

### Libuv based I/O Manager

-- 韩冬(github.com/winterland1989)

### The I/O schedule problem

- Modern OS manage work in unit of process/thread, and running threads often do I/O.
- I/O operation may block, e.g. stop one thread from running for an uncertain time.
- We want to process I/O concurrently as much as possible, e.g. threads doing I/O shouldn't stop others from doing their work.

### The simplest I/O multiplexer

- Start new threads to do concurrent I/O.
- Use blocking system calls, let the kernel do the rest, e.g. auto put thread blocking on I/O to background, wake up when I/O finishes.
- Easy to write programs, perform well under small concurrent numbers.
- OS thread is rather expensive when concurrent number goes up.

#### Enter Event-driven

```
// introduced in 4.2BSD Unix, released in August 1983.
int select(int nfds, fd set *readfds, fd set *writefds
          , fd set *exceptfds, struct timeval *timeout);
void FD CLR(int fd, fd set *set);
void FD SET(int fd, fd set *set);
void FD ZERO(fd set *set);
int FD ISSET(int fd, fd set *set);
// introduced in SVR3 Unix, released 1986.
int poll(struct pollfd *fds, nfds t nfds, int timeout);
struct pollfd {
   int fd; /* file descriptor */
   short events; /* requested events */
   short revents; /* returned events */
```

### Better event(epoll, kqueue...)

```
int epoll create(int size);
int epoll ctl(int epfd, int op, int fd
            , struct epoll event *event);
// op: EPOLL CTL ADD , EPOLL CTL MOD , EPOLL CTL DEL
struct epoll event {
   uint32_t events; /* Epoll events */
   epoll data t data; /* User data variable */
};
int epoll wait(int epfd, struct epoll event *events,
              int maxevents, int timeout);
```

### Overlapped I/O

```
int WSARecv(
 In SOCKET
                     S,
 _Inout_ LPWSABUF lpBuffers,
 In DWORD dwBufferCount,
 Out LPDWORD lpNumberOfBytesRecvd,
 _Inout_ LPDWORD lpFlags,
 _In_ LPWSAOVERLAPPED lpOverlapped,
 In LPWSAOVERLAPPED COMPLETION ROUTINE
                     lpCompletionRoutine
);
DWORD WSAWaitForMultipleEvents(
 In DWORD cEvents,
 _In_ const WSAEVENT *lphEvents,
 _In_ BOOL fWaitAll,
 _In_ DWORD dwTimeout,
 In BOOL fAlertable
```

# IOCP + Overlapped I/O

```
HANDLE WINAPI CreateIoCompletionPort(
         HANDLE FileHandle,
  In
 _In_opt_ HANDLE ExistingCompletionPort,
 _In_ ULONG PTR CompletionKey,
 _In_ DWORD NumberOfConcurrentThreads
BOOL WINAPI GetQueuedCompletionStatus(
 In HANDLE CompletionPort,
 _Out_ LPDWORD lpNumberOfBytes,
  Out PULONG PTR lpCompletionKey,
  Out LPOVERLAPPED *lpOverlapped,
 In DWORD
                  dwMilliseconds
```

### Event-loop thread

- 1. Prepare a watching set of fds with interested event.
- 2. Wait on that set.
- 3. Retrieve events after waiting finish, loop to process events.
- 4. During processing, modify the watching set maybe, then go back to step 2.

## Cross-platform event I/O

- Libuv #ifdef for you so you don't have to.
- · Libuv handles weird corner-cases.
- Libuv supports regular file via a fixed size thread pool.
- Libuv use windows API, and correctly handle UTF-16 conversion, even emulate TTY's ansi control codes!



### Scale the event loop

- · Use multiple threads, each one runs its own event loop.
- Use fork() to pre-fork multiple worker process.
- Use multiple threads, but only one event loop thread, create new worker thread on events.

### Take control back

- Wrap callbacks into Promises, Futures...etc, add await, wait... into language.
- Provide real light-weight thread in runtime system,
   reuse mutex, channel...etc threading concepts.

#### Meet the GHC RTS

- Running STG model, CMM code.
- Light-weight thread(H-thread) on multiple capability, M:N model.
- Two synchronize primitives: MVar, retry(in STM). Basic atomic operation support.
- Cooperative scheduling via yield, Preemptive scheduling via hijack HpLIM.
- · Moving garbage collection, with basic stable pointer support.

### Combine RTS with event loop

- When a lightweight thread encounter a blocking I/O call, allocate an ID, issue event registration, then put the H-thread into background.
- Find proper opportunity to do event polling.
- After polling, get fired IDs back, process them by resuming associated lightweight thread.

### How does MIO do it?

- Use fd as the ID. Save a STG callback into an IntTable, the callback is just a putMVar/writeTVar.
- After event registration, block the H-thread by takeMVar/retry, now the H-thread is removed from run-queue.
- An I/O manager H-thread runs alongside with all the user's H-threads on the same capability. Executes a zero timeout poll or forever poll depend on previous result (If successive two polls return no events, we just enter forever poll directly).
- After poll returns, use event's fd as ID to find registered STG callback, execute them, unblocking the H-thread.

### Libuv's stream API

```
int uv read start(uv stream t* stream
                 , uv alloc cb alloc cb
                 , uv read cb read cb)
void (*uv read cb)(uv stream t* stream
                  , ssize t nread
                  , const uv buf t* buf)
void (*uv_alloc_cb)(uv_handle_t* handle
                   , size t suggested size
                   , uv buf t* buf)
int uv read stop(uv stream t*)
int uv tcp init(uv loop t* loop, uv tcp t* handle)
int uv run(uv loop t* loop, uv run mode mode)
typedef enum {
    UV RUN DEFAULT = 0,
    UV RUN ONCE,
    UV RUN NOWAIT
} uv run mode;
```

#### Pass context with data field

```
void* uv loop t.data
typedef struct {
   size t event counter;
   size_t* event_queue;
  char** buffer table;
   ssize t* buffer size table;
} hs loop data;
// uv handle t is the base type for uv stream t, etc.
void* uv handle t.data
```

#### Pass context with data field

• Core data structure in each I/O manager thread:

```
data UVManager = UVManager
    -- The parking lot
    { uvmBlockTable :: IORef (UnliftedArray (MVar ()))
    -- uv_loop_t*
    , uvmLoop :: Ptr UVLoop
    -- uv_loop_t->data
    , uvmLoopData :: Ptr UVLoopData
    ...
}
```

- Allocate our own integer UVSlot as ID. Hold a growable global MVar table per I/O manager indexed by slot (the parking lot).
- When new I/O handle is initialized, allocated a slot, poke the slot to the uv\_handle\_t's data field. the MVar under slot's index is used to pause/resume I/O H-thread (the parking spot).
- After event registration, block the H-thread by takeMVar.

#### How does stdio allocate slot?

#### A Haskell version

```
data UVManager = UVManager
    { . . .
    , uvmFreeSlotList :: MVar [UVSlot]
allocSlot :: HasCallStack => UVManager -> IO UVSlot
allocSlot (UVManager blockTableRef freeSlotList loop ...) =
   modifyMVar freeSlotList $ \ freeList -> case freeList of
        (s:ss) -> return (ss, s)
        [] -> -- double the blockTable and loop's data
freeSlot :: UVManager -> UVSlot -> IO ()
freeSlot (UVManager _ freeSlotList _ _ _ _ ) slot =
   modifyMVar freeSlotList $ \ freeList -> return (slot:freeList)
```

## How does stdio prepare buffer?

Before registering I/O read, allocate buffers in Haskell,
 poke them into the loop's buffer table.

#### How does libuv side work?

```
void hs alloc cb(uv handle t* handle, size t suggested size, uv buf t* buf) {
    size t slot = (size t)handle->data;
    hs loop data* loop data = handle->loop->data;
    buf->base = loop data->buffer table[slot];
    buf->len = loop data->buffer size table[slot];
void hs read cb (uv stream t* stream, ssize t nread, const uv buf t* buf) {
    size t slot = (size t)stream->data;
    hs loop data* loop data = stream->loop->data;
    if (nread != 0) {
        loop data->buffer size table[slot] = nread;
        loop data->event queue[loop data->event counter] = slot;
        loop data->event counter += 1;
        uv read stop(stream);
int hs uv read start(uv stream t* stream) {
    return uv read start(stream, hs alloc cb, hs read cb);
```

• I/O manager's polling strategy is similar to MIO's one.

```
step :: UVManager -> Bool -> IO CSize
step (UVManager blockTableRef freeSlotList loop loopData _ _ timer _) block =
            blockTable <- readIORef blockTableRef</pre>
    do
            clearUVEventCounter loopData -- clean event counter
            if block
            then if rtsSupportsBoundThreads
                then uvRunSafe loop uV RUN ONCE -- forever poll
                else do
                    uvTimerWakeStart timer 2 -- 2ms timer on non-threaded rts
                    uvRun loop uV RUN ONCE
            else uvRun loop uV RUN NOWAIT -- zero timeout poll
            (c, q) <- peekUVEventQueue loopData</pre>
            forM [0..(fromIntegral c-1)] $ \ i -> do
                slot <- peekElemOff q i</pre>
                lock <- indexArrM blockTable (fromIntegral slot)</pre>
                tryPutMVar lock () -- unlock ghc thread with MVar
            return c
```

- After user's H-thread get resumed, peek the buffer\_size table to get result.
- We apply the similar procedure to write, so each uv\_stream\_t got two slots.

```
readInput :: HasCallStack
         => UVStream -> Ptr Word8 -> Int -> IO Int
readInput (UVStream handle rslot _ uvm) buf len = do
       m <- getBlockMVar uvm rslot
       withUVManager' uvm $ do
           tryTakeMVar m
           pokeBufferTable uvm rslot buf len
           uvReadStart handle
       takeMVar m
       r <- peekBufferTable uvm rslot
       if r > 0 -> return r
            r == fromIntegral uV EOF -> return 0
            r < 0 -> throwUVIfMinus (return r)
```

### Wait, libuv is not thread-safe!

- You can block threadA with epoll\_wait, add more events with epoll\_ctl in threadB, still get your notification.
- But libuv maintain its own data structure, such as callback table, you can't do registration while another thread is still blocking on uv\_run, which is the case for an uvRunSafe loop uv\_RUN\_ONCE.
- · Libuv provides a thread-safe wake up mechanism:

### Wait, libuv is not thread-safe!

```
data UVManager = UVManager
    { . . .
    , uvmRunning :: MVar Bool
    , uvmAsync :: Ptr UVHandle
withUVManager :: HasCallStack
              => UVManager -> (Ptr UVLoop -> IO a) -> IO a
withUVManager uvm f = do
    r <- withMVar (uvmRunning uvm) $ \ running ->
        if running
        then do uvAsyncSend (uvmcAsync uvm)
                return Nothing
        else do r <- f (uvmLoop uvm)
                return (Just r)
    case r of Just r' -> return r'
                      -> yield >> withUVManager uvm f
```

### Wait, libuv is not thread-safe!

```
startUVManager :: HasCallStack => UVManager -> IO ()
startUVManager uvm@(UVManager _ running _ _ ) = loop
 where
   loop = do
       e <- withMVar running $ \ _ -> step uvm False
       if e > 0
       then yield >> loop
       else do
           yield
           e <- withMVar running $ \ -> step uvm False
           if e > 0 then yield >> loop
           else do
               <- swapMVar running True
               <- step uvm True
               <- swapMVar running False
               yield
               loop
```

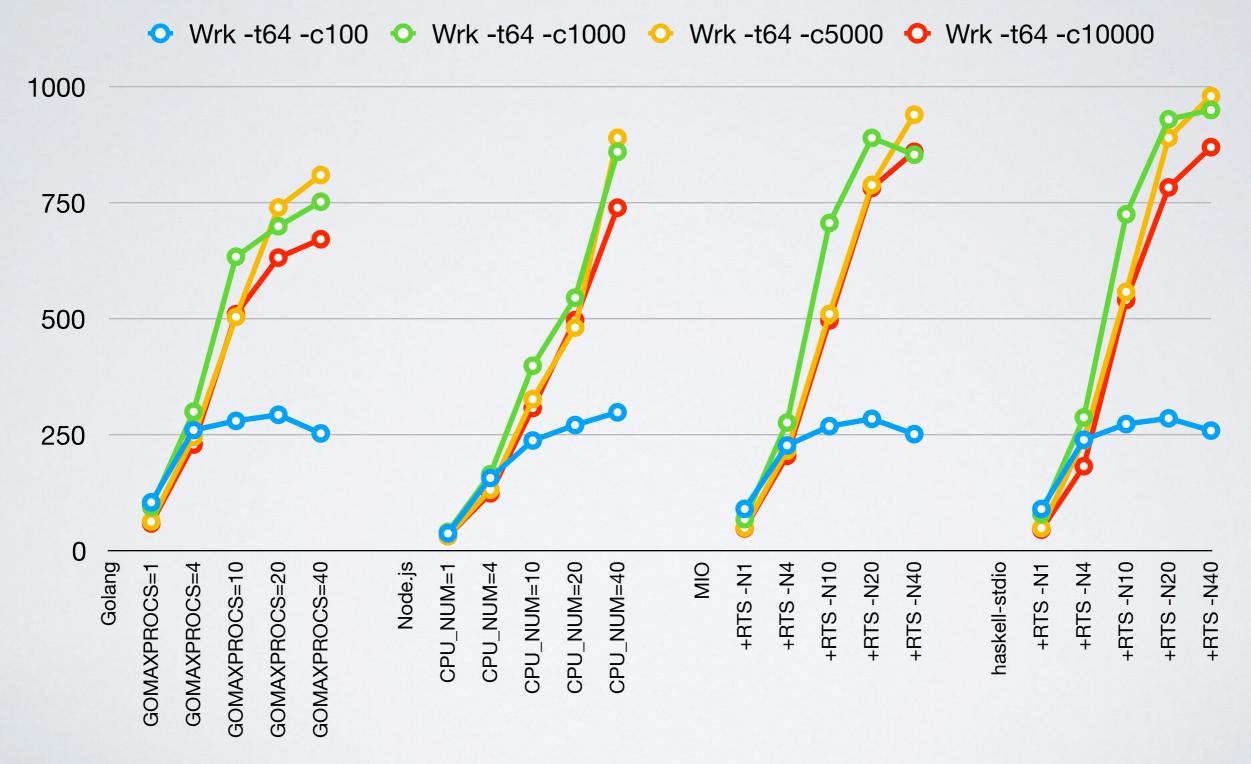
### OK, then. Is this lock expensive?

- Under load, there should be seldom chances to do safe block poll.
- Handle is bound to a certain uv\_loop\_t/UVManager,
  we should stop user's thread from migration, e.g. use
  forkOn instead of forkIO in server's accept loop.
- Without contention, locks perform just the same as a memory read & write.

### Benchmark setup

- Golang vs Node.js Cluster vs MIO vs stdio
- A small server serve 500 bytes in HTTP protocol
- Work load generator: 64 core xeon/128 G ram/10Gb NIC
- Server: 48 core xeon/256G ram/10Gb NIC
- Running wrk with -t64 and -c from 100 to 10000

# So, show me the figures!



# Things we tried

- Use StablePtr PrimMVar as ID, directly resume users' H-thread in C with hs\_try\_putmvar.
- Stable pointers are slow to allocated under contention,
   and make GC slow even on minor ones!
- The implementation of hs\_try\_putmvar asks for malloc/free of a PutMVar struct each invocation, which bring extra overhead.

# Things we tried

Use STM to notify users' H-thread uv\_run is finished.

```
data UVMState = UVM LOCKED | UVM POLLING | UVM FREE
data UVManager = UVManager {
    , uvmState :: TVar UVMState
uvAsyncSendSTM :: Ptr UVHandle -> STM ()
uvAsyncSendSTM = unsafeIOToSTM . uvAsyncSend
withUVManager :: HasCallStack => UVManager -> (Ptr UVLoop -> IO a) -> IO a
withUVManager uvm f = bracket_
    (atomically $ do
        state <- readTVar (uvmState uvm)</pre>
        case state of
            UVM LOCKED -> retry
            UVM POLLING -> uvAsyncSendSTM (uvmAsync uvm) >> retry
            UVM FREE -> swapTVar (uvmState uvm) UVM LOCKED)
    (atomically $ swapTVar (uvmState uvm) UVM FREE)
    (f (uvmLoop uvm))
```

# Things we tried

#### Allocate slot in C

```
typedef struct {
    size t* slot table;
    size_t free_slot;
    size t size;
} hs loop data;
size t alloc slot(uv loop t* loop) {
    hs loop data* loop data = loop->data;
    size t r = loop data->free slot;
    loop data->free slot = loop data->slot table[r];
    if (r == loop data->size - 1) {
         hs uv loop resize(loop, (loop_data->size) * 2);
    return r;
void free slot(uv loop t* loop, size t slot){
    hs loop data* loop data = loop->data;
    loop data->slot table[slot] = loop data->free slot;
    loop data->free slot = slot;
```

#### Future work

- Implement all libuv's I/O operation.
- Improve packed data structures, including text.
- TLS, HTTP protocol, etc.
- Improve stable pointer implementation, and use it.
- Combine H-thread pause/resume with RTS directly.