



*The flexibility and performance of low level, and simplicity of high level*

Takyon is a point to point communication API for C/C++ applications. It's designed for software engineers, in the embedded HPC (High Performance Computing) field, developing complex, multi-threaded applications for heterogeneous compute architectures. They high performance, determinism, and fault tolerance, but also abstraction from the details of the underlying communication protocol.

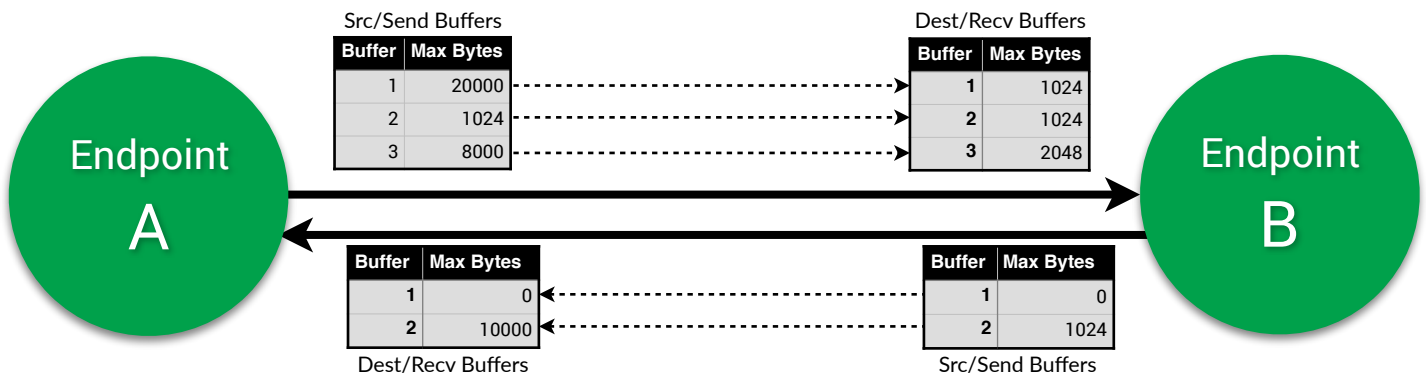
Takyon is designed to be:

- Efficient: A thin wrapper over low level communication APIs that still maintains critical features.
- High Performance: one-way (single underlying transfer), zero-copy (no extra buffering), two-sided transfers (recv gets notification when data arrives).
- Portable: Can work with most any modern interconnect, OS, bit width, and endianness.
- Unified: Same API for any locality (inter-thread/process/processor/application) and point to point model (reliable, connectionless, IO devices).
- Dynamic & Fault Tolerant: Create and destroy paths as needed. Timeouts can be used at any stage. Broken connections can be detected and restarted without effecting other connections.
- Scalable: No inherent limit to the number of connections, connection patterns (mesh, tree, etc.) or physical distance between connections.
- Simple to Use: While low level communications APIs may require weeks or months to learn, Takyon can be understood in a few days.

An application using the core Takyon APIs needs to include "takyon.h"

## Takyon Path Details

Connected interconnects: Support reliable, bi-directional, and multi-buffering. Good for distributed algorithms where every byte matters.



Connectionless interconnects: Support unreliable datagram transfers: Good for live-streaming, multicasting, and IO devices where data could be dropped.



## Transfer Details

The following must be specified when sending data (all are in bytes): size, source offset, and destination offset. The data size can be less than or equal to the max size of the buffer. The source and destination offset do not need to be the same. The destination offset must be zero for connectionless interconnects.





## Create a Path

Create a path to a remote endpoint. If it's a connected interconnect then both end points call this to complete the connection. If it's a connectionless interconnect then the remote endpoint is not required to exist. Attributes are stored in the returned path->attrs for easy reference.

`TakyonPath *takyonCreate(TakyonPathAttributes *attributes)`

Returns NULL on error, and the error\_message field in the attributes structure will be filled in (use TAKYON\_VERBOSITY\_ERRORS to print errors automatically).

TakyonPathAttributes Fields	Description
bool is_endpointA	true - Endpoint A of path. false - Endpoint B of path.
bool is_polling	true - Polling. false - Event driven.
bool abort_on_failure	true - Abort on error. false - Return status.
uint64_t verbosity	Or the bits of the Takyon verbosity flag values.
TakyonCompletionMethod send_completion_method	One of TAKYON_BLOCKING or TAKYON_USE_SEND_TEST (i.e. non blocking)
TakyonCompletionMethod recv_completion_method	Must be TAKYON_BLOCKING
double create_timeout	Timeout to wait for takyonCreate() to complete.
double send_start_timeout	Timeout to wait for takyonSend() to start transferring.
double send_complete_timeout	Timeout to wait for send to complete transferring.
double recv_complete_timeout	Timeout to wait for takyonRecv() to complete receiving.
double destroy_timeout	Timeout to wait for takyonDestroy() to disconnect.
int nbufs_AtoB	Number of buffers from endpoint A to B.
int nbufs_BtoA	Number of buffers from endpoint B to A.
uint64_t *sender_max_bytes_list	List of buffer sizes for the endpoint. Zero or greater.
uint64_t *receiver_max_bytes_list	List of buffer sizes for the endpoint. Zero or greater.
size_t *sender_addr_list	List of pre-allocated send buffers. Set each entry to NULL to auto-allocate.
size_t *receiver_addr_list	List of pre-allocated receive buffers. Set each entry to NULL to auto-allocate.
char interconnect[ MAX_TAKYON_INTERCONNECT_CHARS]	Defines the interconnect and it's properties to be used for transferring. The list of supported interconnects are defined in the implementation's README.txt

Verbosity flags	Value
TAKYON_VERBOSITY_NONE	0x00
TAKYON_VERBOSITY_ERRORS	0x01
TAKYON_VERBOSITY_INIT	0x02
TAKYON_VERBOSITY_INIT_DETAILS	0x04
TAKYON_VERBOSITY_RUNTIME	0x08
TAKYON_VERBOSITY_RUNTIME_DETAILS	0x10

Timeout values	Value
TAKYON_NO_WAIT	0
TAKYON_WAIT_FOREVER	-1
A double value, representing seconds: 1.25 - One and one quarter seconds 0.001 - One millisecond 0.00002 - 20 microseconds	0 or greater

Locality	Interconnect Specification (only the ones portable to all implementations are listed here)
Inter-thread	Memcpy -ID <ID> [-share]
Inter-process	Mmap -ID <ID> [-share] [-app_allocated_recv_mmap] [-remote_mmap_prefix <name>] Socket -local -ID <ID>
Inter-processor	Socket -remoteIP <IP> -port <port> Socket -localIP {<IP>   Any} -port <port> [-reuse]
Connectionless Inter-process/processor	SocketDatagram -unicastSend -remoteIP <IP> -port <port> SocketDatagram -unicastRecv -localIP {<IP>   Any} -port <port> [-reuse] SocketDatagram -multicastSend -localIP <IP> -group <IP> -port <port> [-disable_loopback] [-TTL <n>] SocketDatagram -multicastRecv -localIP <IP> -group <IP> -port <port> [-reuse]

## Transferring

Start sending a contiguous message (blocking and non-blocking):

`bool takyonSend(TakyonPath *path, int buffer_index, uint64_t bytes, uint64_t src_offset, uint64_t dest_offset, bool *timed_out_ret)`

Test if a non-blocking send is complete (only call this if path's send\_completion\_method is TAKYON\_USE\_SEND\_TEST):

`bool takyonSendTest(TakyonPath *path, int buffer_index, bool *timed_out_ret)`

Receive a message:

`bool takyonRecv(TakyonPath *path, int buffer_index, uint64_t *bytes_ret, uint64_t *offset_ret, bool *timed_out_ret)`

### Notes

Returns *true* if successful or *false* if failed. If failed, error text is stored in path->attrs.error\_message, and takyonDestroy() must be called. If a timeout occurs after data starts transferring, then is it considered an unrecoverable error, and *false* will be returned.

If the function returns *true* but *timed\_out\_ret* is *true*, then no bytes were transferred and it is safe to try again.

*timed\_out\_ret* can be NULL if the appropriate timeouts are set to TAKYON\_WAIT\_FOREVER.

When using takyonRecv(), *bytes\_ret* and *offset\_ret* can be NULL if there is no need for that information.

## Destroy a Path

Destroy the path. For connected interconnects, there is a coordination between both endpoints to make sure data in transit is flushed. The path handle will be set to NULL by this function call. Returns NULL on success, otherwise an error message is returned (and this message must be freed by the application).

`char *takyonDestroy(TakyonPath **path_ret)`



## Hello World (multi-threaded, bi-directional transfers)

```
#include "takyon.h"
```

Takyon Header

```
static const char *interconnect = NULL;
```

```
static void *hello_thread(void *user_data) {
    bool is_endpointA = (user_data != NULL);
```

```
TakyonPathAttributes attrs;
attrs.is_endpointA      = is_endpointA;
attrs.is_polling        = false;
attrs.abort_on_failure  = true;
attrs.verbosity         = TAKYON_VERBOSITY_ERRORS;
strncpy(attrs.interconnect, interconnect, MAX_TAKYON_INTERCONNECT_CHARS);
attrs.create_timeout    = TAKYON_WAIT_FOREVER;
attrs.send_start_timeout = TAKYON_WAIT_FOREVER;
attrs.send_complete_timeout = TAKYON_WAIT_FOREVER;
attrs.recv_complete_timeout = TAKYON_WAIT_FOREVER;
attrs.destroy_timeout   = TAKYON_WAIT_FOREVER;
attrs.send_completion_method = TAKYON_BLOCKING;
attrs.recv_completion_method = TAKYON_BLOCKING;
attrs.nbufs_AtoB        = 1;
attrs.nbufs_BtoA        = 1;
uint64_t sender_max_bytes_list[1] = { 1024 };
attrs.sender_max_bytes_list        = sender_max_bytes_list;
uint64_t recver_max_bytes_list[1] = { 1024 };
attrs.recver_max_bytes_list        = recver_max_bytes_list;
size_t sender_addr_list[1]        = { 0 };
attrs.sender_addr_list             = sender_addr_list;
size_t recver_addr_list[1]         = { 0 };
attrs.recver_addr_list             = recver_addr_list;
```

Setup Path Attributes

This could be offloaded to a configuration file: either hand coded or generated from a dataflow design tool

```
TakyonPath *path = takyonCreate(&attrs);
```

Path Creation

```
const char *message = is_endpointA ? "Hello from endpoint A" : "Hello from endpoint B";
for (int i=0; i<5; i++) {
    if (is_endpointA) {
```

```
        strncpy((char *)path->attrs.sender_addr_list[0], message, path->attrs.sender_max_bytes_list[0]);
        takyonSend(path, 0, strlen(message)+1, 0, 0, NULL);
        takyonRecv(path, 0, NULL, NULL, NULL);
        printf("Endpoint A received message %d: %s\n", i, (char *)path->attrs.recver_addr_list[0]);
```

Endpoint A Transfers

```
    } else {
```

```
        takyonRecv(path, 0, NULL, NULL, NULL);
        printf("Endpoint B received message %d: %s\n", i, (char *)path->attrs.recver_addr_list[0]);
        strncpy((char *)path->attrs.sender_addr_list[0], message, path->attrs.sender_max_bytes_list[0]);
        takyonSend(path, 0, strlen(message)+1, 0, 0, NULL);
```

Endpoint B Transfers

```
    }
}
```

```
takyonDestroy(&path);
```

Path Destruction

```
return NULL;
}
```

```
int main(int argc, char **argv) {
    if (argc != 2) { printf("usage: hello <interconnect>\n"); return 1; }
```

```
    interconnect = argv[1];
    pthread_t endpointA_thread_id;
    pthread_t endpointB_thread_id;
    pthread_create(&endpointA_thread_id, NULL, hello_thread, (void *)1LL);
    pthread_create(&endpointB_thread_id, NULL, hello_thread, NULL);
    pthread_join(endpointA_thread_id, NULL);
    pthread_join(endpointB_thread_id, NULL);
```

Run Endpoint A & B Threads

```
return 0;
}
```

### Notes:

- Endpoint A and endpoint B each do 5 sends and 5 receives, but notice how A sends first then B receives, then B is allowed to send then A receives. This is how application induced synchronization removes the need for Takyon to have implicit synchronization to manage correct use of the message buffers.
- No data addresses are passed to the takyonSend() or takyonRecv(). The message buffers are pre-allocated and registered with the transport enabling Takyon to do one-way, zero-copy, two-side transfers. You just can't get faster than that.
- The calls to takyonRecv() use NULL to ignore 'bytes\_received', 'dest\_offset', and 'timed\_out' since they were not needed.



The following Takyon extension functions are provided as open source to be linked into the application as needed. These convenience APIs are provided as source code instead of libraries in order to let the developers duplicate and modify the source to best fit the application needs. It also has the added benefit of simplifying certification.

An application using the Takyon extension functions needs to:

- Include “**takyon\_extensions.h**”, located in "Takyon/extensions/". Reference all the extension data structures in this header file.
- Compile and link the appropriate extension C files, located in "Takyon/extensions/"

## Endian (takyon\_endian.c)

Helpful functions if the endian of all threads are not the same. Swapping is in bytes.

```
bool takyonEndianIsBig();
void takyonEndianSwapUInt16(uint16_t *data, uint64_t num_elements)
void takyonEndianSwapUInt32(uint32_t *data, uint64_t num_elements)
void takyonEndianSwapUInt64(uint64_t *data, uint64_t num_elements)
```

## Time (takyon\_time.c)

Put the thread to sleep for the specified amount of time.

```
void takyonSleep(double seconds)
```

Returns current system wall clock time in seconds.

```
double takyonTime()
```

## Named Memory Allocation (takyon\_mmap.c)

Allocate named memory than can be accessed by remote processes in the same OS. This may be helpful with the collective gather operation when the 'Mmap' interconnect is used.

```
void takyonMmapAlloc(const char *map_name, uint64_t bytes, void **addr_ret, TakyonMmapHandle *mmap_handle_ret)
void takyonMmapFree(TakyonMmapHandle mmap_handle)
```

## Path Attributes (takyon\_attributes.c)

Simplify setting up the attributes for a Takyon path. The attribute structure can still be modified after initially calling **takyonAllocAttributes()**..

```
TakyonPathAttributes takyonAllocAttributes(bool is_endpointA, bool is_polling, int nbufs_AtoB, int nbufs_BtoA, uint64_t bytes, double
timeout, const char *interconnect)
void takyonFreeAttributes(TakyonPathAttributes attrs)
```

## Collective Groups (takyon\_collective.c)

Organize previously created Takyon paths into collective groups.

**Barrier:** (tree based, where each tree node can have any number of children)

```
extern TakyonCollectiveBarrier *takyonBarrierInit(int nchildren, TakyonPath *parent_path, TakyonPath **child_path_list)
extern void takyonBarrierRun(TakyonCollectiveBarrier *collective, int buffer)
extern void takyonBarrierFinalize(TakyonCollectiveBarrier *collective)
```

**Scatter:** One source, multiple destinations

```
TakyonScatterSrc *takyonScatterSrcInit(int npaths, TakyonPath **path_list)
TakyonScatterDest *takyonScatterDestInit(int npaths, int path_index, TakyonPath *path)
void takyonScatterSend(TakyonScatterSrc *collective, int buffer, uint64_t *nbytes_list, uint64_t *soffset_list, uint64_t
*doffset_list)
void takyonScatterRecv(TakyonScatterDest *collective, int buffer, uint64_t *nbytes_ret, uint64_t *offset_ret)
void takyonScatterSrcFinalize(TakyonScatterSrc *collective)
void takyonScatterDestFinalize(TakyonScatterDest *collective)
```

**Gather:** Multiple sources, one destination

```
TakyonGatherSrc *takyonGatherSrcInit(int npaths, int path_index, TakyonPath *path)
TakyonGatherDest *takyonGatherDestInit(int npaths, TakyonPath **path_list)
void takyonGatherSend(TakyonGatherSrc *collective, int buffer, uint64_t nbytes, uint64_t soffset, uint64_t doffset)
void takyonGatherRecv(TakyonGatherDest *collective, int buffer, uint64_t *nbytes_list_ret, uint64_t *offset_list_ret)
void takyonGatherSrcFinalize(TakyonGatherSrc *collective)
void takyonGatherDestFinalize(TakyonGatherDest *collective)
```

**One2One:** a set of paths organized in a useful way, such as a pipeline, mesh, etc.

```
TakyonCollectiveOne2One *takyonOne2OneInit(int npaths, int num_src_paths, int num_dest_paths,
TakyonPath **src_path_list, TakyonPath **dest_path_list);
void takyonOne2OneFinalize(TakyonCollectiveOne2One *collective);
```

**Reduce:** (coming soon)

**All2All:** (coming soon)



## Graph Description (takyon\_graph.c)

The following are helpful functions for creating a scalable application with collective communications. The goal is to make Takyon similar to MPI, but provide fully explicit control over the communication paths. This allows a very efficient use of paths in the application, for example:

- Use the same paths for different collective groups (scatter and gather), or use different paths for each collective group.
- Allow multiple paths between the same endpoints but with different performance properties (fast path for low latency commands, slow path for log data)
- Use a different interconnect for each path (RDMA, TCP sockets, UDP socket, etc).
- Use a different number of buffers for each path (single buffer, double buffer, etc).

Explicitly defining path attributes can make it much easier to fine tune the resources and performance of a real-time application on an embedded HPC application. To see examples of how to use the graph functionality, review the included examples: `hello_world_graph`, `scatter_gather`, and `pipeline`.

Define the communication layout of the application via a graph description file. The file format is as follows:

```
ThreadGroups
ThreadGroup: <name>
  Instances: <integer >= 1>
Processes
Process: <integer ID>
  Threads: <thread_name>[<thread_instance_index>] ...
MemoryBlocks
MemoryBlock: <name><buffer_index>
  ProcessId: <process ID>
  Where: <name>
  Bytes: <integer >= 1>
Paths
Defaults
  IsPolling: {true | false}, {true | false}
  AbortOnFailure: {true | false}, {true | false}
  Verbosity: <or'ing of verbosity flags>, <or'ing of verbosity flags>
  CreateTimeout: <takyon timeout value>, <takyon timeout value>
  SendStartTimeout: <takyon timeout value>, <takyon timeout value>
  SendCompleteTimeout: <takyon timeout value>, <takyon timeout value>
  RecvCompleteTimeout: <takyon timeout value>, <takyon timeout value>
  DestroyTimeout: <takyon timeout value>, <takyon timeout value>
  SendCompletionMethod: {TAKYON_BLOCKING | TAKYON_USE_SEND_TEST}, {TAKYON_BLOCKING | TAKYON_USE_SEND_TEST}
  RecvCompletionMethod: TAKYON_BLOCKING, TAKYON_BLOCKING
  NBufsAtoB: <integer >= 0>, <integer >= 0>
  NBufsBtoA: <integer >= 0>, <integer >= 0>
  SenderMaxBytesList: <space separated list of integers >= 0>, <space separated list of integers >= 0>
  RecvMaxBytesList: <space separated list of integers >= 0>, <space separated list of integers >= 0>
  SenderAddrList: <space separated list of NULL and/or <mem_name>:<offset>>, <list of NULL and/or <mem_name>:<offset>>
  RecvAddrList: <space separated list of NULL and/or <mem_name>:<offset>>, <list of NULL and/or <mem_name>:<offset>>
Path: <integer ID>
  Thread: <thread_name>[<thread_instance_index>], <thread_name>[<thread_instance_index>]
  InterconnectA: <interconnect specification>
  InterconnectB: <interconnect specification>
  <other attributes to override defaults>
CollectiveGroups
CollectiveGroup: <name>
  Type: { BARRIER | ONEZONE | SCATTER | GATHER }
  PathSrcIds: <space separated list of <pathID>:{A|B}>
```

The following must be defined one or more times:

ThreadGroup, Process, Path

The following can be defined zero or more times:

MemoryBlock, Defaults, CollectiveGroup

The 'Defaults' section does not need to specify all defaults, and each item in defaults can be in any order. Each time an item is specified in default, it replaces what the previous default was.

### Load and access a Takyon graph description from a file

```
TakyonGraph *takyonLoadGraphDescription(int process_id, const char *filename)
void takyonFreeGraphDescription(TakyonGraph *graph, int process_id)
```

### Create Takyon paths and collective groups (called by the threads)

```
void takyonCreateGraphPaths(TakyonGraph *graph, int thread_id)
void takyonDestroyGraphPaths(TakyonGraph *graph, int thread_id)
```

### Graph helper functions

```
void takyonPrintGraph(TakyonGraph *graph)
TakyonThreadGroup *takyonGetThreadGroup(TakyonGraph *graph, int thread_id)
int takyonGetThreadGroupInstance(TakyonGraph *graph, int thread_id)
TakyonCollectiveBarrier *takyonGetBarrier(TakyonGraph *graph, const char *name, int thread_id)
TakyonCollectiveOne2One *takyonGetOne2One(TakyonGraph *graph, const char *name, int thread_id)
TakyonScatterSrc *takyonGetScatterSrc(TakyonGraph *graph, const char *name, int thread_id)
TakyonScatterDest *takyonGetScatterDest(TakyonGraph *graph, const char *name, int thread_id)
TakyonGatherSrc *takyonGetGatherSrc(TakyonGraph *graph, const char *name, int thread_id)
TakyonGatherDest *takyonGetGatherDest(TakyonGraph *graph, const char *name, int thread_id)
```

If the `takyonLoadGraphDescription(<file>)` is used, then the user application has to define the following 2 functions to handle memory allocations even if no memory allocations are needed. These functions provided the flexibility for the application to allocated from any type of memory; e.g. CPU ram, MMAs, GPU ram, IO device memory, etc.

```
void *appAllocateMemory(const char *name, const char *where, uint64_t bytes, void **user_data_ret)
void appFreeMemory(const char *where, void *user_data, void *addr)
```



## Strided Message Transfers

**IMPORTANT:** Striding is not supported by most modern interconnects. If software tricks were used to implement striding for interconnects that did not support it, performance would be significantly impacted. Due to this, Takyon does not support strided functions in the core API set.

For the interconnects that do support striding, the following could be used as extensions to the implementation.

Start sending a strided message (blocking and non-blocking, callable only if the interconnect natively supports striding):

```
bool takyonSendStrided(TakyonPath *path, int buffer_index, uint64_t num_blocks, uint64_t bytes_per_block, uint64_t src_offset, uint64_t src_stride, uint64_t dest_offset, uint64_t dest_stride, bool *timed_out_ret)
```

Receive a strided message (callable only if the interconnect natively supports striding):

```
bool takyonRecvStrided(TakyonPath *path, int buffer_index, uint64_t *num_blocks_ret, uint64_t *bytes_per_block_ret, uint64_t *offset_ret, uint64_t *stride_ret, bool *timed_out_ret)
```

## Secure Communication

Still needs investigation if security can be fully specified within the attributes->interconnect[] specification. e.g. add "-SSL" to "Socket" interconnect.