

Pushing the Mark6 to 60Gbps and Beyond with DPDK

J. Barrett

MIT Haystack Observatory
IVTW Meeting IX

October 23, 2024

- ① Mark6 background
- ② Updated dplane architecture and DPDK (Data Plane Development Kit)
- ③ Testing and results
- ④ Future work/outlook

Mark6 Background



- ① Development initiated in the late '00s
- ② Designed to support VLBI recording to readily swapped disk-modules
- ③ Nominally supports continuous recording at 4Gbps per disk-module (up to 16Gbps)
- ④ Hardware is supplied commercially by the Conduant Corp.
- ⑤ Software is split into two components:
 - ① cplane - the control plane handles user interaction via VSI-S command set, and system configuration
 - ② dplane - the data plane handles network-to-disk data movement

- Commodity computer hardware in custom chassis, but without other proprietary hardware
- Data network interface is via 10/25GbE
- Up to 4 disk modules of 8 HDDs each
- Modules can be quickly swapped via SAS cables



Mark6 - Hardware refresh

- ① Every few years, we need a hardware refresh due to component EoL
- ② Last refresh provided the opportunity to explore:
 - ① PCIe4.0
 - ② AMD EPYC CPUs (with 128 PCIe lanes)
 - ③ ATTO HBA
 - ④ Intel network cards



(a) Updated Mark6 interior with AMD EPYC, Supermicro H12SSL motherboard



(b) Mark6 + 2 R2DBEs in hypobaric chamber

- ① Using the existing d-plane software (R. Cappallo), demonstrated we can do $>16\text{Gbps}$.
- ② 32Gbps operation in two modes was tested: ($2 \times 16\text{Gbps}$) and ($4 \times 8\text{Gbps}$)
- ③ This mode also tested at 4500m equivalent in hypobaric chamber

Mark6(+) Hardware selection

- ❶ AMD 16-core CPU - (EPYC 7282, 128 lane PCIe 4.0)
- ❷ PCIe 4.0 capable motherboard (Supermicro H12SSL-CT with 64 GB RAM - supports up to 2TB)
- ❸ 1x 100GbE NIC: Intel E810-CQDA2 (QSFP28)
- ❹ 2x 10/25GbE NIC: Intel XXV710-DA2 (SFP+)
- ❺ 2x HBA: Atto H12F0GT 16-Port 6/12Gb/s SAS (SFF-8644 to SFF-8088)
- ❻ 4x Mark6 disk modules (8 drives each) - Seagate Exos16 16TB HDD
- ❼ Additional high speed cooling fans and custom airflow ducts

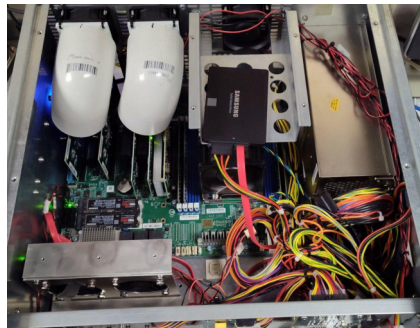


Figure: Interior of Mark6+ with 100GbE NIC

But what about a software refresh?



Several motivating factors for a software refresh:

- ❶ Need to move to a modern OS (Debian Squeeze and CentOS7 are both EoL)
 - ❶ Ubuntu 22.04 LTS selected - support to 2034 (FIPS option available)
- ❷ Need to support new hardware and 100GbE:
 - ❶ ngEHT - future backend (2x 32Gbps VDIF threads on a single 100Gbe interface)
 - ❷ DBEv5 - support for a direct 100Gbe connection (no switch)
- ❸ Need to demonstrate continuous 32Gbps operation of a Mark6 on a single 100Gbe interface.
- ❹ Optionally – demonstrate RX of multiple VDIF threads on a single interface
- ❺ Opportunity to explore newer technologies - particularly kernel bypass packet capture
 - ❶ PF_RING is used by pre-existing software, but various other technologies are available
 - ❷ DPDK looks like the leading candidate for this task
 - ❸ As specialized hardware not needed, and its appropriate for unidirectional UDP capture
- ❻ Need to upgrade cplane to python3

PF_RING vs. DPDK



- ❶ Existing Mark6 dplane software uses PF_RING as packet capture library
 - ❷ Zero-copy mode requires (per-MAC) paid-license, so not utilized
 - ❸ Must manage and tune cpu-interrupts for proper performance! – This can be quite a chore when CPUs (frequently) go EoL
 - ❹ More limited set of supported devices
 - ❺ Device still managed/visible to kernel network stack
 - ❻ Not required, but generally operated in promiscuous mode
- ❶ Moving to DPDK allows for:
 - ❷ Use of poll-mode drivers - removes the need of cpu interrupts entirely.
 - ❸ Takes full control of the network device away from the kernel
 - ❹ Buffered/burst packet capture with DMA from device to host-memory
 - ❺ Zero-copy packet manipulation (no license fee)
 - ❻ However: No kernel-based network management is possible (e.g. no ARP/ICMP support).
 - ❼ Ok, since we are treating it as a purely a point-to-point link
 - ❽ Also not required, but we operate in promiscuous mode

Updated dplane architecture



- ❶ Old architecture revolved around packet ring-buffers (per device)
- ❷ New packet processing relies on DPDK mbufs pool
- ❸ mbufs are pre-allocated – must enable hugepages in kernel params
- ❹ mbufs are filled on device burst-read, and free'd back to the pool on last use
- ❺ Instead of ring buffer, data is processed in chunks by various thread pools, and passed from task-to-task via locking stacks/queues
- ❻ With appropriately sized work items, lock contention is a non-issue
- ❼ For the most part, in order to simplify and preserve packet-ordering on reception, some thread-pools have size 1 (single threaded):
 - ❶ The packet receiver – one per interface
 - ❷ The VDIF thread ID sorter – one per interface (needed when multiple VDIF threads present on a single interface)
 - ❸ The scatter-gather block constructor (one per interface)
- ❽ Exception: The thread pool for SG block writers has size N_{files}

Packet life cycle

- 1 Packets are received by NIC
- 2 Packets are burst read into mbufs (extracted from pool) by poll-mode driver call
- 3 Packet burst is passed to VDIF thread ID sorter, packets are sorted and passed to block constructor
- 4 Each VDIF thread is assigned a scatter-gather block constructor, which collects packets until it is full
- 5 Completed blocks are passed to waiting writers, which write to next available HDD
- 6 Empty blocks are recycled back to the collectors, and used packets to the mbuf pool

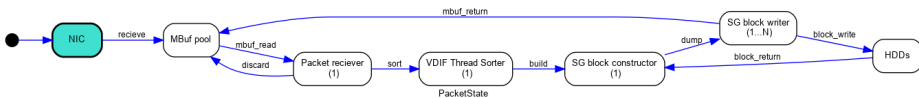


Figure: Packet mbuf cycle.

- ❶ Currently don't have a readily available back-end to serve data at required rates
- ❷ Need a dummy data generator (vdif headers + junk payload) on the cheap
- ❸ DPDK packet burst functionality + nanosleep to gets us to an approximate aggregate data rate
- ❹ Can then feed the 100GbE NIC back to itself
- ❺ After packet capture, we gather the scatter-gather files and check VDIF header validity with the tool dqa.

Loop-back testing

Using VDIF packet simulator we tried the following configurations:

- ① 1x VDIF thread @ 32 Gbps - continuous, 0 dropped ✓
- ② 2x VDIF threads @ 64 Gbps (1 thread recorded) - continuous, 0 dropped ✓
- ③ 1x VDIF thread @ 59 Gbps - continuous, (dropped $<6e-05$) ✓
- ④ 1x VDIF thread @ 64 Gbps 75-80s before memory buffer* exhausted - burst mode only
- ⑤ 2x VDIF threads @ 64 Gbps 20-30s before memory buffer* exhausted - burst mode only



Figure: Back of Mark6 - 100GbE fiber in yellow

- This testing was done with single actuator drives with an individual write speed of 261MB/s
- Results suggest we are achieving about 85-90% of the theoretical aggregate write throughput
- Latest available HDDs have advertised max sustained write speeds up to 285MB/s which might enable us to go a recording rate as high as 64Gbps continuous without resorting to SSDs or multi-actuator drives
- *memory buffer size was 36GB

Some ugly details and dead-ends



- ❶ The kernel parameter `isolcpus` is necessary to keep the cores hosting the packet-receive thread and VDIF sorter thread from handling other tasks
- ❷ Write balancing is tricky - disk performance varies. Forcing equal writes across disks limits the speed to the slowest performer.
- ❸ Alternative: (first-come-first-served) results in $\sim 1\%$ size variance
- ❹ DPDK mbuf buffer pool size must be $2^M \times \text{MTU}$ (which limits RAM buffer configurability)
- ❺ When the RAM buffer is full, data must be dropped.
- ❻ Time-slip between packet simulator and system makes long schedules (24hrs) difficult to test.
- ❼ Some failed experiments:
 - ❶ Async `io_uring` seemed like a good idea, but was found to be slower than standard `fwrite`
 - ❷ Likewise, using unbuffered system call to 'write' (not `fwrite`), requires data to be packaged in 4K blocks
 - ❸ However that doesn't work with the VDIF packet size (8192+32), and since mbufs are not guaranteed to be contiguous, a copy is needed ✗

Future work/experiments



Near term goals:

- 1 Populate the simulated packets with payload data
- 2 Perform integration and verification testing with ngEHT DBE
- 3 The new dplane needs to be integrated with the cplane software (some minor API changes expected)

Future work:

- 1 Examine performance higher speed single-actuator drives and dual actuator HDDs (Seagate MACH.2)
- 2 A large RAM buffer could possibly obviate the need for continuous recording capability for most practical applications, but would be an interesting space to explore
- 3 Test 100GbE system performance at high altitude conditions

Acknowledgements



This work was done with the hardware and support of NASA-SGP, the ngEHT, and the EHT:

