabs.com>) versions. DB-Engines (<https://db-engines.com>) shows that Redis has a significantly higher rank than any of its competitors in this sector.

| Rank | Name | Score |
|------|---------------------------|--------|
| 1. | Redis | 144,26 |
| 2. | Amazon DynamoDB | 56,42 |
| 3. | Microsoft Azure Cosmos DB | 29,08 |
| 4. | Memcached | 27,07 |
| 5. | Hazelcast | 8,27 |
| 6. | Aerospike | 6,59 |
| 7. | Ehcache | 6,56 |
| 8. | Riak KV | 6,06 |
| 9. | OrientDB | 5,69 |
| 10. | ArangoDB | 4,66 |
| 11. | Ignite | 4,26 |
| 12. | Oracle NoSQL | 3,46 |
| 13. | InterSystems Caché | 3,30 |
| 14. | LevelDB | 3,29 |
| 15. | Oracle Berkeley DB | 3,04 |

Figure 9-3. DB-Engines ranking of key-value stores (July 2019). Scoring method: https://db-engines.com/en/ranking_definition. Source: <https://db-engines.com/en/ranking/key-value+store>

Pmemkv was created as a separate project not only to complement PMDK's set of libraries with cloud-native support but also to provide a key-value API built for persistent memory. One of the main goals for pmemkv developers was to create friendly environment for open source community to develop new engines with the help of PMDK and to integrate it with other programming languages. Pmemkv uses the same BSD 3-Clause permissive license as PMDK. The native API of pmemkv is C and C++. Other programming language bindings are available such as JavaScript, Java, and Ruby. Additional languages can easily be added.

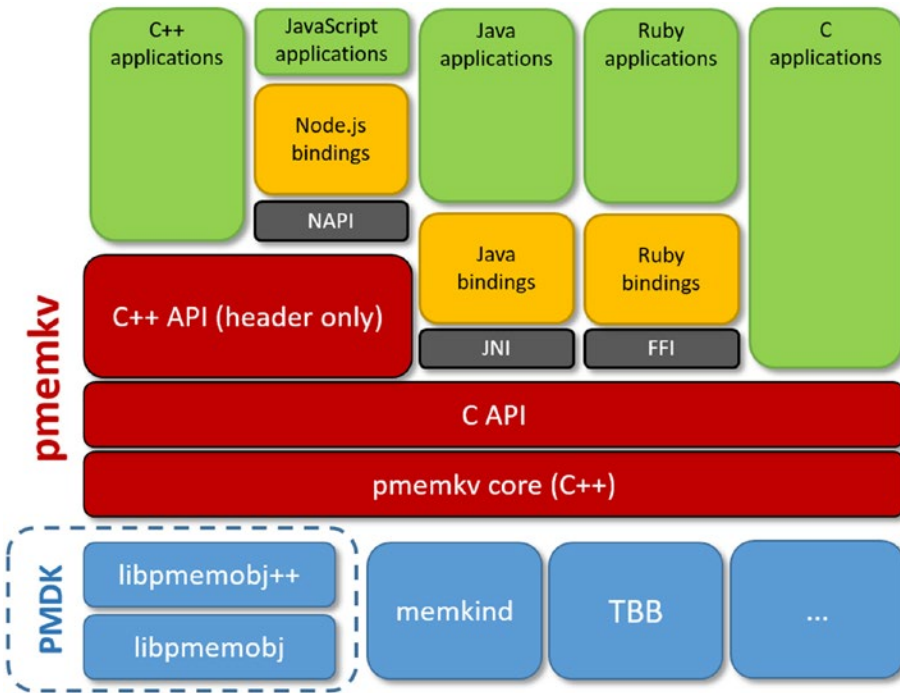


Figure 9-4. The architecture of pmemkv and programming languages support

The pmemkv API is similar to most key-value databases. Several storage engines are available for flexibility and functionality. Each engine has different performance characteristics and aims to solve different problems. Because of that, the functionality provided by each engine differs. They can be described by the following characteristics:

- Persistence: Persistent engines guarantee that modifications are retained and power-fail safe, while volatile ones keep its content only for the application lifetime.
- Concurrency: Concurrent engines guarantee that some methods such as `get()`/`put()`/`remove()` are thread-safe.
- Keys' ordering: "Sorted" engines provide range query methods (like `get_above()`).

What makes pmemkv different from other key-value databases is that it provides direct access to the data. This means reading data from persistent memory does not require a copy into DRAM. This was already mentioned in Chapter 1 and is presented again in Figure 9-5.

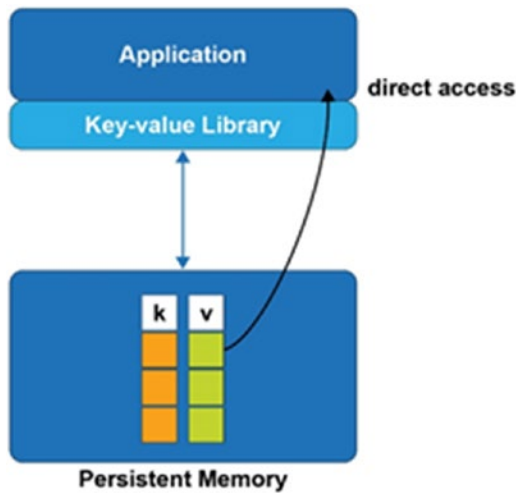


Figure 9-5. Applications directly accessing data in place using pmemkv

Having direct access to the data significantly speeds up the application. This benefit is most noticeable in situations where the program is only interested in a part of the data stored in the database. In conventional approaches, this would require copying the whole data in some buffer and returning it to the application. With pmemkv, we provide the application a direct pointer, and the application reads only as much as it is needed.

To make the API fully functional with various engine types, a flexible `pmemkv_config` structure was introduced. It stores engine configuration options and allows you to tune its behavior. Every engine has documented all supported config parameters. The pmemkv library was designed in a way that engines are pluggable and extendable to support the developers own requirements. Developers are free to modify existing engines or contribute new ones (<https://github.com/pmem/pmemkv/blob/master/CONTRIBUTING.md#engines>).

Listing 9-1 shows a basic setup of the `pmemkv_config` structure using the native C API. All the setup code is wrapped around the custom function, `config_setup()`, which will be used in a phonebook example in the next section. You can see how error handling is solved in pmemkv – all methods, except for `pmemkv_close()` and `pmemkv_errormsg()`, return a status. We can obtain error message using the `pmemkv_errormsg()` function. A complete list of return values can be found in pmemkv man page.

Listing 9-1. `pmemkv_config.h` – An example of the `pmemkv_config` structure using the C API

```

1  #include <stdio>
2  #include <cassert>
3  #include <libpmemkv.h>
4
5  pmemkv_config* config_setup(const char* path, const uint64_t fcreate,
6  const uint64_t size) {
7      pmemkv_config *cfg = pmemkv_config_new();
8      assert(cfg != nullptr);
9
10     if (pmemkv_config_put_string(cfg, "path", path) != PMEMKV_STATUS_OK) {
11         fprintf(stderr, "%s", pmemkv_errormsg());
12         return NULL;
13     }
14
15     if (pmemkv_config_put_uint64(cfg, "force_create", fcreate) !=
16         PMEMKV_STATUS_OK) {
17         fprintf(stderr, "%s", pmemkv_errormsg());
18         return NULL;
19     }
20
21     if (pmemkv_config_put_uint64(cfg, "size", size) != PMEMKV_STATUS_OK) {
22         fprintf(stderr, "%s", pmemkv_errormsg());
23         return NULL;
24     }
25
26     return cfg;
27 }
```

- Line 5: We define custom function to prepare config and set all required params for engine(s) to use.
- Line 6: We create an instance of C config class. It returns `nullptr` on failure.

- Line 9-22: All params are put into config (the `cfg` instance) one after another (using function dedicated for the type), and each is checked if was stored successful (`PMEMKV_STATUS_OK` is returned when no errors occurred).

A Phonebook Example

Listing 9-2 shows a simple phonebook example implemented using the `pmemkv` C++ API v0.9. One of the main intentions of `pmemkv` is to provide a familiar API similar to the other key-value stores. This makes it very intuitive and easy to use. We will reuse the `config_setup()` function from Listing 9-1.

Listing 9-2. A simple phonebook example using the `pmemkv` C++ API

```

37  #include <iostream>
38  #include <cassert>
39  #include <libpmemkv.hpp>
40  #include <string>
41  #include "pmemkv_config.h"
42
43  using namespace pmem::kv;
44
45  auto PATH = "/daxfs/kvfile";
46  const uint64_t FORCE_CREATE = 1;
47  const uint64_t SIZE = 1024 * 1024 * 1024; // 1 Gig
48
49  int main() {
50      // Prepare config for pmemkv database
51      pmemkv_config *cfg = config_setup(PATH, FORCE_CREATE, SIZE);
52      assert(cfg != nullptr);
53
54      // Create a key-value store using the "cmap" engine.
55      db kv;
56
57      if (kv.open("cmap", config(cfg)) != status::OK) {
58          std::cerr << db::errmsg() << std::endl;

```

```

59         return 1;
60     }
61
62     // Add 2 entries with name and phone number
63     if (kv.put("John", "123-456-789") != status::OK) {
64         std::cerr << db::errmsg() << std::endl;
65         return 1;
66     }
67     if (kv.put("Kate", "987-654-321") != status::OK) {
68         std::cerr << db::errmsg() << std::endl;
69         return 1;
70     }
71
72     // Count elements
73     size_t cnt;
74     if (kv.count_all(cnt) != status::OK) {
75         std::cerr << db::errmsg() << std::endl;
76         return 1;
77     }
78     assert(cnt == 2);
79
80     // Read key back
81     std::string number;
82     if (kv.get("John", &number) != status::OK) {
83         std::cerr << db::errmsg() << std::endl;
84         return 1;
85     }
86     assert(number == "123-456-789");
87
88     // Iterate through the phonebook
89     if (kv.get_all([](string_view name, string_view number) {
90         std::cout << "name: " << name.data() <<
91         ", number: " << number.data() << std::endl;
92         return 0;
93     }) != status::OK) {

```



```

94         std::cerr << db::errmsg() << std::endl;
95         return 1;
96     }
97
98     // Remove one record
99     if (kv.remove("John") != status::OK) {
100         std::cerr << db::errmsg() << std::endl;
101         return 1;
102     }
103
104     // Look for removed record
105     assert(kv.exists("John") == status::NOT_FOUND);
106
107     // Try to use one of methods of ordered engines
108     assert(kv.get_above("John", [](string_view key, string_view
109         value) {
110         std::cout << "This callback should never be called" <<
111             std::endl;
112         return 1;
113     }) == status::NOT_SUPPORTED);
114
115     // Close database (optional)
116     kv.close();
117
118     return 0;
119 }

```

- Line 51: We set the `pmemkv_config` structure by calling `config_setup()` function introduced in previous section and listing (imported with `#include "pmemkv_config.h"`).
- Line 55: Creates a volatile object instance of the class `pmem::kv::db` which provides interface for managing persistent database.
- Line 57: Here, we open the key-value database backed by the *cmap* engine using the config parameters. The *cmap* engine is a persistent concurrent hash map engine, implemented in `libpmemobj-cpp`.

You can read more about *cmap* engine internal algorithms and data structures in Chapter 13.

- Line 58: The `pmem::kv::db` class provides a static `errmsg()` method for extended error messages. In this example, we use the `errmsg()` function as a part of the error-handling routine.
- Line 63 and 67: The `put()` method inserts a key-value pair into the database. This function is guaranteed to be implemented by all engines. In this example, we are inserting two key-value pairs into database and compare returned statuses with `status::OK`. It's a recommended way to check if function succeeded.
- Line 74: The `count_all()` has a single argument of type `size_t`. The method returns the number of elements (phonebook entries) stored in the database by the argument variable (`cnt`).
- Line 82: Here, we use the `get()` method to return the value of the "John" key. The value is copied into the user-provided `number` variable. The `get()` function returns `status::OK` on success or an error on failure. This function is guaranteed to be implemented by all engines.
- Line 86: For this example, the expected value of variable `number` for "John" is "123-456-789". If we do not get this value, an assertion error is thrown.
- Line 89: The `get_all()` method used in this example gives the application direct, read-only access to the data. Both key and value variables are references to data stored in persistent memory. In this example, we simply print the name and the number of every visited pair.
- Line 99: Here, we are removing "John" and his phone number from the database by calling the `remove()` method. It is guaranteed to be implemented by all engines.
- Line 105: After removal of the pair "John, 123-456-789", we verify if the pair still exists in database. The API method `exists()` checks the existence of an element with given key. If the element is present, `status::OK` is returned; otherwise `status::NOT_FOUND` is returned.

- Line 108: Not every engine provides implementations of all the available API methods. In this example, we used the *cmap* engine, which is unordered engine type. This is why *cmap* does not support the `get_above()` function (and similarly: `get_below()`, `get_between()`, `count_above()`, `count_below()`, `count_between()`). Calling these functions will return `status::NOT_SUPPORTED`.
- Line 114: Finally, we are calling the `close()` method to close database. Calling this function is optional because `kv` was allocated on the stack and all necessary destructors will be called automatically, just like for the other variables residing on stack.

Bringing Persistent Memory Closer to the Cloud

We will rewrite the phonebook example using the JavaScript language bindings. There are several language bindings available for `pmemkv` – JavaScript, Java, Ruby, and Python. However, not all provide the same API functionally equivalent to the native C and C++ counterparts. Listing 9-3 shows an implementation of the phonebook application written using JavaScript language bindings API.

Listing 9-3. A simple phonebook example written using the JavaScript bindings for `pmemkv v0.8`

```

1  const Database = require('./lib/all');
2
3  function assert(condition) {
4      if (!condition) throw new Error('Assert failed');
5  }
6
7  console.log('Create a key-value store using the "cmap" engine');
8  const db = new Database('cmap', '{"path":"/daxfs/
kvfile","size":1073741824, "force_create":1}');
9
10 console.log('Add 2 entries with name and phone number');
11 db.put('John', '123-456-789');
12 db.put('Kate', '987-654-321');
13

```

```
14 console.log('Count elements');
15 assert(db.count_all == 2);
16
17 console.log('Read key back');
18 assert(db.get('John') === '123-456-789');
19
20 console.log('Iterate through the phonebook');
21 db.get_all((k, v) => console.log(` name: ${k}, number: ${v}`));
22
23 console.log('Remove one record');
24 db.remove('John');
25
26 console.log('Lookup of removed record');
27 assert(!db.exists('John'));
28
29 console.log('Stopping engine');
30 db.stop();
```

The goal of higher-level pmemkv language bindings is to make programming persistent memory even easier and to provide a convenient tool for developers of cloud software.

Summary

In this chapter, we have shown how a familiar key-value data store is an easy way for the broader cloud software developer audience to use persistent memory and directly access the data in place. The modular design, flexible engine API, and integration with many of the most popular cloud programming languages make pmemkv an intuitive choice for cloud-native software developers. As an open source and lightweight library, it can easily be integrated into existing applications to immediately start taking advantage of persistent memory.

Some of the pmemkv engines are implemented using `libpmemobj-cpp` that we described in Chapter 8. The implementation of such engines provides real-world examples for developers to understand how to use PMDK (and related libraries) in applications.



Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.