NSM2 for Telescope Test

Mikihiko Nakao (KEK-IPNS)

mikihiko.nakao@kek.jp

April 10, 2013 現場 (Gemba) meeting at DESY



Goal

"Master Control" to control everybody

- TTD, COPPER, EVB and HLT (NSM2 native speakers)
- PXD, SVD (EPICS speakers)
- EUDET (foreign language)

Control means

- Set up and start subsystem in the correct sequence
- Actions to be made by the Master Control is high-level ones only
- Correct all the log messages with proper time stamp (msec unit) to understand the sequence and origin of a problem

"Logging Database"

- Events initiated by the Master
- Log messages from subsystems
- Beta-testing for the Belle II control system

NSM2 status

NSM2 alpha release (2013.3)

- Low-level technology choices were made and implemented (corelib)
- Very preliminary Belle II dedicated library is made (b2lib)

NSM2 beta release plan

- Database interface for log messages and NSM shared memory data
- mysql at the beginning, to be switched to what database group offers
- Belle II dedicated library for state handling
- First target system: CDC/Belle2link/COPPER system at KEK
- Currently lower priority w.r.t. FTSW firmware development

Running EUDET Run Control

- Get the zip file and unzip (I had a problem in getting from SVN)
- cd main and make (to compile GUI, go to gui and make)
- Open four terminals, and cd bin
- Make an empty configuration file: touch ../conf/default.conf
- Run ./TestRunControl.exe, ./TestLogCollector.exe,
 ./TestDataCollector.exe, ./TestProducer.exe (in this order)
- In run control window, type c to config, then b to begin run
- In producer window, type r to generate one event

GUI version (run control and log collector) looks similar

TestRunControl.exe

TestLogCollector.exe

donau:p6	· r 🗍 donautp12 · r
<pre>nakao@donau(31)% ./TestRunControl.exe DEBUG: listenaddress=tcp://localhost:44000 Commands p Print connections 1 [msg] Send log message c [cnf] Configure clients (with configuration 'cnf') r Reset s Status b [msg] Begin Run (with run comment 'msg') e End Run x Terminate clients q Quit ? Display this help Connect: (waiting) (127.0.0.1:52397) LogServer responded: tcp://localhost:44002 Connect: LogCollector (127.0.0.1:52397) Receive: OK from LogCollector (127.0.0.1:52397) Connect: (waiting) (127.0.0.1:52400) DataServer responded: RunNumber = 2 DataServer responded: Server = 'tcp://localhost:44001'</pre>	<pre>nakao@donau(21)% ./TestLogCollector.exe ###################################</pre>
Connect: DataCollector (127.0.0.1:52400) Receive: OK from LogCollector (127.0.0.1:52397) Receive: OK from DataCollector (127.0.0.1:52400) Receive: OK from DataCollector (127.0.0.1:52400) Receive: OK from DataCollector (127.0.0.1:52400) Connect: (waiting) (127.0.0.1:52402) Connect: Producer.Test (127.0.0.1:52402) Receive: OK from Producer.Test (127.0.0.1:52402) Receive: OK from Producer.Test (127.0.0.1:52402) c Receive: Receive: NONE: Wait from LogCollector (127.0.0.1:52402) c Receive: NONE: Wait from DataCollector (127.0.0.1:52402) Receive: NONE: Wait from Producer.Test (127.0.0.1:52402)	donau:p16 r nakao@donau(7)% ./TestProducer.exe Commands s s data Send StringEvent with 'data' as payload r size Send RawDataEvent with size bytes of random data (default=1k) 1 msg Send log message o msg Set status=0K w msg Set status=ERROR q Quit TestProducer.exe
Receive: NONE: Wait from LogCoffector (127.0.0.1:52397) Receive: OK: Configured () from DataCollector (127.0.0.1:52400) Receive: OK: Configured () from Producer.Test (127.0.0.1:52400) B Receive: NONE: Wait from DataCollector (127.0.0.1:52400) Receive: NONE: Wait from Producer.Test (127.0.0.1:52402) Receive: NONE: Wait from LogCollector (127.0.0.1:52397) Receive: OK from DataCollector (127.0.0.1:52400) Receive: NONE: Wait from LogCollector (127.0.0.1:52397) Receive: OK from Producer.Test (127.0.0.1:52397) Receive: OK from Producer.Test (127.0.0.1:52402)	Configuring. Start Run: 3 r r r connect: Producer.Test (127.0.0.1:35857) Configuring () Configured () Complete Event: 3.0 Complete Event: 3.1 Complete Event: 3.2 Complete Event: 3.3

TestDataCollector.exe

single host, 4 processes

Controlling EUDET

NSM2 to control EUDET

- CUI version of EUDET run control is easy to understand and modify
- EUDET can control NSM2, too, but in my point of view, EUDET is a much simpler program and there's no reason to do so

A quick hack

- CONFIG/START/STOP from NSM2 to EUDET was coded in 1 hour
- It is not yet working in a coordinated way, e.g., run numbers or error messages, but no technical problem forseen

NSM2RunControl.exe

NSM2 master

donau‡p6	· 🗉 🗌	donau:p14	· • 🗆
p Print connections		nakao@donau(67)% ./master MASTER	
1 [msg] Send log message		nsmlib_checkif: checking interface <lo></lo>	
c [cnf] Configure clients (with configuration 'cnf')		nsmlib_checkif: interface <lo> does not support broadcast</lo>	
r Reset		nsmlib_checkif: checking interface <wlan0></wlan0>	
s Status		shmkey=8120 size=616680	
b [msg] Begin Run (with run comment 'msg')			
e End Run		MASIER>config EUDEI	
X Jerminate clients		ac=2 av[0]=contig	
		MASIEK>start LUDEL 7	
// Display this help		ac=5 avLØ]=start ▼MASTFR>	
Connect: (waiting) (127.0.0.1:52837)		donauto5	· • [
LogServer responded: tcp://localhost:44002		halesoftenau(66)% /neminfo?	
Connect: LogCollector (127.0.0.1:52857)		1000000000000000000000000000000000000	
Receive: OK from LogCollector (127.0.0.1:52857)		NSMD daemon started at 2013.04.10 12:23:21 (SinCas ago)	
Lonnect: (waiting) (127.0.0.1:52840)		MASTER(1) at 172 31 0 116 DEPUTY(-1) is missing. I'm the MASTER nemd	
Lonnect: Producer.lest (127.0.0.1:52840)		NSME shukev = 81201 in $= 172.31 // 116$ pid $= 12015$	
Kecelve, UK Trom Producer.lest (127.0.0.1.02840)		created 5052s ago last updated 0s ago (up-to-date)	
Lonnect, (waiting) (127.0.0.1.52642)		CONN 2 nid= 1 pid=12046 sock= 6 rtim=unknow nodeid=1	
DataServer responded: Server = 5		CONN 3 nid= A pid=12048 sock= 5 rtim=unknow nodeid=0	
Databerver responded, berver – tcp://tocarnost.44001		NODE Ø MASTER pid/uid=12048/21022 @172.31.0.116	5 -1
Bacaiva: OK from LagCollector (127.0.0.1.52042)		NODE 1 EUDET pid/uid=12046/21022 @172.31.0.116	δ - <u>1</u>
Receive: OK from DataCollector (127.0.0.1:52007)		REQ Ø OK code=1000 hash=1801(262184,1000)	
Receive: OK from DataCollector (127.0.0.1:52842)		REQ 1 ERROR code=1001 hash=1810(262220,1001)	
Receive: OK from Producer Test (127.0.0.1:52840)		REQ 2 CONFIG code=1002 hash=173(262256,1002)	
12:32:47 999 CONFIG<=MASTER		REQ 3 START code=1003 hash=481(262292,1003)	
CONFIG message received		REQ 4 STOP code=1004 hash=984(262328,1004)	
Receive: NONE: Wait from DataCollector (127.0.0.1:52842)		alist = -1	
Receive: NONE: Wait from LogCollector (127.0.0.1:52837)		🖥 nakao@donau(67)% 📋 🛛 🚺 🚺 🚺 🚺 🚺 🚺	
Receive: NONE: Wait from Producer.Test (127.0.0.1:52840)		donauto17	· • []
12:32:48.499 OK=>MASTER (IDLE) configured		▲ 12:32:59 234 DBG: topwrited: write = 20/20 9013c00	
Receive: NONE: Wait from LogCollector (127.0.0.1:52837)		12:32:59 234 DBG: repairing willer 28/28 3813688	Я
Receive: OK: Configured () from DataCollector (127.0.0.1:52842)		12:32:59 234 DBG: writeg after f=00000000(00000000) =00000000 g=09013be0	ľ,
Receive: OK: Configured () from Producer.Test (127.0.0.1:52840)		12:53:01.237 DBG: topredy redylen=16/16 redyp=9006610	
12:32:59.234 START<=MASTER		12:33:01.237 DBG: tcprecy_recylen=28/28_recyp=9006620	
CONFIG message received		12:33:01.237 DBG: tcprecy reg=1000 len=20 npar=2 p=[4099.1] coni=0	
Receive: NONE: Wait from LogCollector (127.0.0.1:52857)		12:33:01.237 DBG: command: reg = 1000 len = 20 npar = 2 datap = RUNNING cd	on = 3
Keceive: NUNE: Wait from Producer. lest (127.0.0.1:52840)		12:33:01.237 DBG: tcpsend f=09014028(00000000) =09014028 q=09014028 1	
Keceive: NUNE: Wait from DataLollector (127.0.0.1:52842)		12:33:01.237 DBG: tcpwriteg reg=1000 len=20 npar=2 p=[4099,1] buf=9014048.	writep
Receive: UK from Datalollector (127.0.0.1:52842)		=9014048	
IZ-55-01.257 UK=>MASTEK (KUNNING) run started		12:33:01.237 DBG: tcpwriteq: write = 44/44 9014048	
Receive: NUNE: Wait from Logioffector (127.0.0.1.52857)		12:33:01.237 DBG: writeg before f=09014028(00000000)	3
Reserve: OK from Dataconjector (127.0.0.1.52042)		12:33:01.237 DBG: writeq after f=00000000(00000000) l=00000000 q=09014028;	

nsmd2 daemon

7 processes (4 above +3 EUDET processes)

NSM2EUDET plan

- Proper response and error forwarding back to NSM2
- Log message forwarding to NSM2
- Configuration parameters passed from NSM2
- Rest of the world is kept as it is from NSM2 point of view, but DataCollector will be modified to connect to EB2

Conclusion

- EUDAQ Run Control system is easy to understand (to me)
- EUDAQ Run Control can be easily controlled by NSM2, while leaving the rest of the EUDAQ world as it is
- Still a lot more to do on NSM2 side

Backup

Introduction to NSM



- nsmd2 daemon process is running on each host
- user (client) program communicate with other user program through nsmd via NSM API
- network shared memory is visible as POSIX shared memory

NSM2

Full rewriting, instead of modifying existing code

Aims of changes

- Making the system more robust against a flood of messages
- Reducing / removing hard coded features
- Make the behaviour more deterministic
- Make the allowed node-name longer, memory size bigger, etc
- Multi-network in one process: supporting a network bridge
- Factorize NSM-generic code and Belle(-II) specific features
- (A better support on database, GUI)
- Still minimizing the look-and-feel difference from NSM1
- Make it free from 2038 (unix-year) problem 64-bit time_t

Code writing started in November/2012



Today, alpha version of NSM2 package is released. https://belle2.cc.kek.jp/~twiki/bin/view/Detector/DAQ/NSM2

NSM2 and Belle II

NSM2 has two library layers: corelib and belle2lib

NSM1 has some Belle-dependent features embedded in the library

corelib

- Nothing is assumed about usage, message or node type
- All functions start with nsmlib_
- Flexibility at the cost of allowing wrong combinations

b2lib

- Heavily assumes the run/slow control model of Belle II
- All functions start with b2nsm_
- All necessary functions are (re)defined to avoid direct call to corelib
- Limited flexibility to minimize mistakes

NSM messages

How to represent a message

- A message code is a text-string, with a 16-bit hash code (In NSM1, it was a hard-coded symbol in an include file)
- In NSM2, the node name also has a hash code (in NSM1, a linear search was used)
- Hash is maintained in the MASTER nsmd, similarly to the nodename

How does it work (implemented, slightly modified since Jan)

- Message is define on-the-fly, by registering a callback function
- Number of parameters is variable (before it was fixed to 2)

(master)

(CDC)

```
int p[2]; void startfunc(NSMmsg *m,NSMcontext *c) {
p[0] = runno; :
p[1] = runtype; b2nsm_ok(m,"RUNNING","run %d",m->pars[0]); }
b2nsm_sendreq("CDC","START",2,p); :
b2nsm_callback("START", startfunc);
```

How to wait for NSM messages

Signal handler (default)

- User program can do anything else while waiting (e.g. GUI)
- May receive another message while processing the previous one
- This was the only method in NSM1

Wait function (by explicitly calling b2nsm_wait)

- Cannot do anything else until timeout
- Message order is guaranteed
- Has to be decided when calling nsm_init2

Writing a code to use select

- Socket is available in the NSMcontext object (returned by nsm_init)
- One can wait for NSM and something else at the same time
- Once there's something to read, one can call b2nsm_wait
- Need a bit more programming skill

NSM data

data structure definition (implemented)

- In NSM1, it was a hard-coded struct in an include file
- In NSM2, the include file is parsed on the fly (only a simple struct definition can be used)
- NSM has to know it to convert into network-byte-order
- Revision number to make sure the same data structure is used

```
(ttd_data.h)
const int
  ttd_data_revision = 1;
struct ttd_data {
    int32 runno;
    int32 evtno;
    double trigrate;
    byte status[128];
};
```

NSM data sharing

UDP broadcast (implemented)

- Allocated data structure is shared between nsmd nodes
- Single UDP packets if the size is below 1484 bytes, very minimal header, 8-byte for UDP and only 8-byte for NSM (but non-8-byte-boundary may cause a problem, probably will switch to 1480 byte)
- Above this size, data update time may not be uniform, but time stamp difference is recorded and can be monitored (it was not clear in NSM1 when a particular part of data is updated)
- Data size in 16-bit word, up to 65,296 byte (44×1484)

NSM2 screen shot

sonoma:p19 · 「	napa;p20 · r
10:16:45.387 DBG: writeq after f=0000000(0000000) l=00000000 q=08fda6c0 10:16:45.389 DBG: tcprecv recvlen=16/16 recvp=8fd0d00 10:16:45.389 DBG: tcprecv recvlen=41/41 recvp=8fd0d10 10:16:45.389 DBG: ccprecv req=1000 len=33 npar=2 p=[4098,0] coni=0 10:16:45.389 DBG: tcprecv req=1000 len=33 npar=2 p=[4098,0] buf=8fdab08 1 10:16:45.389 DBG: tcpsend f=08fdab08(00000000) l=08fdab08 q=08fdab08 1 10:16:45.390 DBG: tcpwriteq req=1000 len=33 npar=2 p=[4098,0] buf=8fdab28 writep =8fdab28 10:16:45.390 DBG: tcpwriteq: write = 57/57 8fdab28 10:16:45.390 DBG: writeq after f=00000000(00000000) l=00000000 q=08fdab08 10:16:45.390 recv USRCPYMEM(4608)<=130.87.227.110 len=20 npar=2 10:16:55.398 recv USRCPYMEM(5120)<=130.87.227.110 len=20 npar=2 10:16:55.398 recv USRCPYMEM(5376)<=130.87.227.110 len=20 npar=2 10:16:55.398 recv USRCPYMEM(532)<=130.87.227.110 len=20 npar=2 10:16:58.401 recv USRCPYMEM(5388)<=130.87.227.110 len=20 npar=2 10:17:01.405 recv USRCPYMEM(6440)<=150.87.227.110 len=20 npar=2 10:17:07.409 recv USRCPYMEM(6440)<=130.87.227.110 len=20 npar=2 10:17:07.409 recv USRCPYMEM(6400)<=130.87.227.110 len=20 npar=2	<pre>10:16:45.387 DBG: tcprecv recvlen=16/16 recvp=15cc7b0 10:16:45.387 DBG: tcprecv recvlen=4/4 recvp=15cc7c0 10:16:45.387 DBG: tcprecv req=1002 len=0 npar=1 p=[1,2] coni=0 10:16:45.387 DBG: tcpsend f=01620200(00000000) l=01620200 q=01620200 1 10:16:45.387 DBG: tcpwriteq req=1002 len=0 npar=1 p=[1,0] buf=1620230 writep=162 0230 10:16:45.387 DBG: tcpwriteq: write = 20/20 1620230 10:16:45.387 DBG: writeq before f=01620200(00000000) l=01620200 q=01620200 10:16:45.387 DBG: writeq after f=00000000(0000000) l=00000000 q=01620200 10:16:45.387 DBG: writeq after f=00000000(0000000) l=00000000 q=01620200 10:16:45.387 DBG: tcprecv recvlen=16/16 recvp=161dad0 10:16:45.387 DBG: tcprecv recvlen=41/41 recvp=161dae0 10:16:45.387 DBG: tcprecv req=1000 len=33 npar=2 p=[4098,0] coni=1 10:16:45.387 DBG: tcpwriteq req=1000 len=33 npar=2 p=[4098,0] buf=1620660 1 10:16:45.387 DBG: tcpwriteq: write = 57/57 1620690 10:16:45.387 DBG: tcpwriteq: write = 57/57 1620690 10:16:45.387 DBG: writeq after f=00000000(0000000) l=001620660 q=01620660 10:16:45.387 DBG: tcpwriteq: write = 57/57 1620690 10:16:45.387 DBG: tcpwriteq: write = 57/57 1620690 10:16:45.387 DBG: writeq after f=000000000(0000000) l=00000000) l=01620660 q=01620660 10:16:45.387 DBG: tcpwriteq: write = 57/57 1620690 10:16:45.387 DBG: writeq after f=000000000(0000000) l=00000000) l=01620660 q=01620660 10:16:45.387 DBG: writeq after f=000000000(0000000) l=00000000) l=01620660 q=01620660 10:16:45.387 DBG: writeq after f=000000000(0000000) l=00000000) l=01620660 q=01620660 10:16:45.387 DBG: writeq after f=000000000(0000000) l=00000000) l=000000000 </pre>
sonoma:p21 · 「 🗌	napato25 · · · ·
<pre>hakao@sonoma(2)% ./master master nsmlib_checkif: checking interface <lo> nsmlib_checkif: interface <lo> does not support broadcast nsmlib_checkif: checking interface <eth0> shmkey=8120 size=616680 master>start SLAVE 1 ac=3 av[0]=start master></eth0></lo></lo></pre>	nakao@napa(2)% ./client slave nsmlib_checkif: checking interface <lo> nsmlib_checkif: interface <lo> does not support broadcast nsmlib_checkif: checking interface <eth0> shmkey=8120 size=616680 allocmem: ret=0 dtpos=24 10:15:54.579 allocmem(slave, rev.1) at 0x7fc96a141018 10:15:54.619 callback(STAPE) registered</eth0></lo></lo>
sonoma:p29 nsmlib_checkif: interface <lo> does not support broadcast nsmlib_checkif: checking interface <eth0> shmkey=8120 size=616680</eth0></lo>	10:15:54.658 callback(STOP) registered 10:16:45.387 START<=MASTER START message received 10:16:45.387 OK=>MASTER (RUNNING) run 1 started as 1th run
RUNNING run=1 (total 1) evt=2 (total 2) RUNNING run=1 (total 1) evt=5 (total 5) RUNNING run=1 (total 1) evt=8 (total 8) RUNNING run=1 (total 1) evt=11 (total 11) RUNNING run=1 (total 1) evt=14 (total 14) RUNNING run=1 (total 1) evt=20 (total 20) RUNNING run=1 (total 1) evt=20 (total 20)	

2 hosts (2 nsmd2), 3 processes (master, client, readdat)

Missing items in the alpha release

- Missing library functions (e.g., nsmlib_closemem, ...)
- Remote logging function (e.g., b2nsm_warning, ...)
- Belle II run scheme design
- Various stress tests
 - Many nodes (up to limit)
 - Large datasize, fast data update
 - DoS attack-like flood-of-message handling
 - Colliding two NSM networks
 - Handling incorrect network setting, disconnecting network, etc
- **Full document** (especially reference manual for the library section)

Run state model for a readout subsystem



- State transitions, but states are kept inside belle2nsm library
- The system can be controlled by somebody other than MASTER
- Making a subsystem READY is the nasty part of the procedure...

Start-up model for COPPER



Run control wishlist

Constraints

 Cold start, firmware programming, parameter downloading, link establishing have to be done

Wish

- Readout cycle to be all the time running, regardless the HV / accelerator condition (idle running)
- In case of an error in one sub-system, the rest should not stop
- Transition from idle running to real running is a quick run number and run type change in the fast TT link command

"Problems"

- Detectors are willing to initialize everything every time regardless the deadtime caused by it
- Not easy to synchronize at the event builder
- FPGA reprogramming has to be included as we expect SEU

Status and Plan

Now testing with 4 PCs on general network

(a few more can be quickly added if needed)

Many new key features are now implemented, but still not ready for a test use

Target dates

- January DAQ WS for alpha release (for PocketDAQ) did not happen
- March B2GM for beta alpha release
- March B2GM for a first proposal of the state model and run start sequence or something alike for Pocket-DAQ

After March B2GM

- Throughput measurement with a dedicated fast (GbE) network
- Serious error handling tests with a large system (\gg 10 nodes)
- Master COPPER/HSLB/FTSW control for PocketDAQ