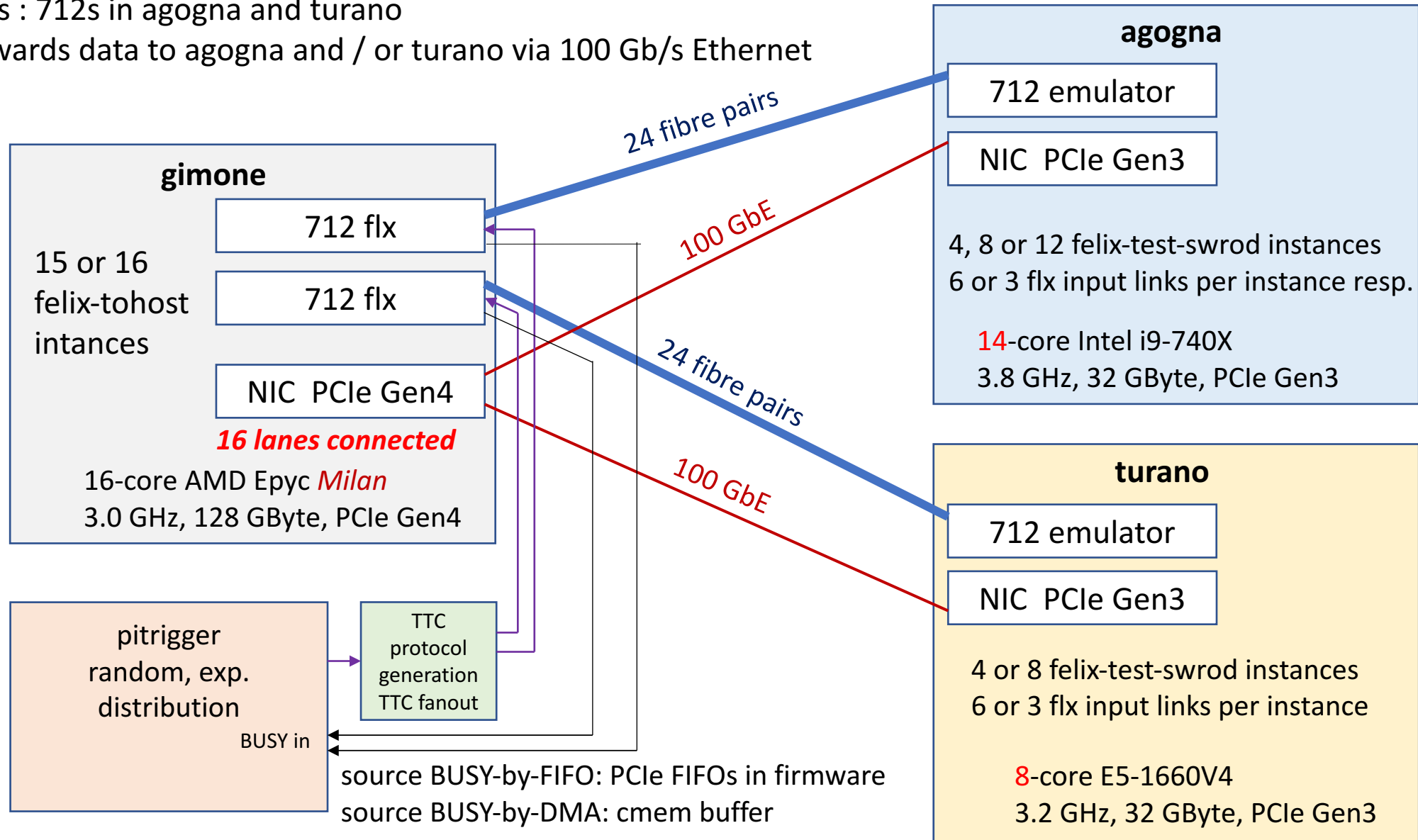# Test setup at Nikhef

FELIX server: gimone
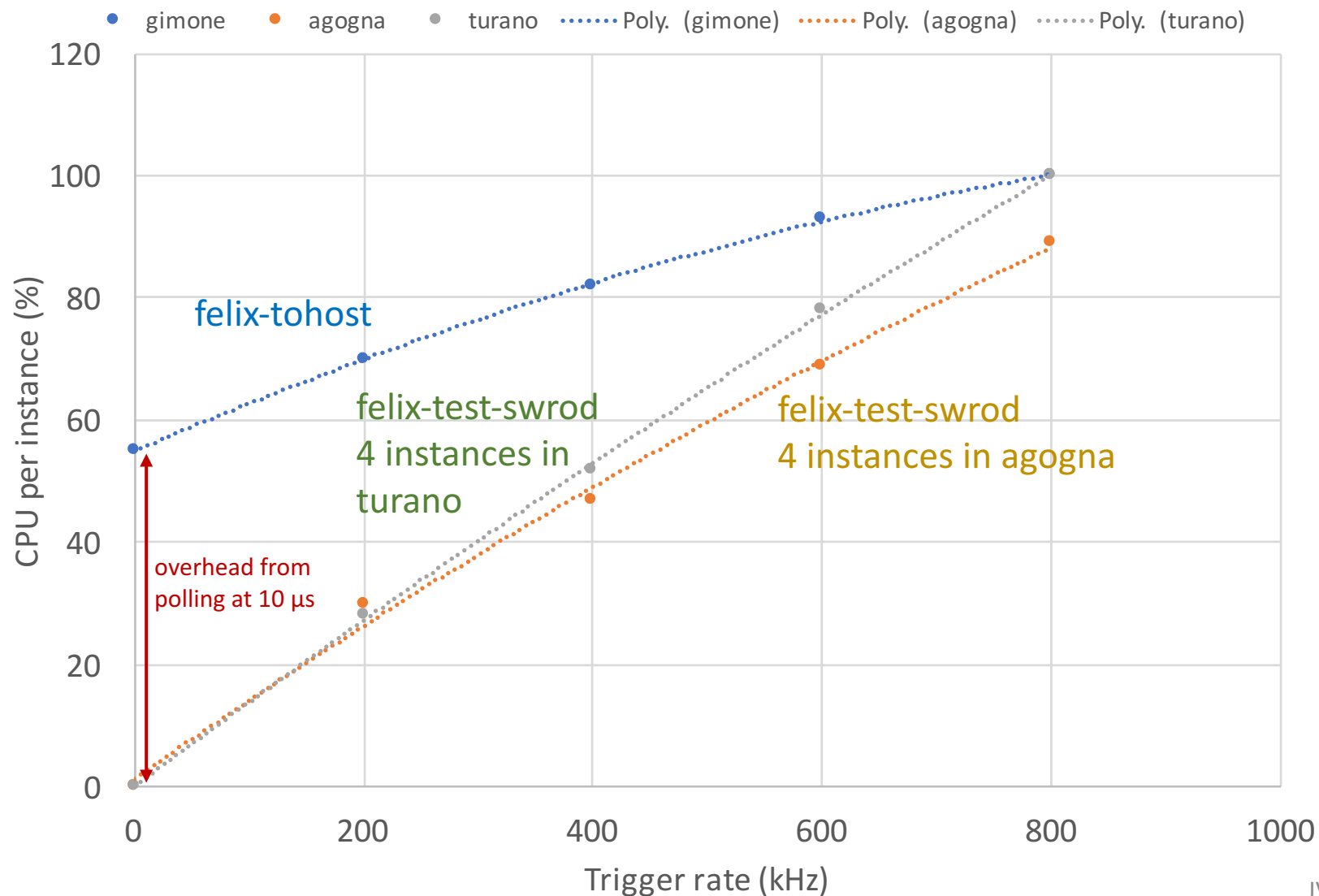data sources : 712s in agogna and turano
gimone forwards data to agogna and / or turano via 100 Gb/s Ethernet

**gimone**

712 flx

712 flx

NIC  PCIe Gen4

15 or 16 felix-tohost intances

*16 lanes connected*

16-core AMD Epyc *Milan*
3.0 GHz, 128 GByte, PCIe Gen4

24 fibre pairs

100 GbE

24 fibre pairs

100 GbE

**agogna**

712 emulator

NIC  PCIe Gen3

4, 8 or 12 felix-test-swrod instances
6 or 3 flx input links per instance resp.

14-core Intel i9-740X
3.8 GHz, 32 GByte, PCIe Gen3

**turano**

712 emulator

NIC  PCIe Gen3

4 or 8 felix-test-swrod instances
6 or 3 flx input links per instance

8-core E5-1660V4
3.2 GHz, 32 GByte, PCIe Gen3

pitrigger
random, exp.
distribution

TTC protocol generation TTC fanout

BUSY in

source BUSY-by-FIFO: PCIe FIFOs in firmware
source BUSY-by-DMA: cmem buffer

GBT mode, BUSY-by-FIFO off (no truncations), BUSY-by-DMA on, 0% dead time for rate ≤ 800 kHz for 15 felix-tohost instances

- eight 8-bit E-links per GBT link, 48 GBT links, 3 GBT links per DMA channel, 16 DMA channels in total, 20 bytes per trigger per E-link
- 16 felix-tohost instances, one per core: instable running, instance pinned to core 0 cannot handle the data
- 15 felix-tohost instances, one per core, servicing 45 GBT links: stable running if 4 felix-test-swrod instances on turano are pinned to core 1, 2, 3 and 4.
- Max. rate about 800 kHz: felix-tohost ~100% CPU, felix-test-swrod turano ~100% CPU, felix-test-swrod agogna ~88-90% CPU



felix-tohost

felix-test-swrod
4 instances in
turano

felix-test-swrod
4 instances in agogna

overhead from
polling at 10 µs

microcode.service masked
mitigations = off

# GBT mode, BUSY-by-FIFO off (no truncations), BUSY-by-DMA on, 0% dead time for rate ≤ 800 kHz for 15 felix-tohost instances

- eight 8-bit E-links per GBT link, 48 GBT links, 3 GBT links per DMA channel, 16 DMA channels in total, 20 bytes per trigger per E-link
- 16 felix-tohost instances, one per core: instable running, instance pinned to core 0 cannot handle the data
- 15 felix-tohost instances, one per core, servicing 45 GBT links: stable running if 4 felix-test-swrod instances on turano are pinned to core 1, 2, 3 and 4.
- Max. rate about 800 kHz: felix-tohost ~100% CPU, felix-test-swrod turano ~100% CPU, felix-test-swrod agogna ~88-90% CPU

Maximum rate GBT mode determined by:

- felix-to-host CPU usage per instance per instance 100% at 800 kHz, but polling interval can be made larger:
  - Polling interval 10 µs: 55% CPU usage for 0 kHz trigger rate
  - Polling interval 100 µs: 5.8 – 7.1% CPU usage for 0 kHz trigger rate
  - Polling interval 1 ms: 0.7% CPU usage for 0 kHz trigger rate

- 100% CPU not possible for 16[th] felix-to-host instance, but increase of polling interval length also reduces Linux overhead:
  - Polling interval 10 µs: /sbin/rngd –f: 18.8%
  - Polling interval 100 µs: /sbin/rngd –f: 4.2%
  - Polling interval 1 ms: /sbin/rngd –f: 0.7%

- felix-test-swrod runs on turano at 100% CPU at 800 kHz, handling data from 6 GBT links / 48 E-links
  - ➢ 16 instances and therefore 16 cores needed if 3 GBT links per felix-test-swrod instance are handled for reduction of CPU usage
  - ➢ turano has 8 cores and -> run at max. 7 instances, maybe 8 (if allowed by Linux overhead)
  - ➢ Tests done with 12 instances on agogna (14-core CPU) and 4 on turano
  - ➢ Could have tested with less instances on agogna and more on turano, would have resulted in more even bandwidth usage of both 100 GbE links (short test of this afternoon: 8 instances on agogna and 8 on turano: same result as for 12 instances on agogna and 4 on turano)

Result with 16 felix-to-host instances, 28 bytes per E-link per trigger for first 712, 24 bytes per E-link for second 712:
1 MHz rate just possible with 0% BUSY-by-DMA during 1 hour. 28 bytes per E-link for second 712 causes sometimes BUSY-by-DMA. felix-to-host for last link DMA channel of second 712 runs on core 0.
Bandwidth usage reported by felix-test-swrod : 55 Gb/s for 100 GbE link to agogna, 16 Gb/s for 100 GbE link to turano

FULL mode tests done in buffered mode, polling interval of 100 us
- Mellanox software used, previous observed low maximum rates not seen again
- BUSY-by-FIFO switched off, no truncations
- 4 felix-test-swrod instances on agogna and 4 on turano, each handling data from 6  flx input links
- felix-tohost instances running below ~ 45% CPU
- felix-test-swrod instances running at about 30% CPU
- Max. rate ~ 1050 kHz, 0% BUSY-by-DMA
- about 72 Gb/s via each 100 GbE link (PCIe Gen 4 interface!): network bandwidth seems to limit the rate, CPU usage felix-tohost not if polling interval is not too short (~100% at 10 us felix-tohost)

Two snapshots of cmem buffer (1 GB per felix-tohost instance) usage in bytes @ 1050 KHz

```
device 0 [bytes]:  660480       614400      1124352      1152000
device 1 [bytes]: 1513472      1003520       884736      1430528
device 2 [bytes]: 1539072       815104       947200       620544
device 3 [bytes]: 1230848       712704       535552      1564672

device 0 [bytes]: 1049600       992256      1505280      1529856
device 1 [bytes]:  828416       392704      1368064       728064
device 2 [bytes]:  742400      1252352      1355776       994304
device 3 [bytes]: 1229824       734208       619520       429056
```

Snapshot of cmem buffer (1 GB per felix-tohost instance) usage in bytes @ 1100 KHz after 30 s

```
device 0 [bytes]: 840698880    818032640    817379328    818610176
device 1 [bytes]: 841504768    816556032    816988160    817250304
device 2 [bytes]: 839153664    813294592    813891584    814288896
device 3 [bytes]: 841615360    814907392    814271488    814206976
```
-> BUSY-by-DMA (threshold @ about 800 MB)

Remarks:

- Earlier report: small dead time fraction at high rates seems to be difficult to avoid.
  - Dead time fraction observed caused by BUSY-by-FIFO, thresholds for setting and removing the BUSY were set, upon the basis of trial and error, to avoid truncations
  - Thresholds used effective for avoiding truncations, but cause a small dead-time fraction also at 1 MHz or lower rates
  - BUSY-by-FIFO has been disabled as no truncation occurs with the firmware currently used and if the firmware is properly functioning.
  - Resetting the firmware or rebooting the server is sometimes necessary, in particular for the GBT mode firmware, to get the firmware properly functioning

- The FELIG emulator consistently outputs a somewhat smaller number of event fragments than the number of triggers, as seen in the totals reported by fdaqm. For different GBT links the number of fragments may be somewhat different (of the order of 0.05%), but for the E-links of one GBT link they are the same, example on next slide

- Pinning of felix-tohost and felix-test-swrod instances to cores (with taskset) ineeded for best performance if CPU load of each instance approaches 100%. If the CPU load is less pinning may not be needed, in that case reduction of /proc/sys/kernel/sched_migration_cost_ns, as found for the ROS, could help. Some experimentation (GBT mode, 1 ms polling, 1 MHz rate) with setting it to 1000 instead of 50000 without core pinning to felix-tohost instances did not show a clear improvement, and peaks in cmem buffer filling appeared to be considerably higher than for pinning, caused sometimes BUSY-by-DMA

(From RHEL7 doc: /proc/sys/kernel/sched_migration_cost_ns specifies the amount of time after the last execution that a task is considered to be "cache hot" in migration decisions. Increasing this variable reduces task migrations. Adjust by factor of 2-10x. Task migrations may be irrelevant depending on any configured task affinity settings).

```
-> Data checked @Dev-DMA=0-0: Blocks 3458957
-> Elink chunks @Dev-DMA=0-0:
Elink Lnk-i  Chunks
----- -----  ---------
0x008 00-08  5013397
0x00c 00-12  5013397
0x010 00-16  5013397
0x014 00-20  5013397
0x018 00-24  5013397
0x01c 00-28  5013397
0x020 00-32  5013397
0x024 00-36  5013397
0x048 01-08  5021221
0x04c 01-12  5021221
0x050 01-16  5021221
0x054 01-20  5021221
0x058 01-24  5021221
0x05c 01-28  5021221
0x060 01-32  5021221
0x064 01-36  5021221
0x088 02-08  4956131
0x08c 02-12  4956131
0x090 02-16  4956131
0x094 02-20  4956131
0x098 02-24  4956131
0x09c 02-28  4956131
0x0a0 02-32  4956131
0x0a4 02-36  4956131
0x600 24-00  5023743
```

GBT link
with 8 E_links

TTC-to-host

pitrigger: 5040962

```
-> Data checked @Dev-DMA=0-0: Blocks 34262333
-> Elink chunks @Dev-DMA=0-0:
Elink Lnk-i  Chunks
----- -----  ---------
0x008 00-08  49648974
0x00c 00-12  49648974
0x010 00-16  49648974
0x014 00-20  49648974
0x018 00-24  49648974
0x01c 00-28  49648974
0x020 00-32  49648974
0x024 00-36  49648974
0x048 01-08  49724208
0x04c 01-12  49724208
0x050 01-16  49724208
0x054 01-20  49724208
0x058 01-24  49724208
0x05c 01-28  49724208
0x060 01-32  49724208
0x064 01-36  49724208
0x088 02-08  49118639
0x08c 02-12  49118639
0x090 02-16  49118639
0x094 02-20  49118639
0x098 02-24  49118639
0x09c 02-28  49118639
0x0a0 02-32  49118639
0x0a4 02-36  49118639
0x600 24-00  49748082
```

pitrigger: 49867493