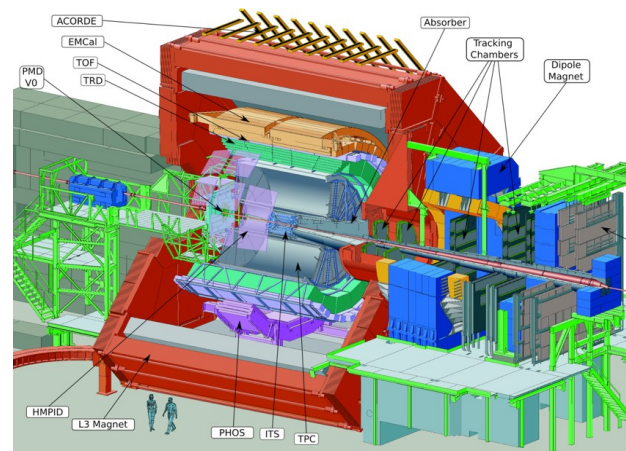
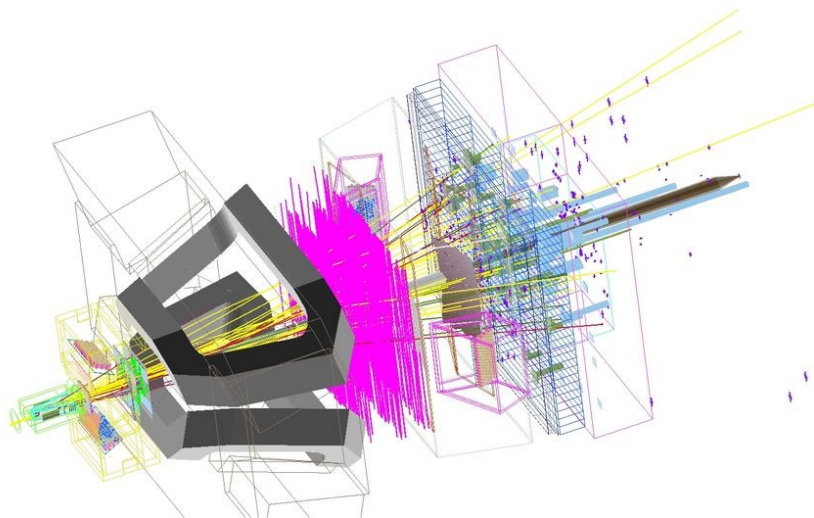




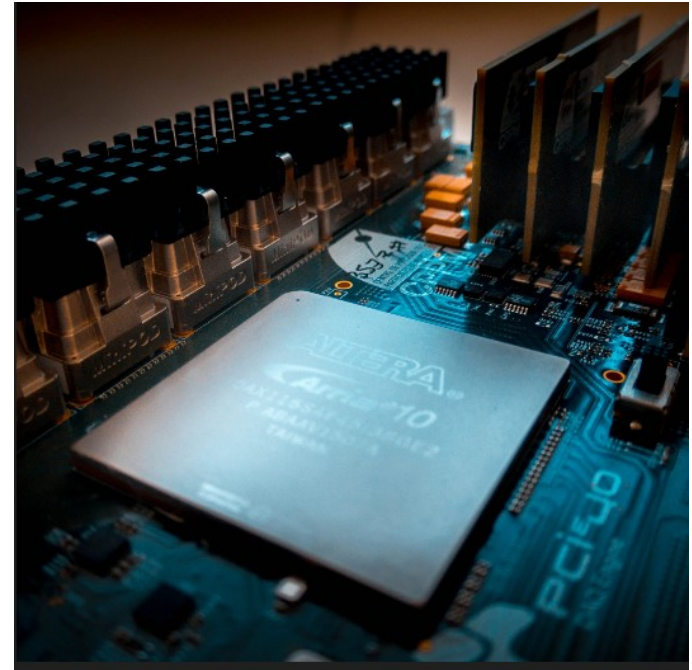
The PCIe40 card and the importance of efficient production tests



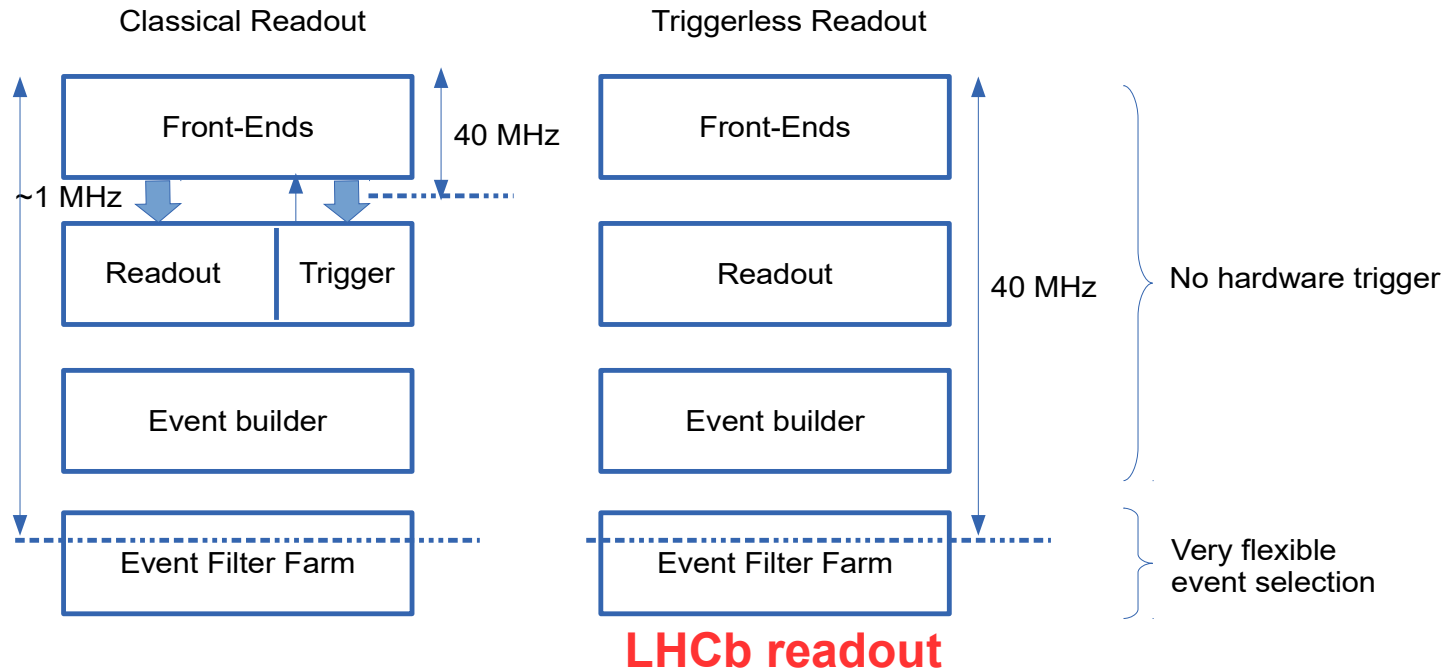
J.P. Cachemiche, on behalf of the LHCb collaboration

Outline

- The PCIe40 card
 - o LHCb and ALICE Readout architecture
 - o Card main features
 - o Measurements
 - o Production
- Testing to the limits



LHCb Upgrade key features

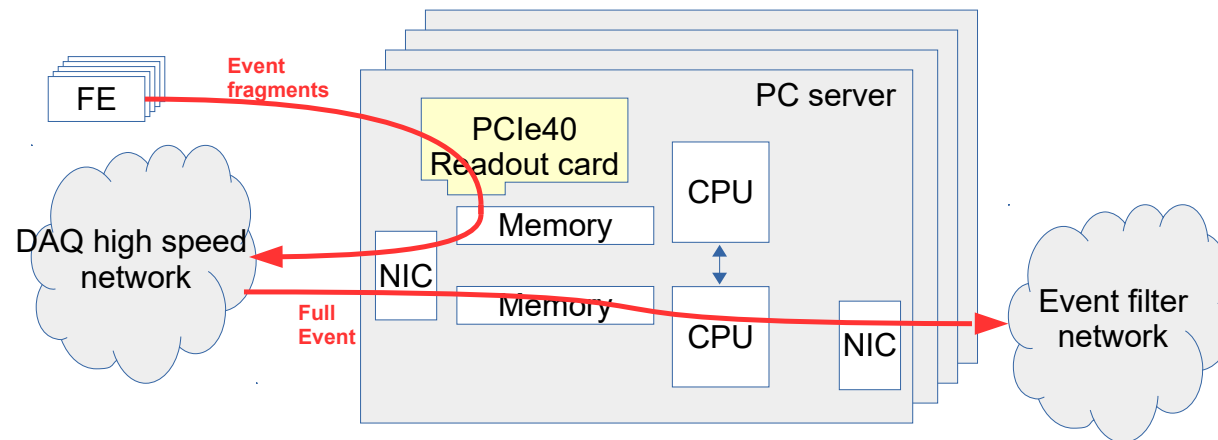


- LHCb uses a **triggerless readout**
- All event fragments routed at 40 MHz up to the farm

LHCb Upgrade key features

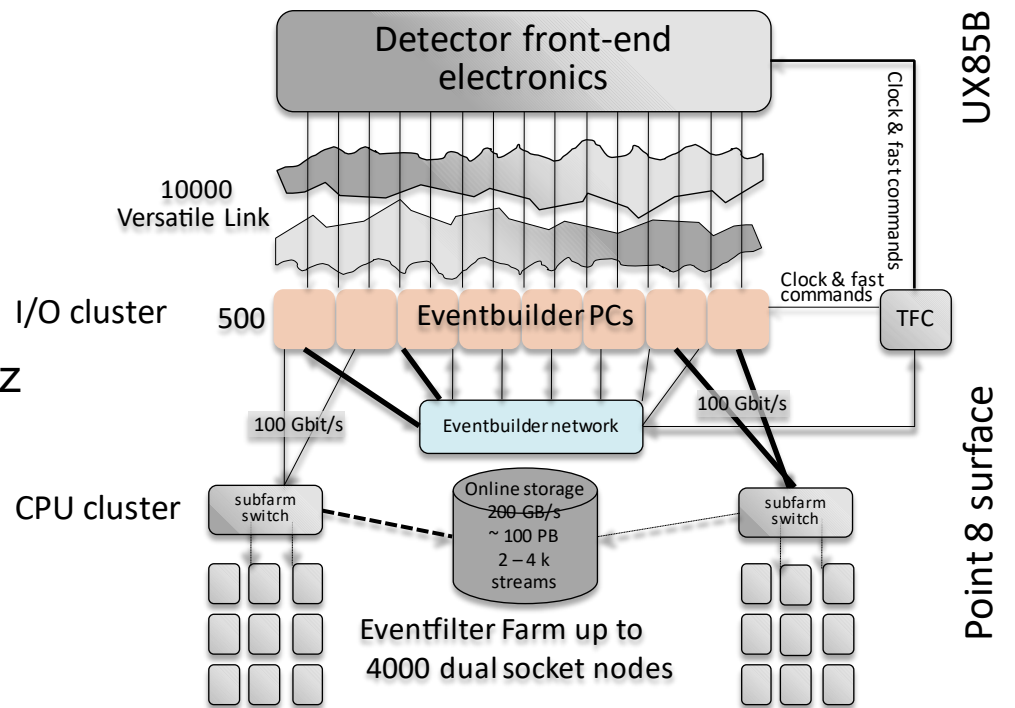
Principle

- Event building done by tightly coupled acquisition boards, CPUs and high speed network
- No intermediate back-end stage
 - ➔ Readout card implemented as a PCIe module
- Event building through servers in real time
 - ➔ Now possible due to internal CPU architecture evolution
- Event reconstruction **with offline quality** in real time
- Triggering replaced by **filtering of reconstructed events**



LHCb architecture

- Readout located on surface
 - o Distance between FE and RO : ~350m
- ~ 10000 optical links
- ~ 500 readout boards
- ~ 50 TFC/ECS cards
- ~ 100 kBytes per event at 40 MHz
- ~ 32 Tb/s aggregate bandwidth
- ~ 4000 dual CPU nodes



Alice upgrade key features

- Event topology too complex for electronics trigger
- 60% of events are kept
 - ➔ Continuous triggerless readout + Low interaction rate (50 kHz)
- CRU (Common Readout Unit) based on the PCIe40 card
- Acquires and compresses data on the fly

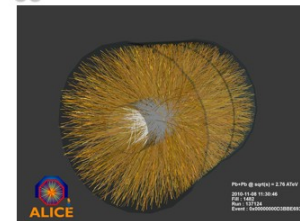
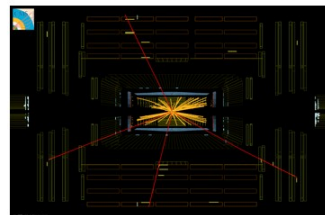


At present (Run1 & 2)

- Interaction rate 8 kHz (Not all LHC bunches have collisions) → max. trigger rate < 3.5 kHz

Why low interaction rate?

- Event topology too complex for simple electronics triggers



3 TB/s data in Run 3

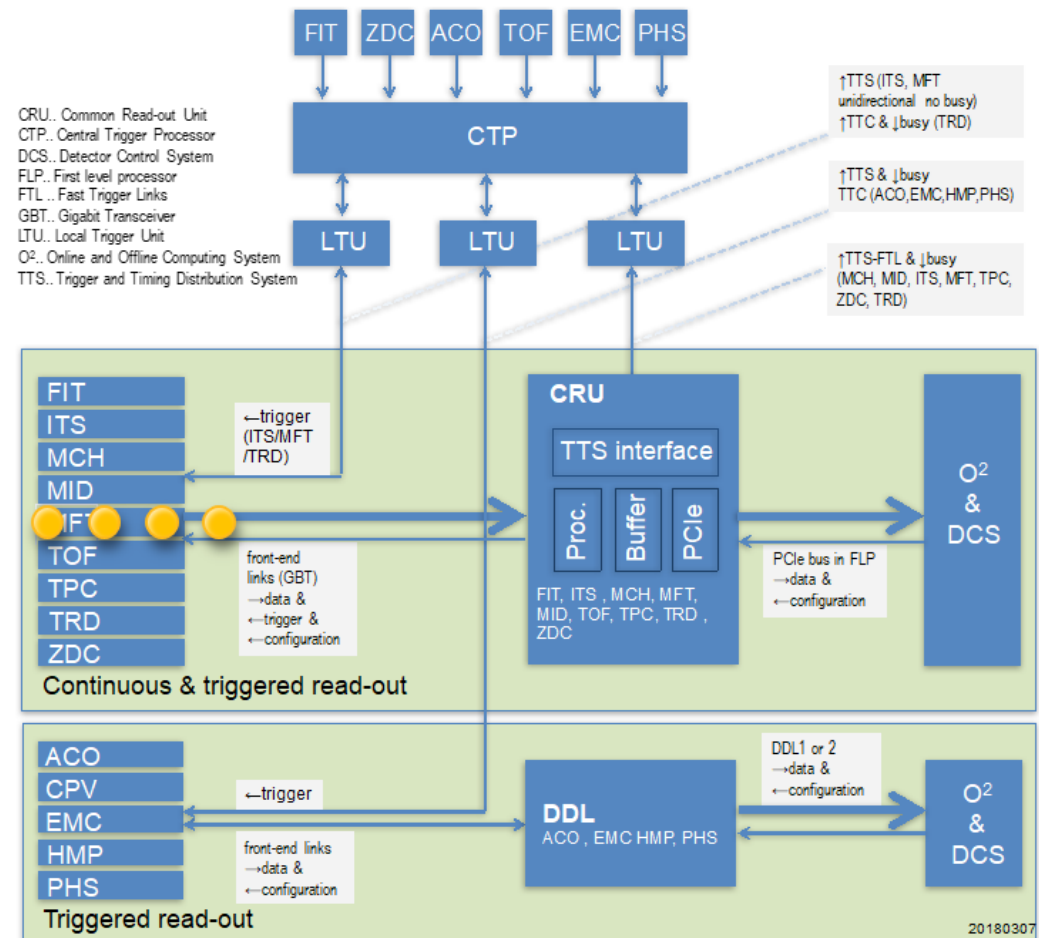
After upgrade (≥ Run 3)

- Target
 - Pb-Pb $\geq 10 \text{ nb}^{-1}$ → 9×10^{10} events
 - pp (@5.5 TeV) $\geq 6 \text{ pb}^{-1}$ → 1.4×10^{11} events
 - Gain factor 100 in statistics
- Interaction rate 50 kHz (PbPb) → continuous triggerless read-out

Courtesy Alex Kluge

ALICE architecture

- Readout located on surface
 - o Distance between FE and RO : ~120m
- ~ 9000 optical links
- ~ 540 readout boards
- ~ 68 MBytes per event at 50 KHz
- ~ 27 Tb/s aggregate bandwidth
- ~ 1500 GPU based event processing nodes

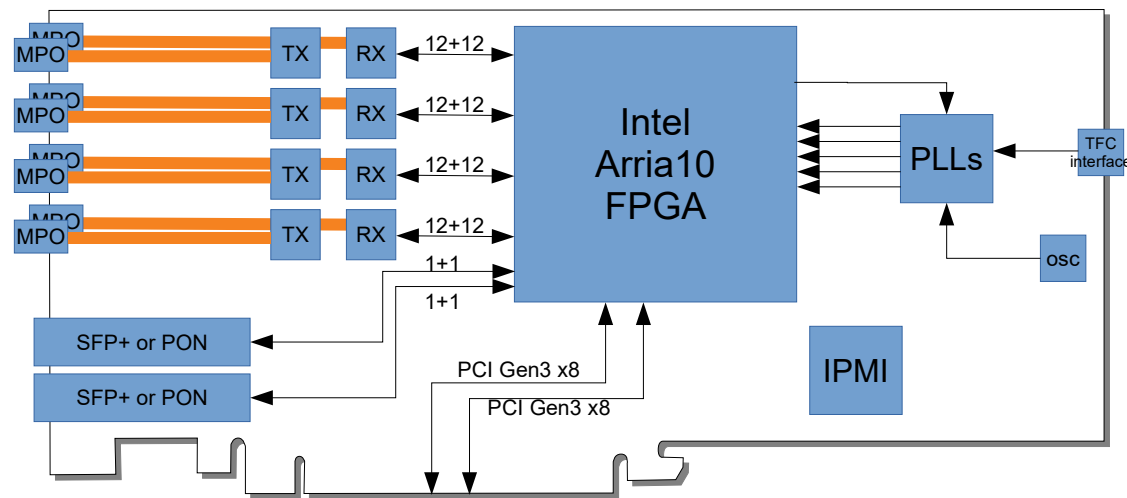


Courtesy Alex Kluge

The readout board : PCIe40

- Features :

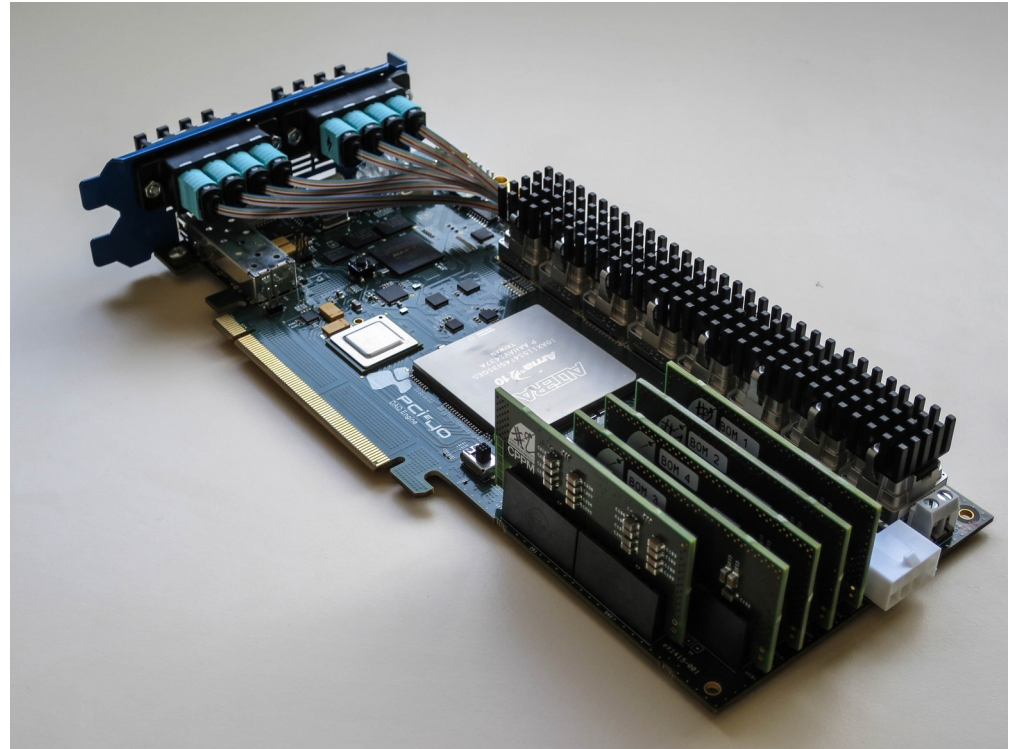
- 1 large FPGA 1.15 million cells (Arria10 10AX115S3F45E2SG)
- 48 bidirectional links running at up to 10 Gbits/s each (minipods)
- 2 bidirectional links running at up to 10 Gbits/s devoted to time distribution (can use SFP+ or 10G PON devices)
- Sustained 112 Gbits/s interface with CPU memory through PCIe
- No on-board buffer memory : we use the PC memory instead
- Remote reconfiguration of all the programmable devices
- Fully instrumented: all voltages, currents and temperatures measured



Hardware design

PCle40 prototype

- First prototype developed in 2016
- 24 copies manufactured for both the LHCb and Alice collaboration
 - o Used as « mini DAQ » for debugging front-end cards
 - o Programmed to provide acquisition, ECS and TFC in a single firmware



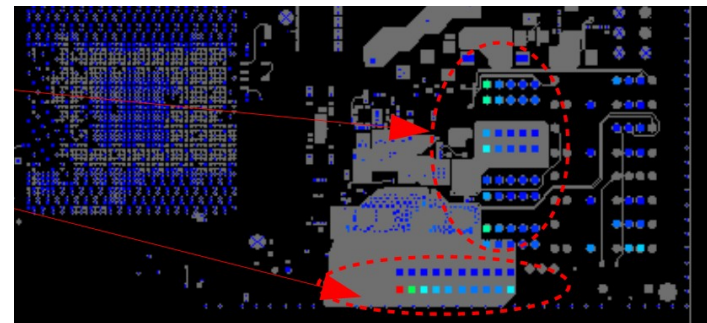
Preparing the final module

Power consumption of large FPGAs very high

- Up to **52 A** on the core !
- Power consumption
 - o FPGA estimated at $\sim 80 \text{ W}$
 - o Card estimated at $\sim 150 \text{ W}$ with Engineering Sample
 - o Limited thickness for the stackup

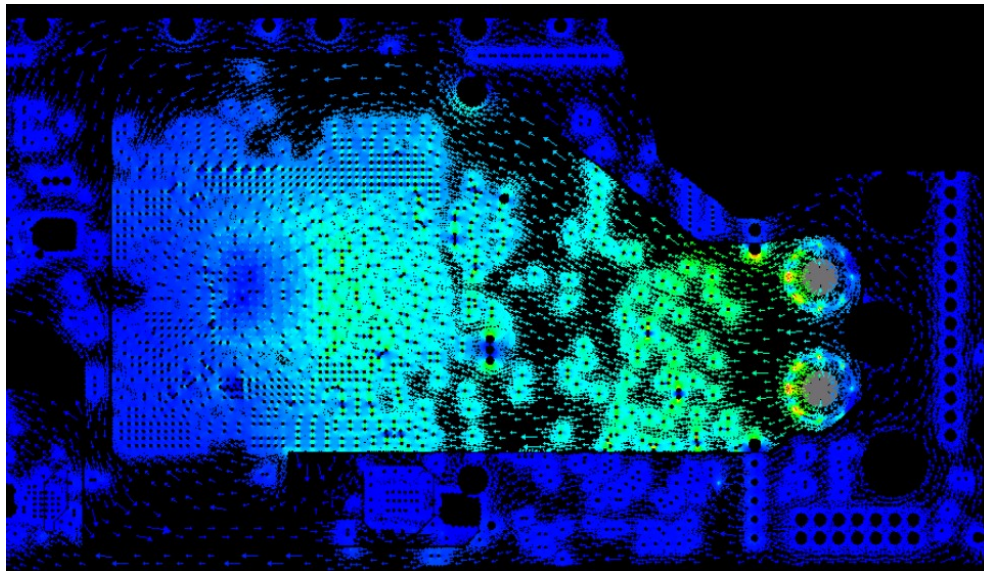
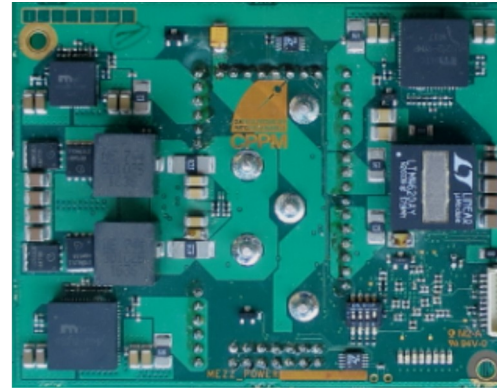
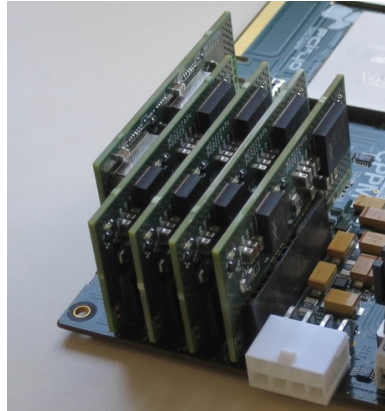
Refining of current flow simulations

- Simulations of current flow showed dangerous hot spots at full load
 - ➡ Power planes have been redesigned and vias placement has been optimized
- Current flow through power mezzanine connections not symmetric

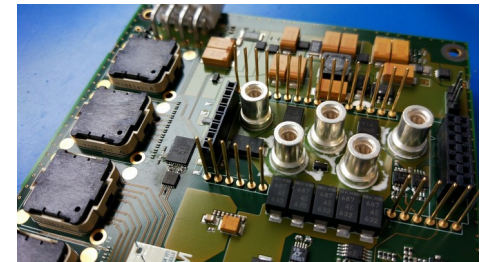


Preparing the final module

Replacement of the 5 vertical mezzanines by a single flat one



Current flow between mezzanine and FPGA with new design



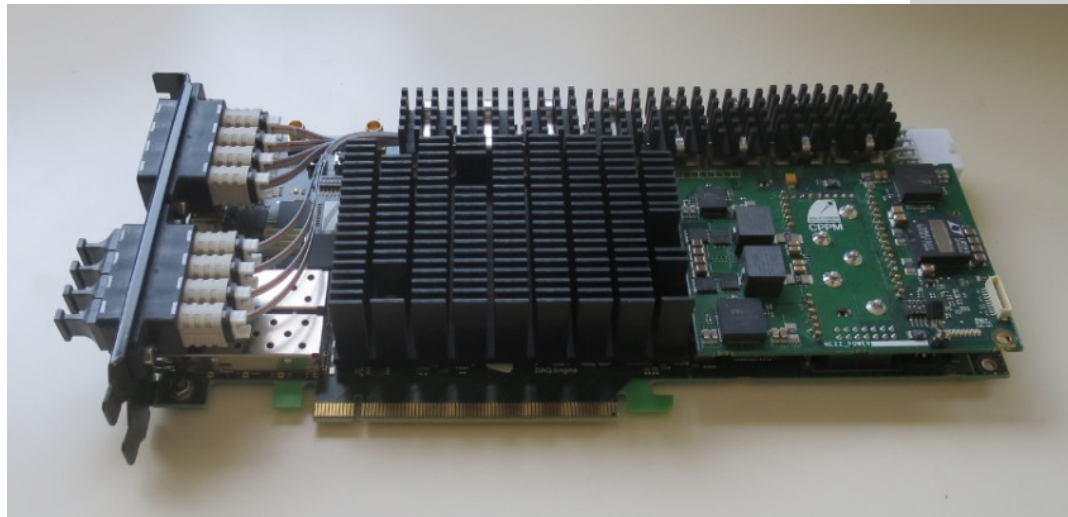
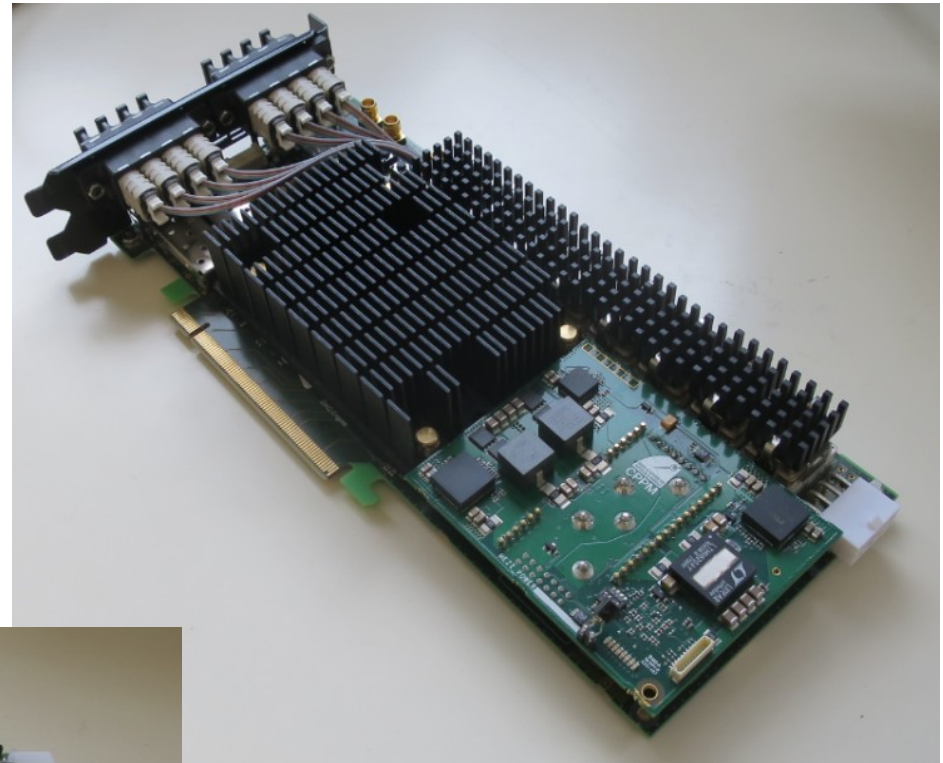
Optimizations

Many improvements

- Cost savings
 - Removal of expensive components (PCIe bridge, Serial Flash and corresponding power supply)
 - One additional SFP+ or PON cage added → less TFC/ECS modules
- Performance improvement
 - Use of new PLLs with a very low jitter compared to previous ones
- Reliability
 - Complete redesign of the power supply due to buggy DCDC converters
 - Optimisation of current flows → avoids local over-heating in the PCB
→ Single power mezzanine now horizontal for symmetrical current flow
 - Improvement of power sequencing to ease maintenance and guaranty a longevity of the module → manages now power down
 - Optimization of decoupling → less noise
 - Heat sink redesign for better cooling
- New functionalities
 - Programming speed multiplied by factor 4 with a new embedded USB Blaster II
 - Serial flash for identificating modules during production
 - IPMI management : allows the system to adjust the fan speed in function of the temperature or automatically cut the power supply if temperature is too high

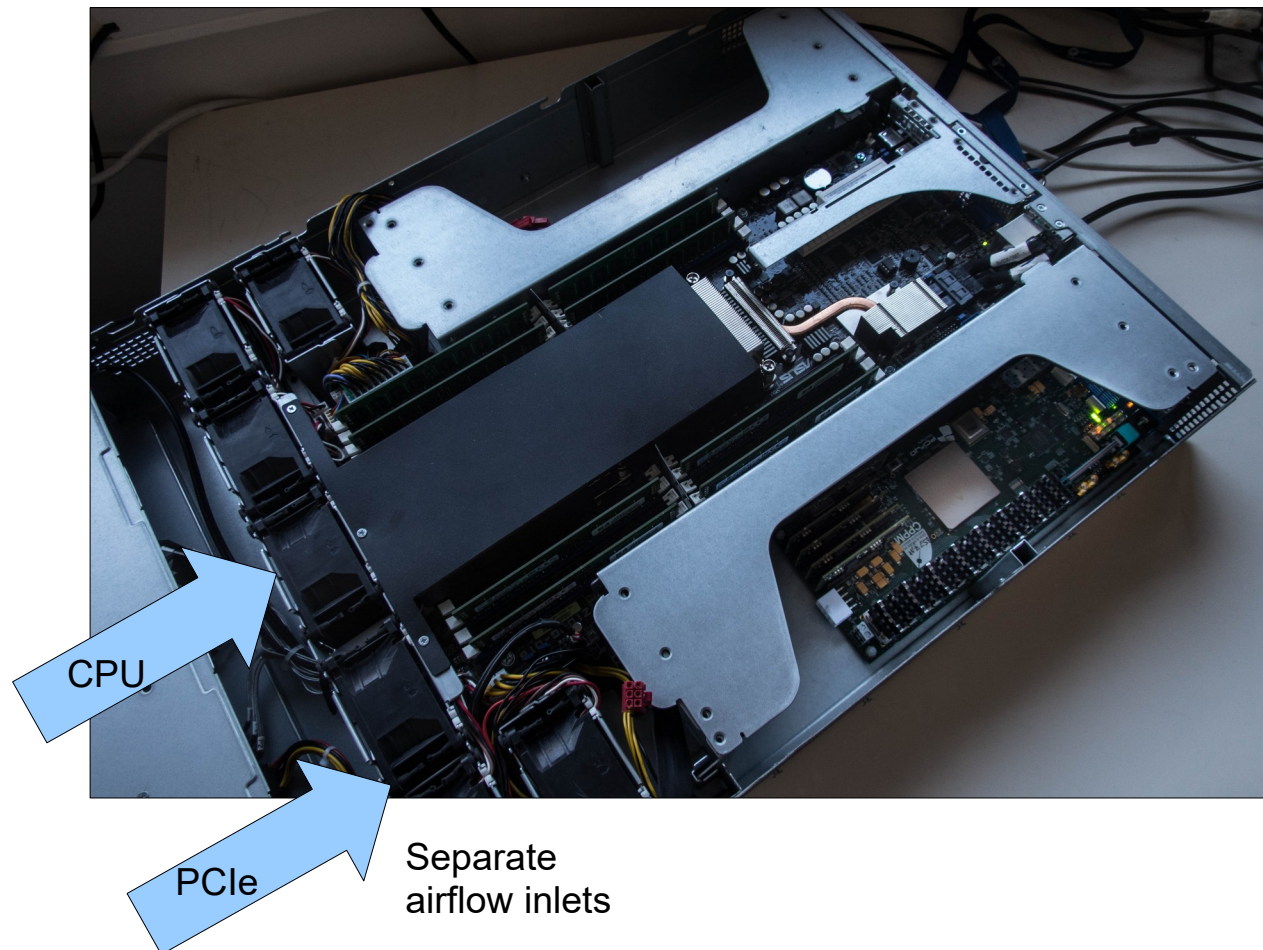
Final module

- Two first modules validated end 2017
- Early duplication by Alice of 28 modules to speed up first production



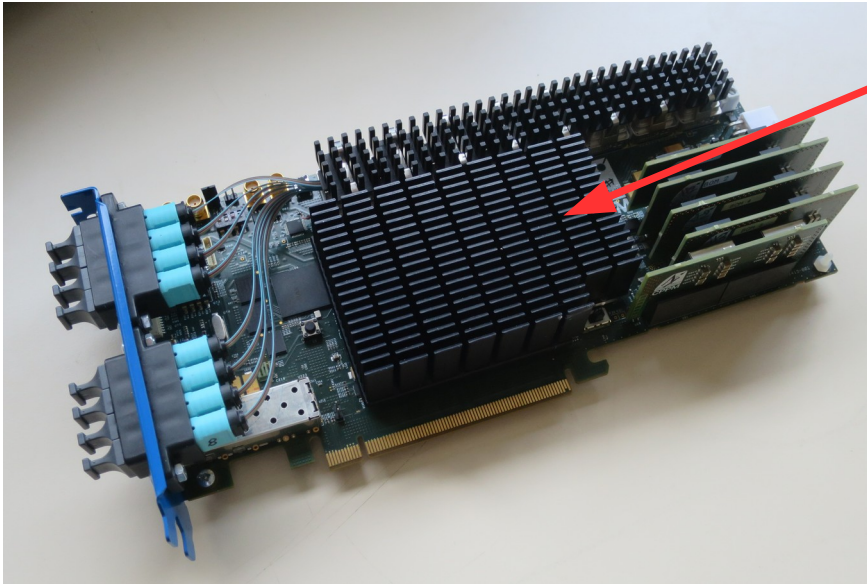
Cooling

- PC environment not as well defined as xTCA systems
- Very well cooled PC server has been selected



Cooling solution

Use of a custom passive cooling



Custom passive heatsink

Power consumption and cooling

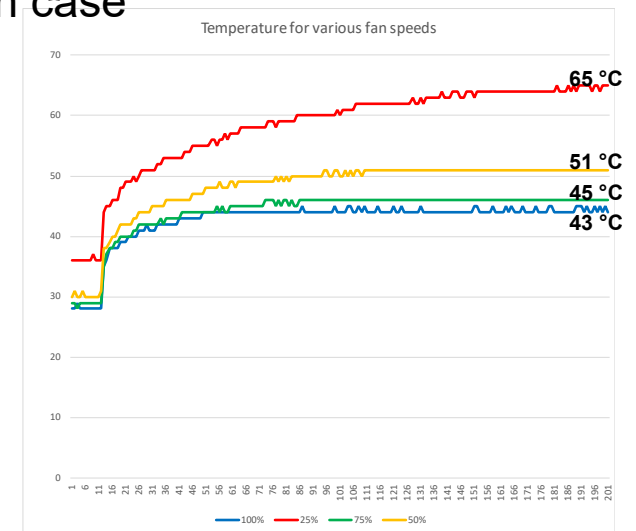
Power consumption and cooling

- Push the module at the limit of power dissipation
- Principle:
 - Use a « heating function» replicated thousands of times to get an FPGA occupancy of 89%
 - Inject a clock with programmable frequency between 10 MHz and 600 MHz
- Automatic power off if the FPGA temperature overpasses 82°C
- Vary the speed of server fans (25%, 50%, 75%, 100%)
- Measure voltages, currents and temperature in each case



Results obtained with ASUS server

- 2 cards on same side
- Passive cooling seems sufficient



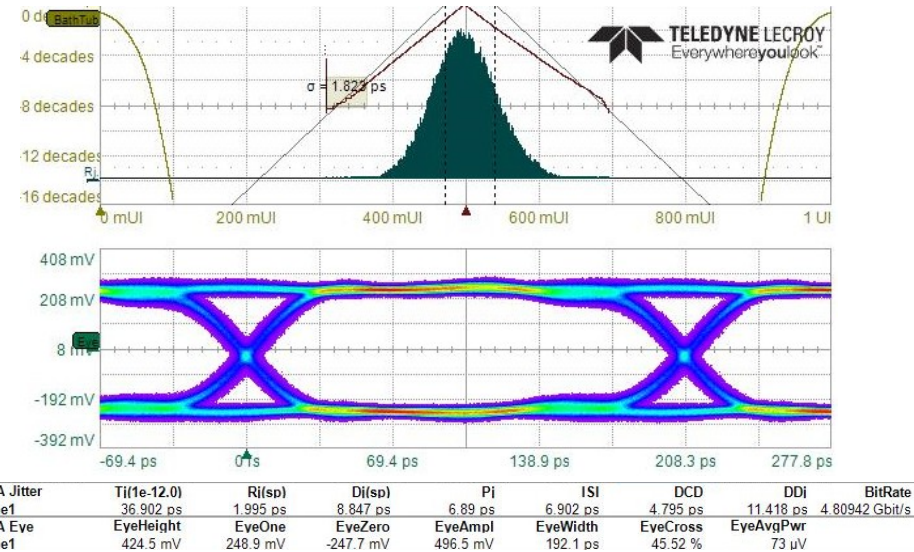
FPGA temperature for several fan speeds in ASUS server

Links measurements

BER << 10⁻¹⁶

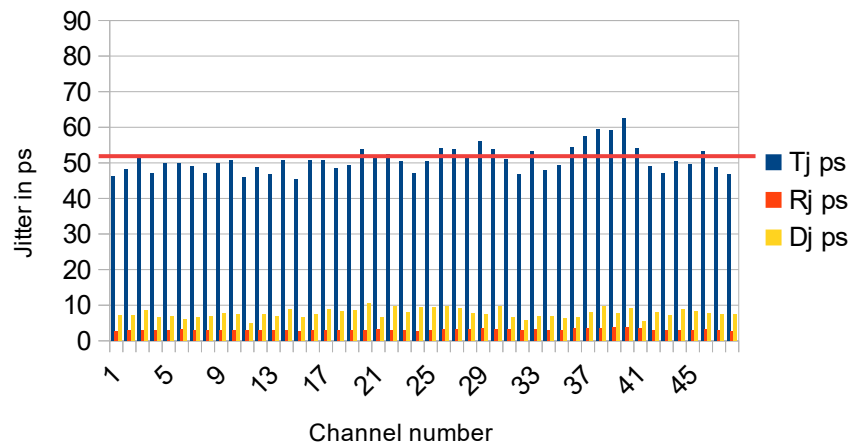
Jitter

- Final card jitter improved vs prototype
Total jitter goes from 51 ps → 38 ps



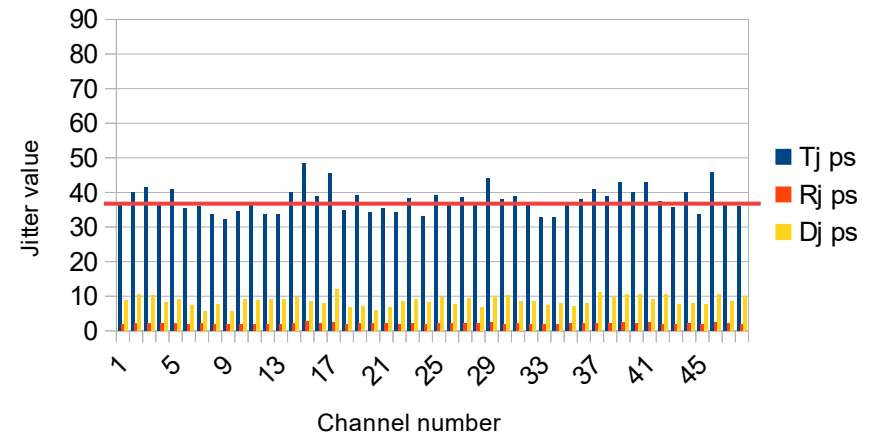
**Measurements at reception stage
for a PRBS31 pattern
running at 4.8 Gbits/s**

Jitter measurement over 48 links



Prototype

Jitter measurement over 48 links

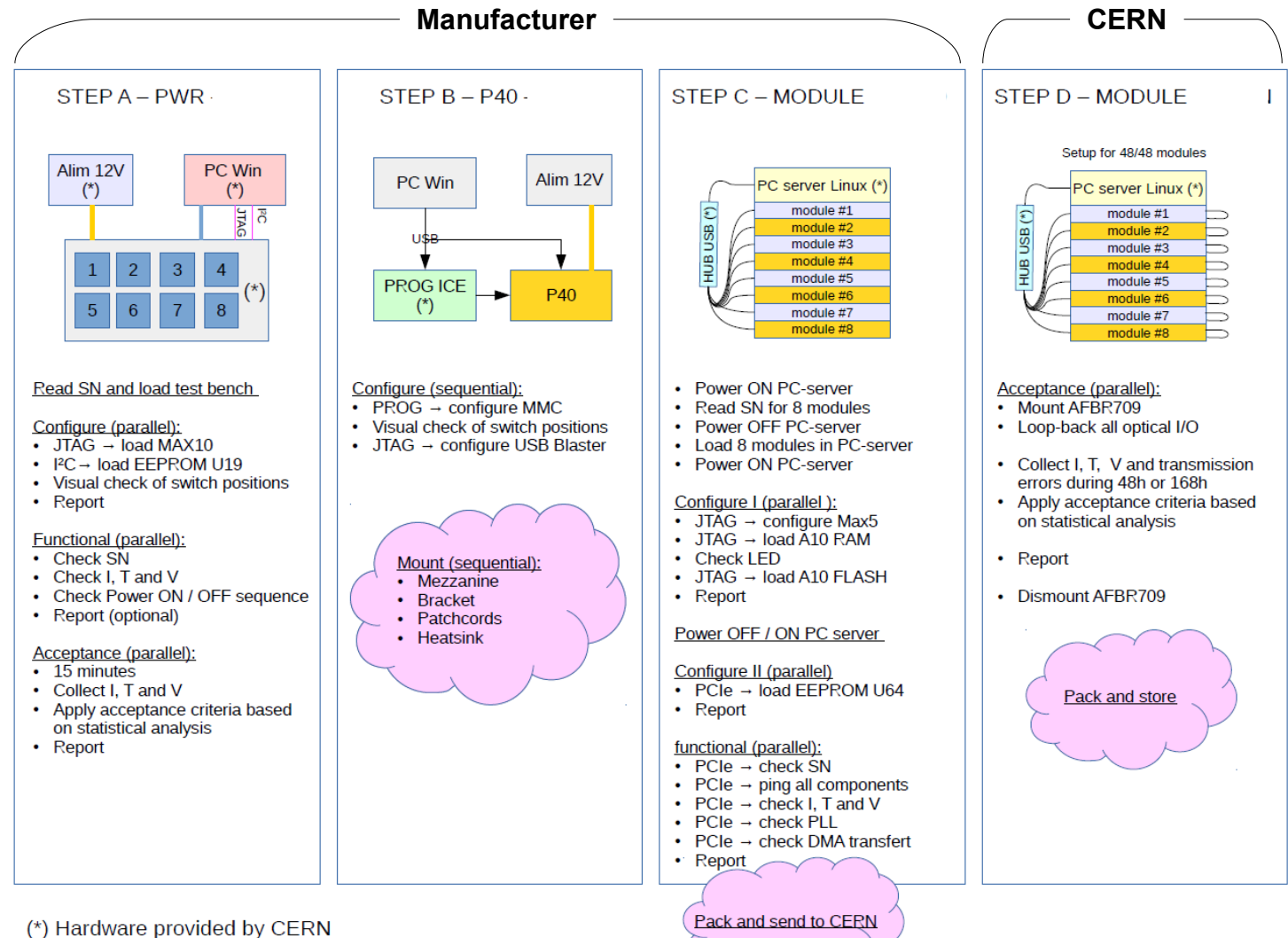


Final card

Production

Testing methodology

4 steps



Production tests

Run in assembly company

- **Based on Pytest**
 - Very flexible command line testing tool
 - Able to test target sub-set of components
 - Object oriented design
 - Can be driven by a GUI or ncurses
- **Fully tests the board**
 - ~ 146 unitary tests ran in a few minutes on 8 cards at a time
 - Check the operation of all the devices on the modules
 - Measure voltages, currents, temperatures, frequencies, etc.
 - Produces test reports for each module
- **Centralized management of reports**
 - Reports directly sent to CERN data base

```
[upgrade@marupgrade10 p40_functional]$ pytest
===== test session starts =====
platform linux2 -- Python 2.7.14, pytest-3.3.2, py-1.5.2, pluggy-0.6.0 -- /shared-PCIE40/Miniconda2/bin/python
cachedir: .cache
rootdir: /shared-PCIE40/PYD_FOR_V2/LLT_PClE40_devices/SCRIPTS_FC0/TOOLS/p40_functional, infile: pytest.ini
plugins: profiling-1.2.11, hypothesis-3.38.5
collected 109 items

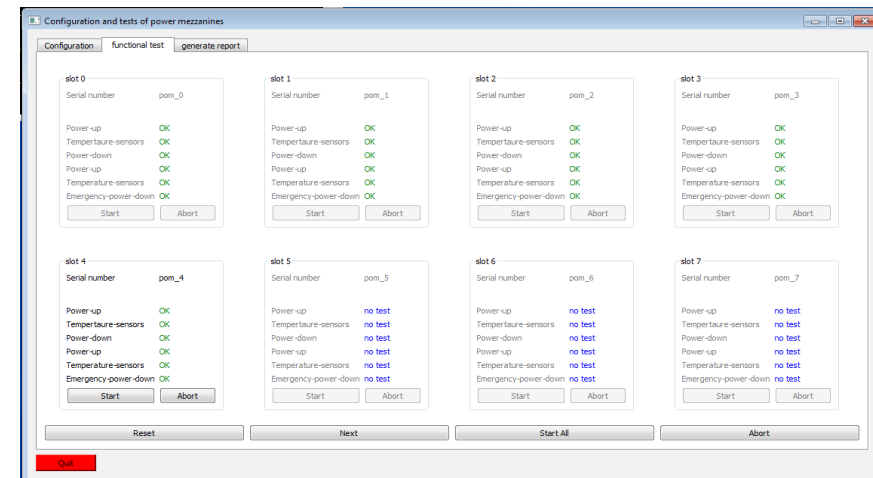
test_01_base.py::test_arria10_u1_ping_pcie_50101 PASSED
test_01_base.py::test_arria10_u1_ping_gen3_50102[0] PASSED
test_01_base.py::test_arria10_u1_ping_gen3_50102[1] PASSED
test_01_base.py::test_max1619_u16_ping_50104 FAILED
test_01_base.py::test_si5344_u54_ping_50105 PASSED
test_01_base.py::test_si5345_u23_ping_50106 PASSED
test_01_base.py::test_si5345_u48_ping_50107 PASSED
test_01_base.py::test_minipod_ping_50108[mipd0] SKIPPED
test_01_base.py::test_minipod_config_50109 ERROR
test_01_base.py::test_si53154_u11_ping_50110 PASSED
test_01_base.py::test_afbr709_ping_50111[u19] FAILED
test_05_io.py::test_afbr709_tx_fault_50508[u19] FAILED
test_05_io.py::test_afbr709_rx_loss_50509[u19] FAILED
test_05_io.py::test_afbr709_data_ready_50510[u19] FAILED
test_01_base.py::test_afbr709_ping_50111[u219] FAILED
test_05_io.py::test_afbr709_tx_fault_50508[u219] FAILED
test_05_io.py::test_afbr709_rx_loss_50509[u219] FAILED
test_05_io.py::test_afbr709_data_ready_50510[u219] FAILED
test_01_base.py::test_eeeprom_pwr_u19_part_number_50112 ERROR
test_01_base.py::test_eeeprom_u64_part_number_50113 ERROR
test_02_pll.py::test_si5344_u54_program_50201 ^C

KeyboardInterrupt

to show a full traceback on KeyboardInterrupt use --fulltrace
/shared-PCIE40/PYD_FOR_V2/LLT_PClE40_devices/FC0/devices_ll1/components/si534x_comp.py:336: KeyboardInterrupt

===== 9 failed, 7 passed, 1 skipped, 3 error in 12.03 seconds =====
[upgrade@marupgrade10 p40_functional]$
```

Expert interface



Example of operator interface

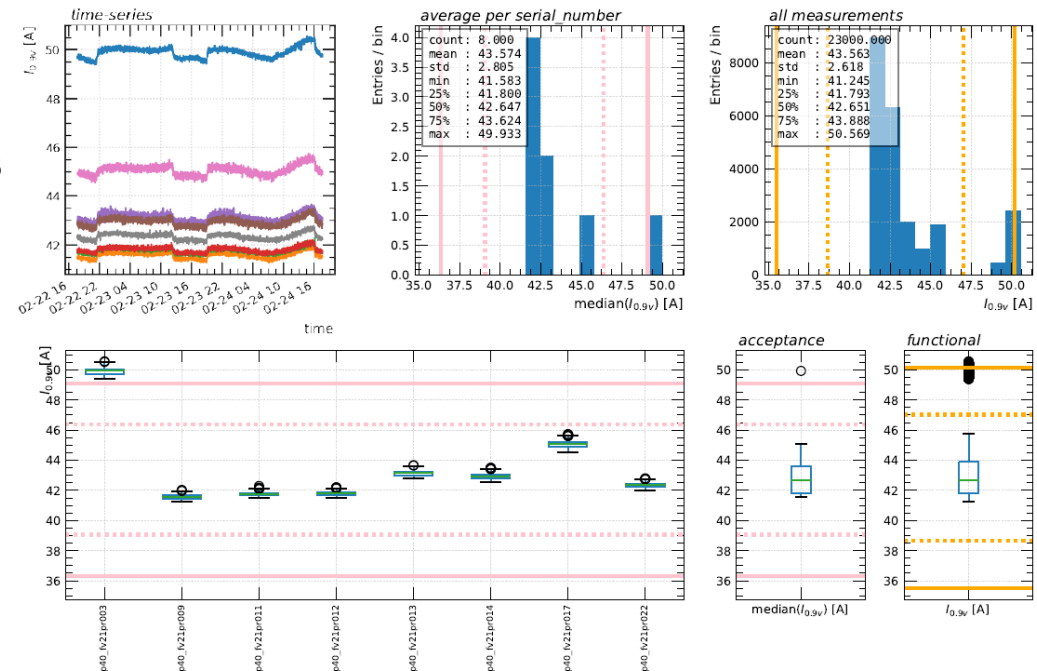
Acceptance tests

Run at CERN

- Duration 24 or 168 hours
Allow to eliminate **early failures**
- Rely on Pytest
- Possible post processing of results
- Logged in data bases

Raw $I_{0.9V}$ -- 190222 cern 300mhz 8fv21 2d0h1m

Feb 25, 2019 15:18
day 056 -- week 08



FPGA core current measurement on 8 cards

Testing to the limits



Not everything is pink

Which firmware for testing ?

Target design

- Up to 100% occupancy
- Average clock 250 MHz
- Average toggle rate : 50%

How to test the design at maximum load ?

- Final firmware was not ready (will it be one day ?)
- Since then there is a preliminary one but, very difficult to handle
 - o Requires WIN-CC
 - o Complicate initialization
 - o Not scriptable
- Fixed configuration

We decided to emulate it

- Not a perfect emulation because many unknown features
- But scalable design allowing to explore the limits set when designing the card

Firmware emulation

LLI (Low Level Interface) + programmable load emulation

- $n \times 80$ blocs of 16 random pattern generators

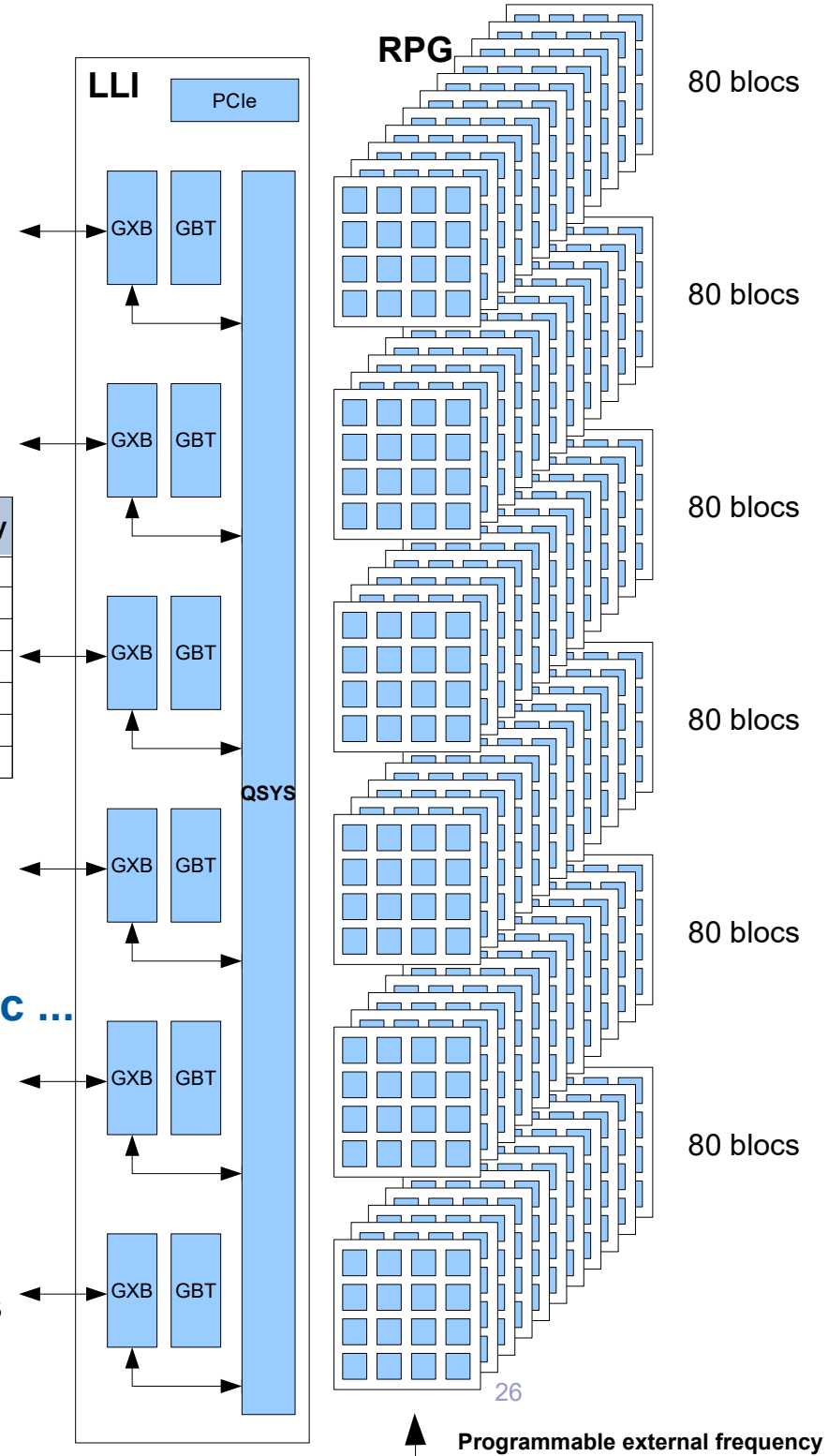
Number of Macro blocks	Number of blocks	Block size	Individual RPG size	Total number of RPG	FPGA occupancy
0	80	16	128 bits	0	14%
1	80	16	128 bits	1280	27%
2	80	16	128 bits	2560	39%
3	80	16	128 bits	3840	52%
4	80	16	128 bits	5120	65%
5	80	16	128 bits	6400	78%
6	80	16	128 bits	7680	89%

- Programmable frequency injected in random pattern generators
 - ➔ SI5344 PLL embedded on the card

Initial goal: checking power supply, cooling, etc ...

Extended to checking errors at full load

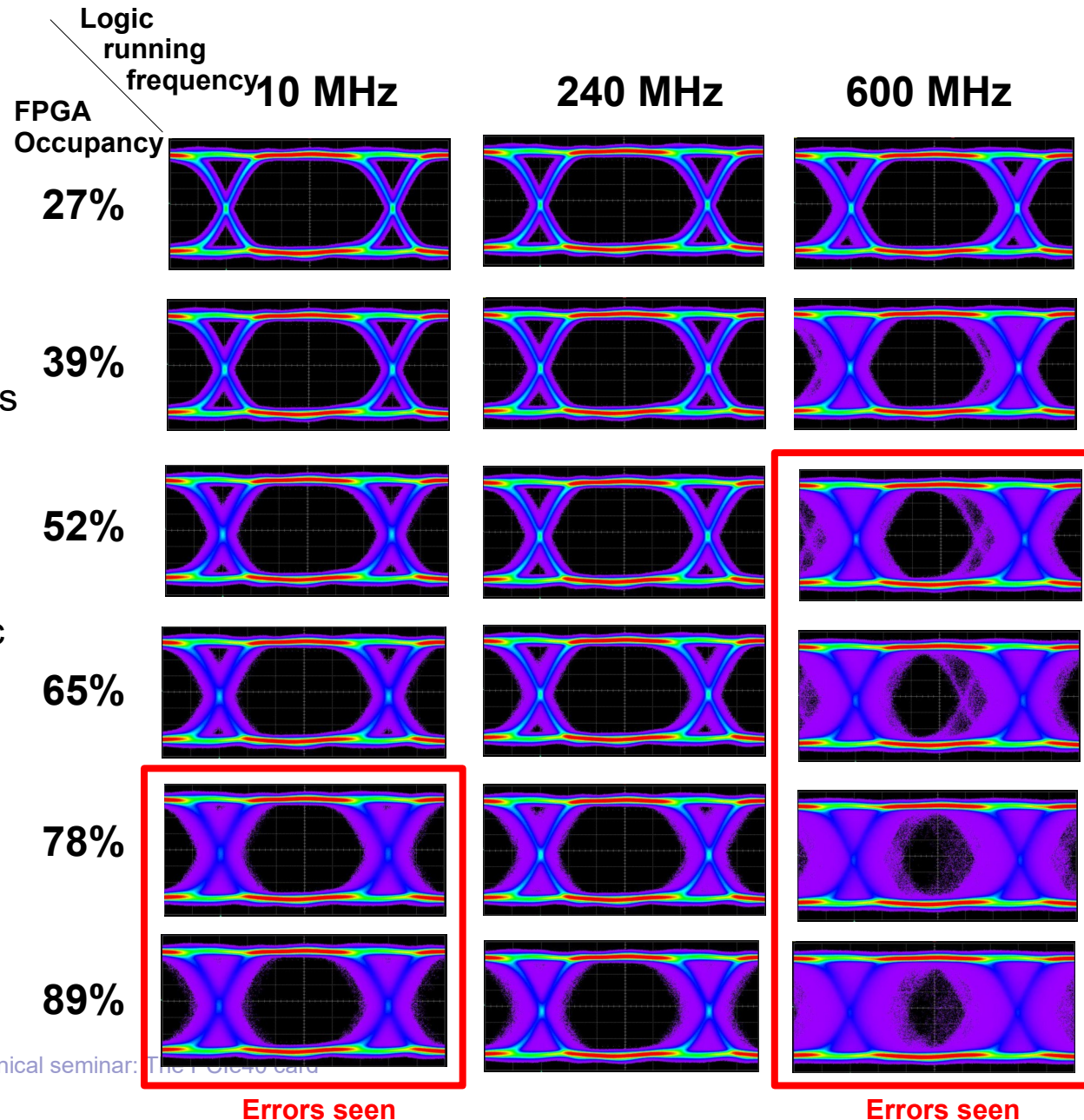
- BER tests made with GBX internal PRBS generators and checkers
 - ➔ TTK-like test by addressing GXB registers by software



Error checking vs occupancy

Conditions

- 48 links at 4.8 Gbits/s
- PRBS31
- Internal loopback mode
- Load emulator made of pseudo-random generators
- Increasing occupancy by varying the number of random generators
- Programmable clock frequency for internal logic
 - Design runs at 10 240 or 600 MHz



Unexpected errors seen at high occupancy!

Probable cause

VCCR/VCCT plane and VCC plane proximity

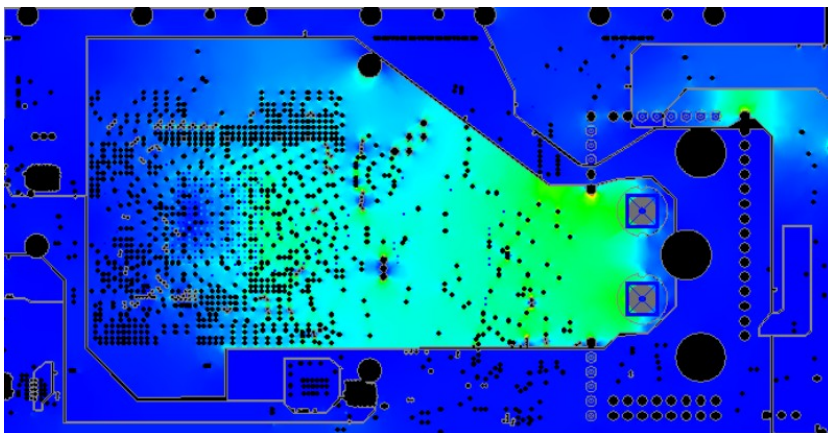
- Both capacitive and inductive effects

PCIe40V2						
Layer	FPGA				Thickness	
					um	mills
TOP	Diff Sig			Diff Sig	55	2,17
					75	2,95
L2	GND				40	1,57
					90	3,54
L3	Diff Sig			Diff Sig	17	0,67
					75	2,95
L4	GND				35	1,38
					85	3,35
L5	Sig			Sig	17	0,67
					75	2,95
L6	VCCR VCCT				35	1,38
					80	3,15
L7	VCC, VCCP, VCCE, RAM				70	2,76
					75	2,95
L8	GND				70	2,76
					80	3,15
L9	VCCPT, 1.8V, VCCH, VCC, PLL				35	1,38
					75	2,95
L10	Sig			Sig	17	0,67
					85	3,35
L11	GND				35	1,38
					75	2,95
L12	Diff Sig			Diff Sig	17	0,67
					90	3,54
L13	GND				40	1,57
					75	2,95
BOTTOM	Diff Sig			Diff Sig	55	2,17
					1573	61,93

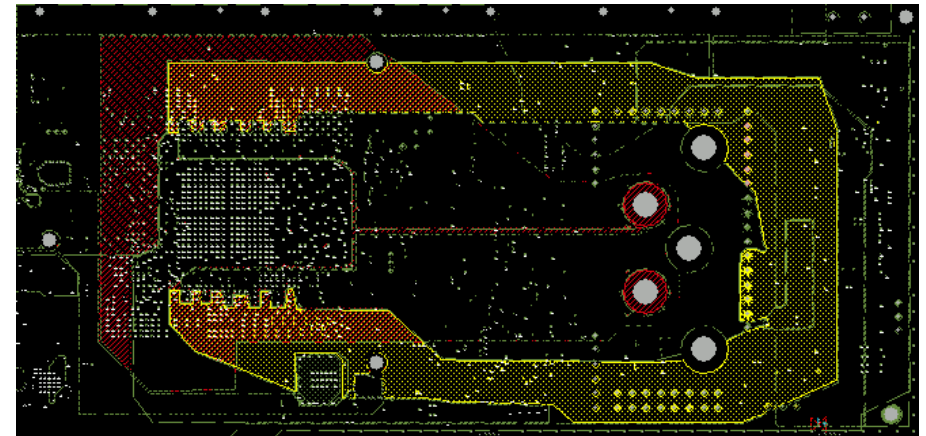
Probable cause

Overlap mostly of VCCT plane and VCC

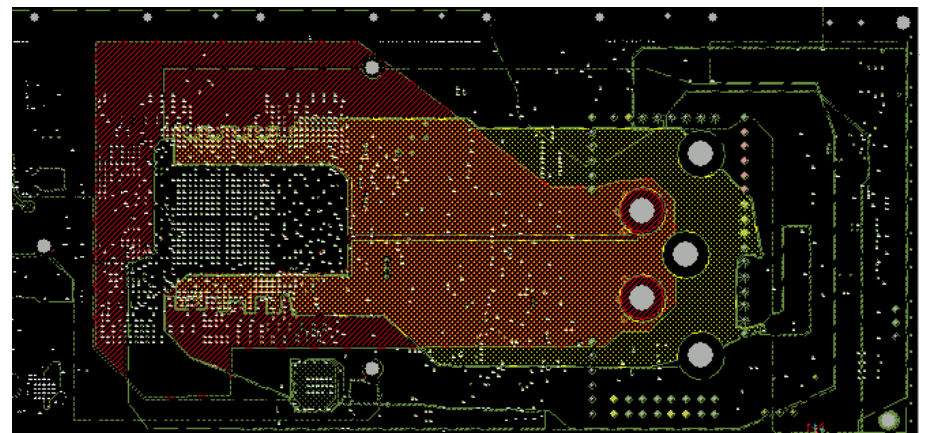
- Partial overlap between VCCR and VCC
 - ➡ Weak because in an area with nearly no current
- Large overlap between VCCT and VCC
 - ➡ Strong : high currents here



Current density in VCC (0.9V)



VCCT overlap (orange)



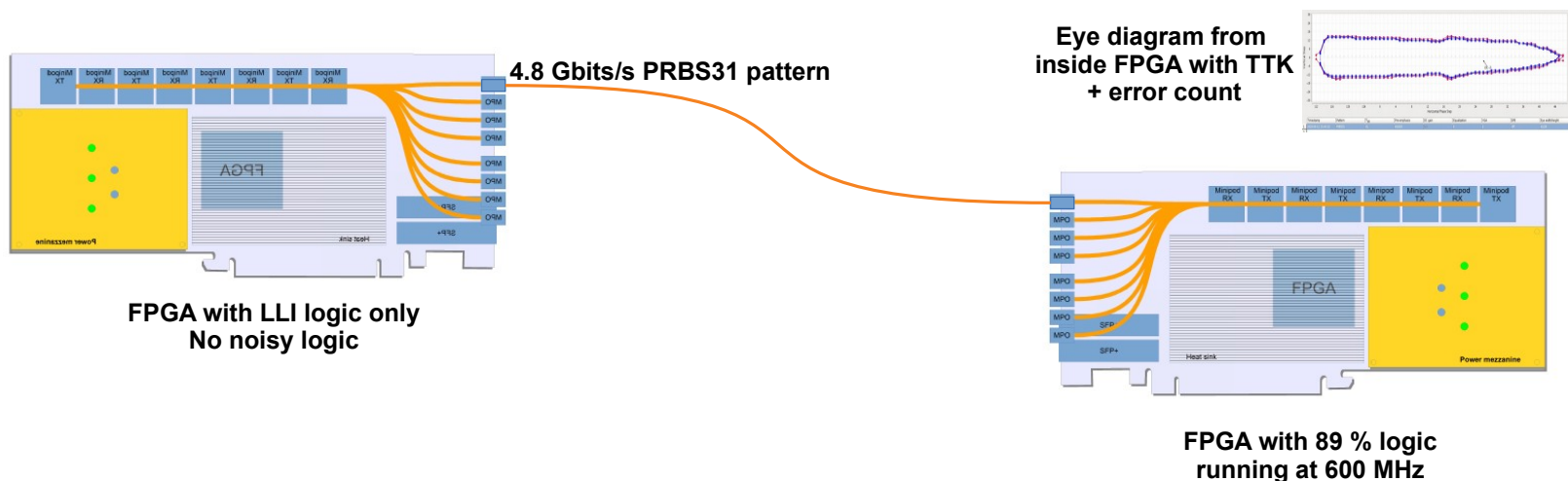
VCCR overlap (orange)

Checking hypothesis

If weak overlap between VCCR and VCC receiving side should not be affected

Verification

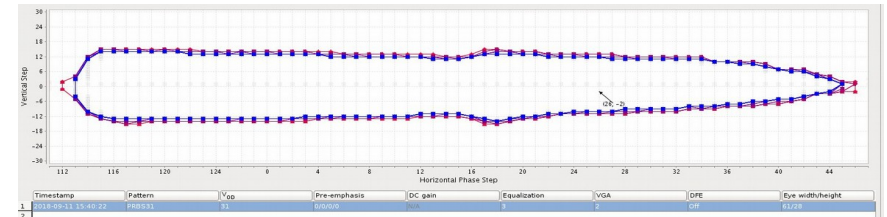
- Use of two cards:
 - o Emitting card not loaded → emission of a quiet signal
 - o Receiving card fully loaded (89%) + injection of frequencies 10, 240 and 600 MHz
- Check error count with the Transceiver Tool Kit in Quartus
- Check eye diagram from inside of receiving FPGA



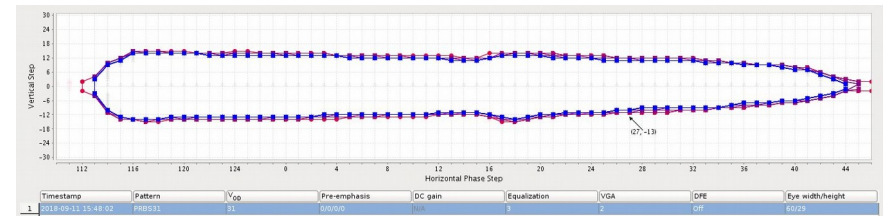
Receiving side

Measurement results

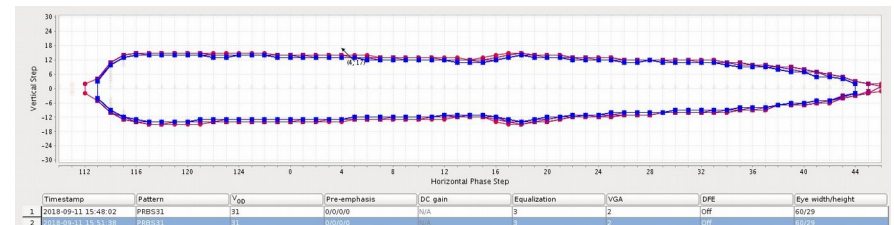
- ➡ **No error, even at full load**
(89% occupancy, 600 MHz)
- ➡ **No degradation of eye diagram**



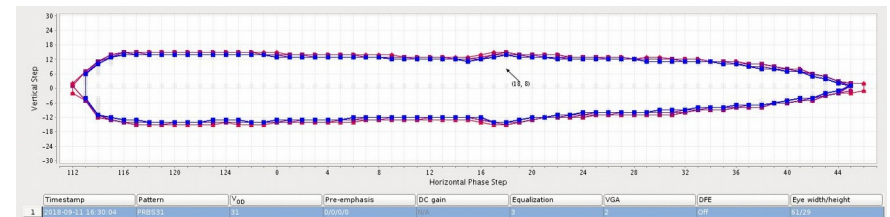
Reference measurement : LLI only - no load



89% occupancy – running frequency = 10 MHz



89% occupancy – running frequency = 240 MHz



89% occupancy – running frequency = 600 MHz

Corrective action

Invert power and ground planes

- 3 possibilities
- chosen 2 and 3

PCIE40V2					
Layer	FPGA				Thickness
					um mills
TOP	Diff Sig		Diff Sig		55 2,17
					75 2,95
L2	GND				40 1,57
					90 3,54
L3	Diff Sig		Diff Sig		17 0,67
					75 2,95
L4	GND				35 1,38
					85 3,35
L5	Sig		Sig		17 0,67
					75 2,95
L6	VCCR VCCT				35 1,38
					80 3,15
L7	GND				70 2,76
					75 2,95
L8	VCC, VCCP, VCCE_RAM				70 2,76
					80 3,15
L9	VCCPT,1.8V,VCCH,VCC_PLL				35 1,38
					75 2,95
L10	Sig		Sig		17 0,67
					85 3,35
L11	GND				35 1,38
					75 2,95
L12	Diff Sig		Diff Sig		17 0,67
					90 3,54
L13	GND				40 1,57
					75 2,95
BOTTOM	Diff Sig		Diff Sig		55 2,17
					1573 61,93

- ① Invert GND and VCCT/VCCR
Inconvenient : coupling between VCC (0.9) and VCCPT/1.8V/H

PCIE40V2					
Layer	FPGA				Thickness
					um mills
TOP	Diff Sig		Diff Sig		55 2,17
					75 2,95
L2	GND				40 1,57
					90 3,54
L3	Diff Sig		Diff Sig		17 0,67
					75 2,95
L4	GND				35 1,38
					85 3,35
L5	Sig		Sig		17 0,67
					75 2,95
L6	VCC_PLL VCCR VCCT				35 1,38
					80 3,15
L7	GND				70 2,76
					75 2,95
L8	VCC, VCCP, VCCE_RAM				70 2,76
					80 3,15
L9	GND				35 1,38
					75 2,95
L10	Sig		Sig		17 0,67
					85 3,35
L11	VCCPT,1.8V,VCCH				35 1,38
					75 2,95
L12	Diff Sig		Diff Sig		17 0,67
					90 3,54
L13	GND				40 1,57
					75 2,95
BOTTOM	Diff Sig		Diff Sig		55 2,17
					1573 61,93

- ② Invert VCC (0.9V) and VCCT/VCCR
Invert bottom GND and VCCPT/1.8V/VCCH
Inconvenient : diff signals over cut power plane

PCIE40V2					
Layer	FPGA				Thickness
					um mills
TOP	Diff Sig		Diff Sig		55 2,17
					75 2,95
L2	GND				40 1,57
					90 3,54
L3	Diff Sig		Diff Sig		17 0,67
					75 2,95
L4	VCC_PLL VCCR VCCT				35 1,38
					85 3,35
L5	Sig		Sig		17 0,67
					75 2,95
L6	GND				35 1,38
					80 3,15
L7	GND				70 2,76
					75 2,95
L8	VCC, VCCP, VCCE_RAM				70 2,76
					80 3,15
L9	GND				35 1,38
					75 2,95
L10	Sig		Sig		17 0,67
					85 3,35
L11	VCCPT,1.8V,VCCH				35 1,38
					75 2,95
L12	Diff Sig		Diff Sig		17 0,67
					90 3,54
L13	GND				40 1,57
					75 2,95
BOTTOM	Diff Sig		Diff Sig		55 2,17
					1573 61,93

- ③ Invert VCC (0.9V) and VCCT/VCCR
Invert top GND and VCCT/VCCR and
Invert bottom GND and VCCPT/1.8V/VCCH
Inconvenient : more diff signals over cut power plane

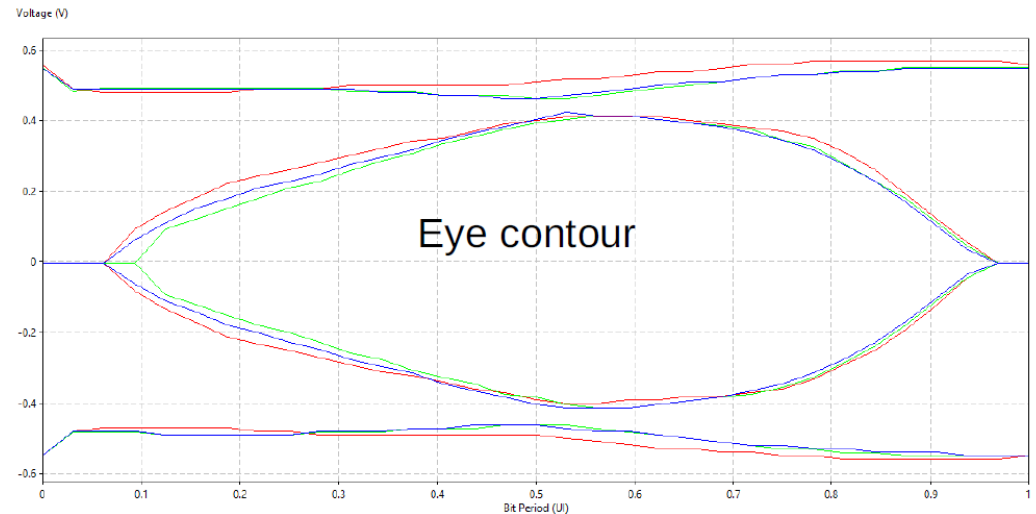
Diff signals simulation

Simulation to mitigate the risk of signal integrity issue if solution 2 and 3 are chosen

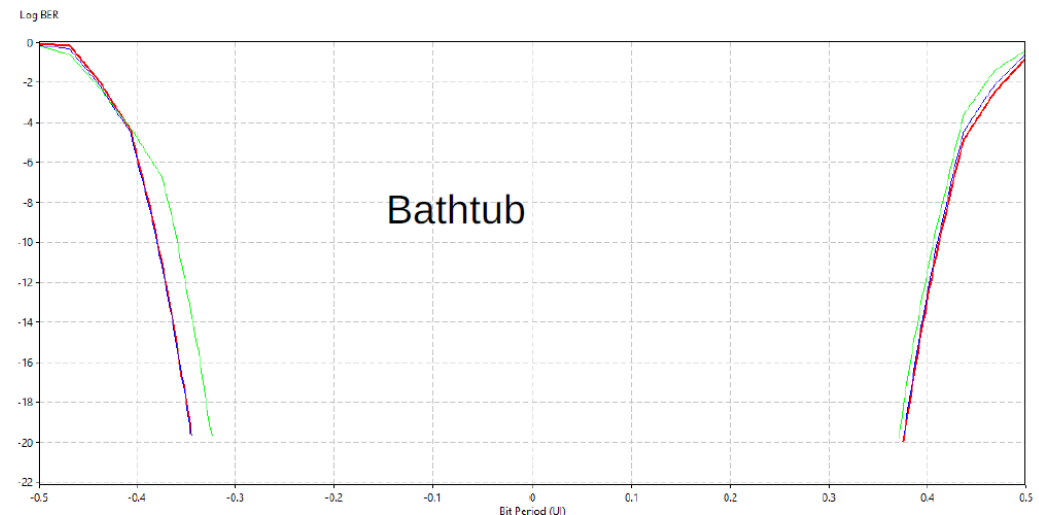
3 cases tested with same diff track:

- **Initial**: diff signal between 2 continuous GND planes
- **Simple plane inversion**: diff signal between continuous GND plane and power plane (5 cuttings)
- **Improved plane inversion**: power plane rearranged (2 cuttings only)

→ Negligible loss



10 Gbps simulation results with Sigrity



Error measurements

Building of tools to measure errors at the 3 maximum occupancy rates

- Conditions:
 - o Serial links: PRBS31 at 4.8 Gbits/s
 - o Test duration 100s per frequency
 - o Internal serial loop back
 - o Measurement of total number of errors on all 48 links
- **There exists a quite wide domain of operation without error**
- Correct operation range decreases with FPGA occupancy

	Frequency	Time	Average temperature	Average current	Errors found		Frequency	Time	Average temperature	Average current	Errors found		Frequency	Time	Average temperature	Average current	Error found
SOL40	10	100	51	5693	0		10	100	45	5323	188961960		10	100	51	5905	432194662
	20	100	48	5769	1083694		20	100	45	5708	1182867725		20	100	47	6050	127539359177
	40	100	47	6395	0		40	100	46	6523	0		40	100	48	6983	0
	80	100	48	7768	0		80	100	47	8174	0		80	100	49	8952	0
	120	100	49	9156	0		120	100	48	9828	0		120	100	50	10927	0
	160	100	50	10552	0		160	100	49	11495	0		160	100	51	12915	0
TELL40 CRU	200	100	51	11931	0		200	100	51	13141	0		200	100	52	14879	0
	240	100	52	13076	0		240	100	52	14808	0		240	100	54	16868	0
	280	100	53	14721	0		280	100	53	16478	0		280	100	55	18863	0
	320	100	54	16130	0		320	100	55	18162	0		320	100	57	20881	0
	360	100	55	17516	0		360	100	56	19847	0		360	100	59	22895	0
	400	100	57	18900	0		400	100	58	21551	0		400	100	60	24923	271
	440	100	58	20337	0		440	100	59	23229	162		440	100	62	26955	2863
	480	100	59	21778	0		480	100	61	24997	1960		480	100	63	28997	2693403
	520	100	61	23218	683		520	100	62	26665	6432		520	100	64	30941	11665646
	560	100	62	24674	1365		560	100	63	28342	86571		560	100	66	33051	354496016
	600	100	63	26114	24595		600	100	65	30110	2927490		600	100	68	35212	7818417758
65%						78%						89%					

Domain of operation around 40 MHz

Same measurements with a finer granularity (5 MHz)

- Peak of errors centered on 25 MHz
- No resurgence of errors after 30 MHz

	Frequency	Time	Average temperature	Average current	Errors found
	5	100	41	5267	72
	10	100	42	5896	0
	15	100	43	6537	0
	20	100	44	7137	205117
	25	100	45	7660	18604767359276
	30	100	47	8399	305779605112
	35	100	48	9152	0
SOL40	40	100	48	9699	0
	45	100	49	11730	0
	50	100	48	11071	0
	55	100	49	11759	0
	60	100	52	12465	0
	65	100	54	13153	0
	70	100	55	13826	0
	75	100	55	14508	0
	80	100	54	15189	0
	85	100	53	15857	0

Modified cards

Same results with 2.1 (initial), 2.2.2 and 2.2.3 modified cards !

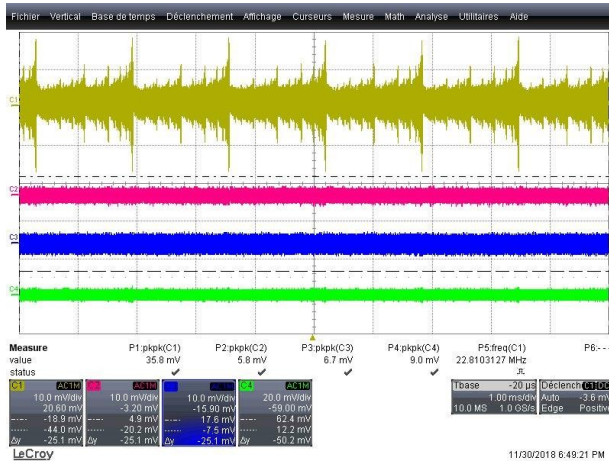
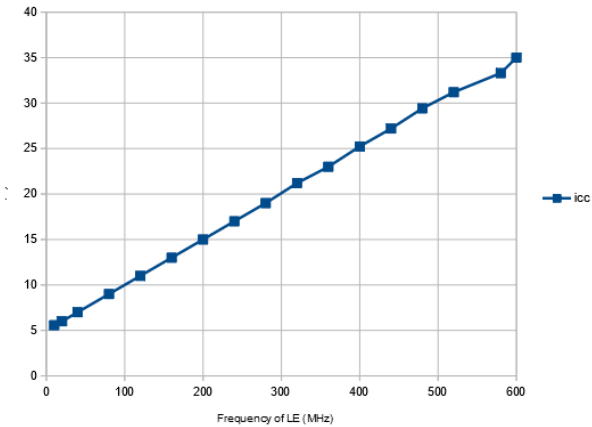


Decoupling ?

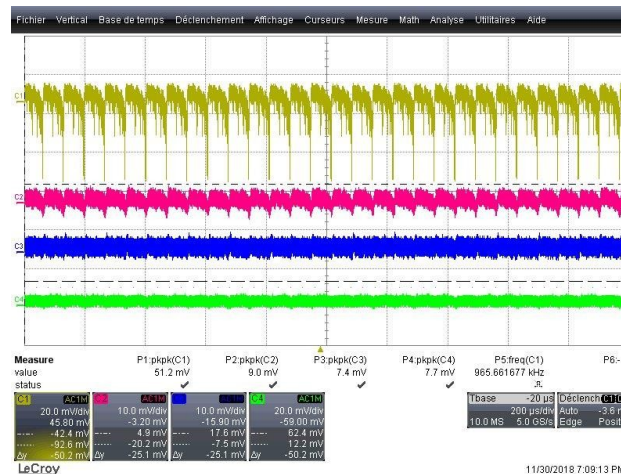
Significant noise on VCC

- **High frequency** cyclic noise
- Bursts of **low frequency** noise
- Proportional to injected frequency
- VCCR and VCCT preserved

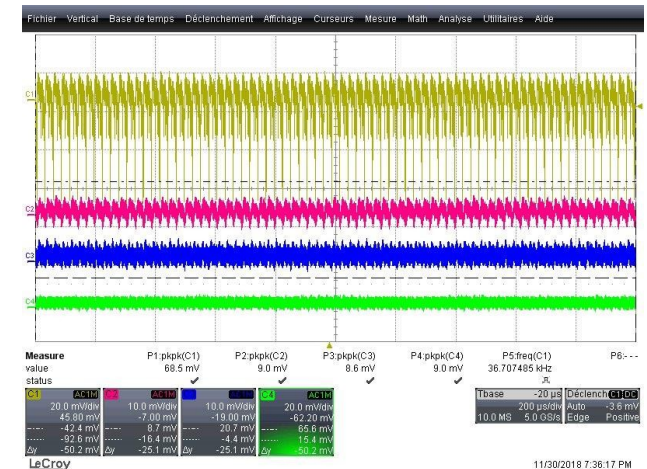
Icc Vs LE Frequency (89% FPGA occupancy)



Ripple noise at 10 MHz



Ripple noise at 240 MHz

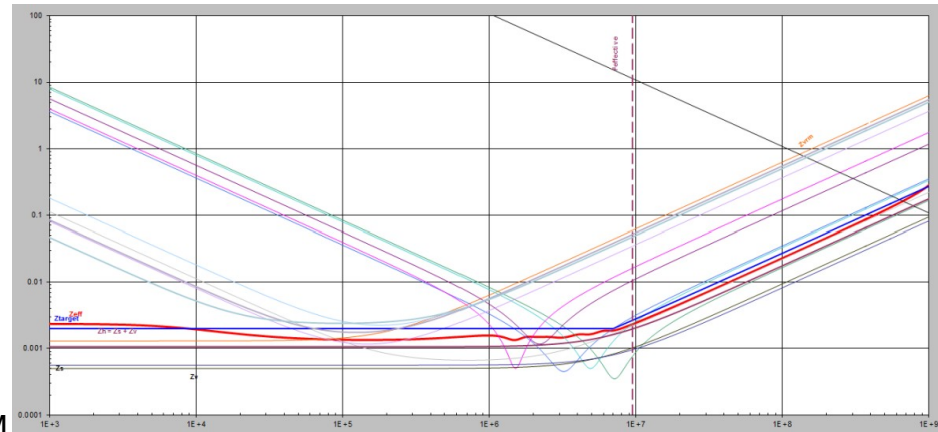


Ripple noise at 600 MHz

Decoupling simulations and measured ripple noise

Simulations

- Made with Intel PDN tool
- Impedance below Z_{target} except very low frequencies
- Simulations cross checked with similar results with:
 - Sigrity: takes into account board geometry and exact BOM
 - Ansys by a CERN expert from Alice (Michel Morel)

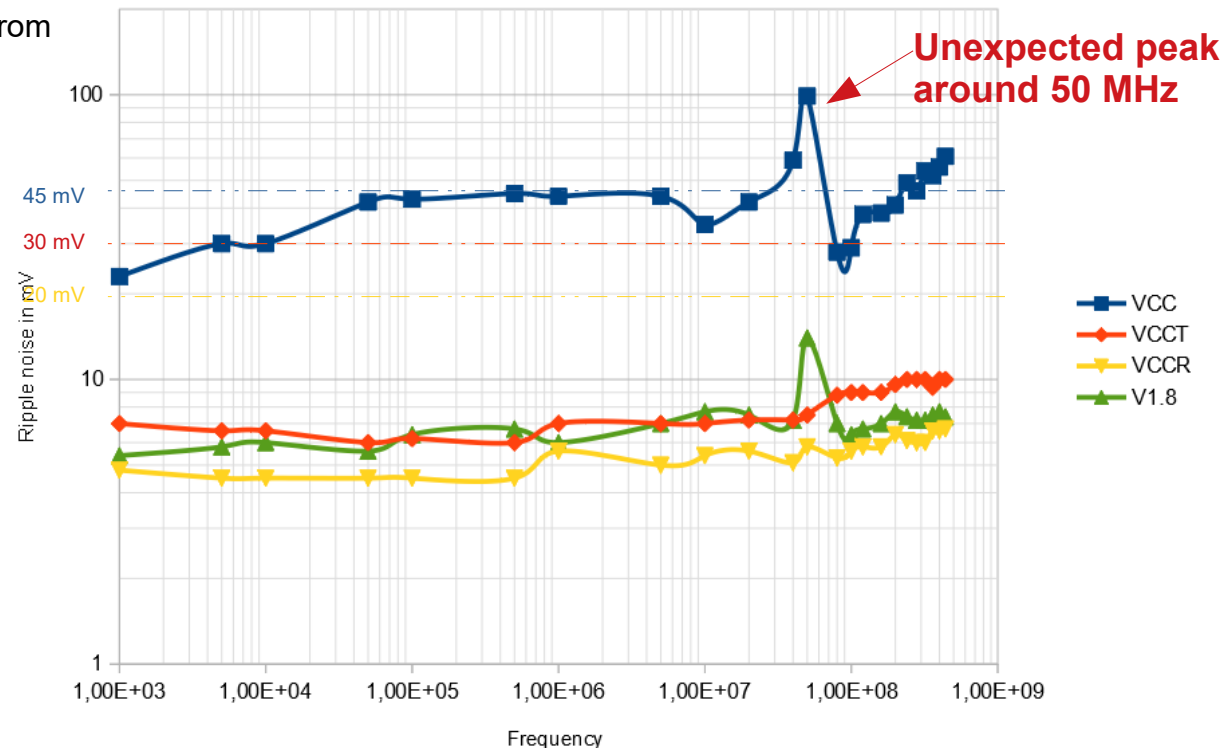


Measured ripple noise

- Authorized max values:
 - VCC : 45 mV
 - VCCT: 20 mV
 - VCCR: 30 mV
 - 1.8V : 54 mV

➡ VCC above threshold

➡ VCCR and VCCT OK

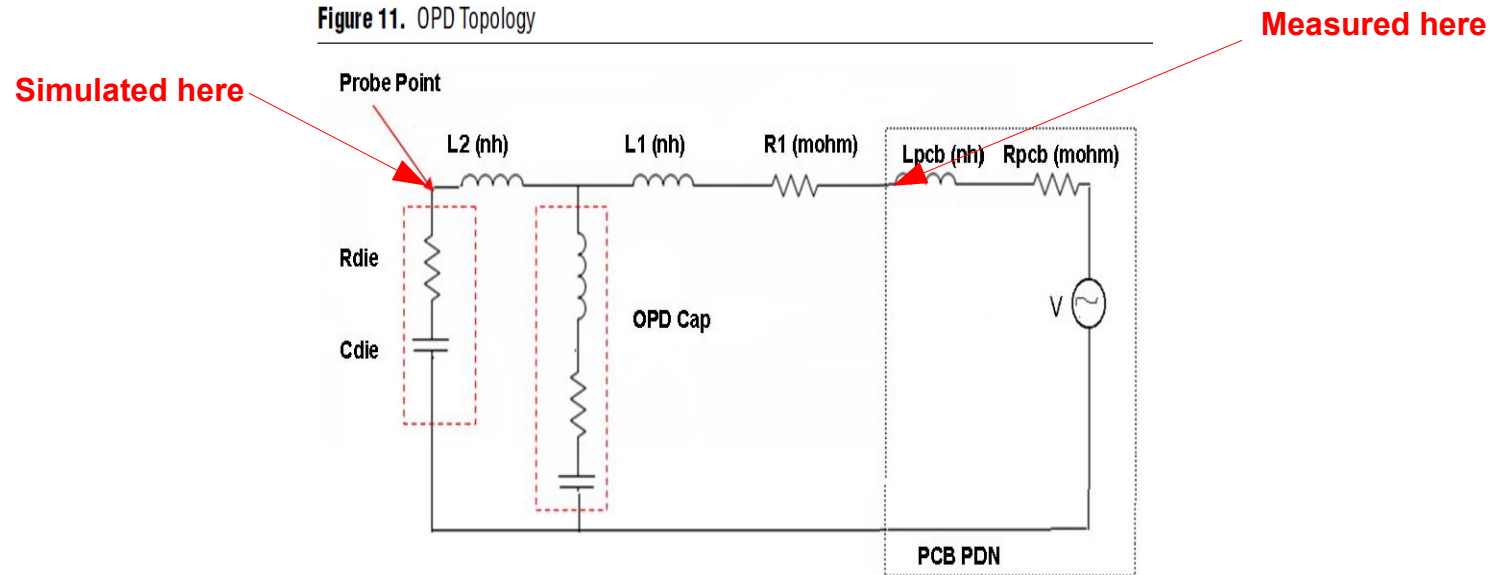


PDN simulation tools accuracy

PDN tool from Intel does not modelize accurately the leading inductance of capacitors

- Lots of approximations
- Use of external tools for determining values (lead inductances, etc ...)

Figure 11. OPD Topology



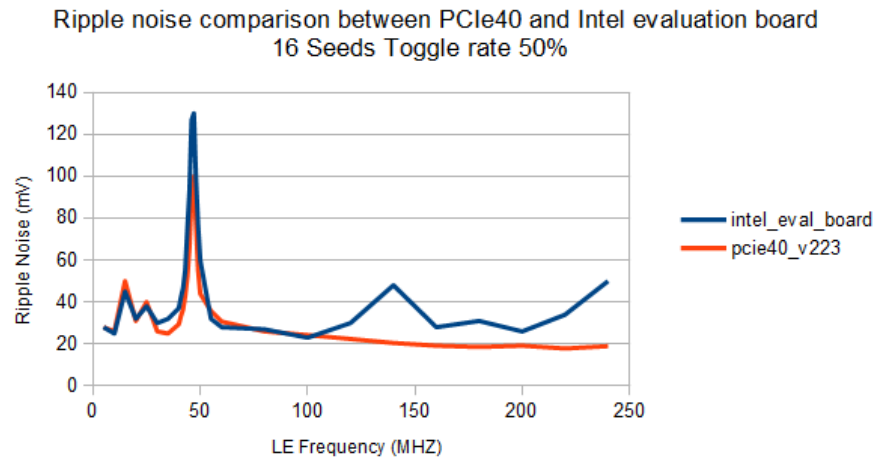
Sigrity or Ansoft do not modelize the FPGA on-package decoupling caps, nor the on-die caps

- Intel does not provide any model

Comparison with Intel SDK

Porting of a simple version of firmware on an Intel SDK (RPG only)

- Ripple noise measurements and comparison

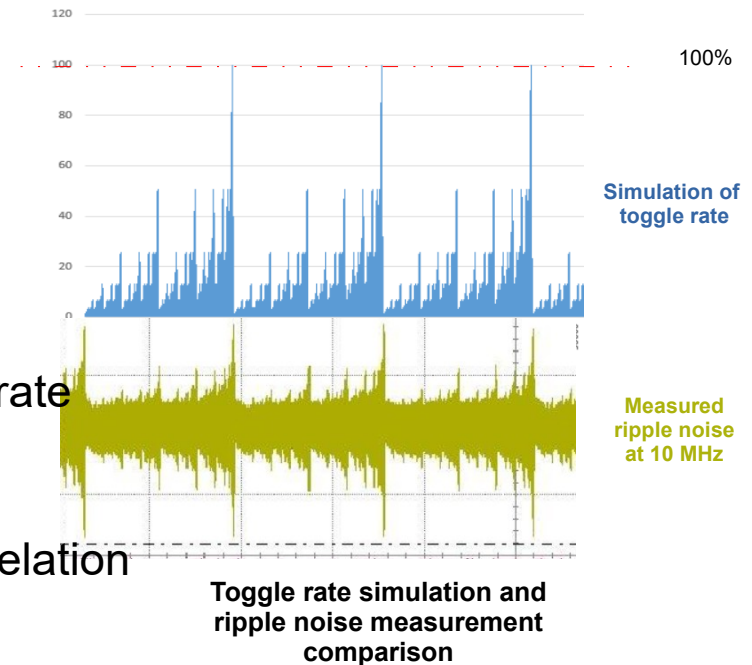


- Same peak around 50 MHz and even more ripple noise on high frequencies
 - ➡ Decoupling certainly not the reason

Noise bursts causes

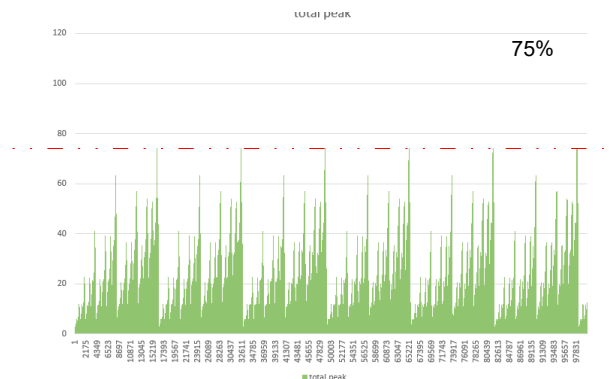
Correct firmware emulation ?

- Random pattern theoretically gives 50% of toggle rate in average ...
- ... but locally can be much higher than this.
- Simulations of toggle rate showed an obvious correlation with observed noise
- Average toggle rate 12.5% with peaks at 100% !



➡ new tests made with 16 different seeds

- Toggle rate amplitude decreases: 75%
- Average toggle rate 20%
- **Nearly no more errors at high speed,** but errors around 20 MHz remained

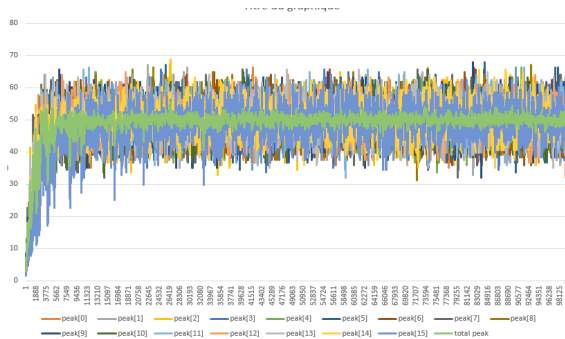


Frequency	Time	Average temperature	Average current	Errors found
10	100	44	5393	914342
20	100	42	5899	750598840
40	100	43	7139	0
80	100	46	9698	0
120	100	46	12251	0
160	100	48	14824	0
200	100	50	17384	0
240	100	52	19986	0
280	100	54	22584	0
320	100	57	25218	0
360	100	59	27861	0
400	100	61	30527	0
440	100	63	33210	0
480	100	65	35919	0
520	100	67	38652	0
560	100	70	41384	0
600	100	72	44139	992

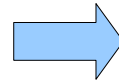
BER at 89%
16 different seeds

Fixing the issue

Improvement of RPG by adding more feed backs



Simulated toggle rate
(total contribution in green)



SOL40

TELL40
CRU

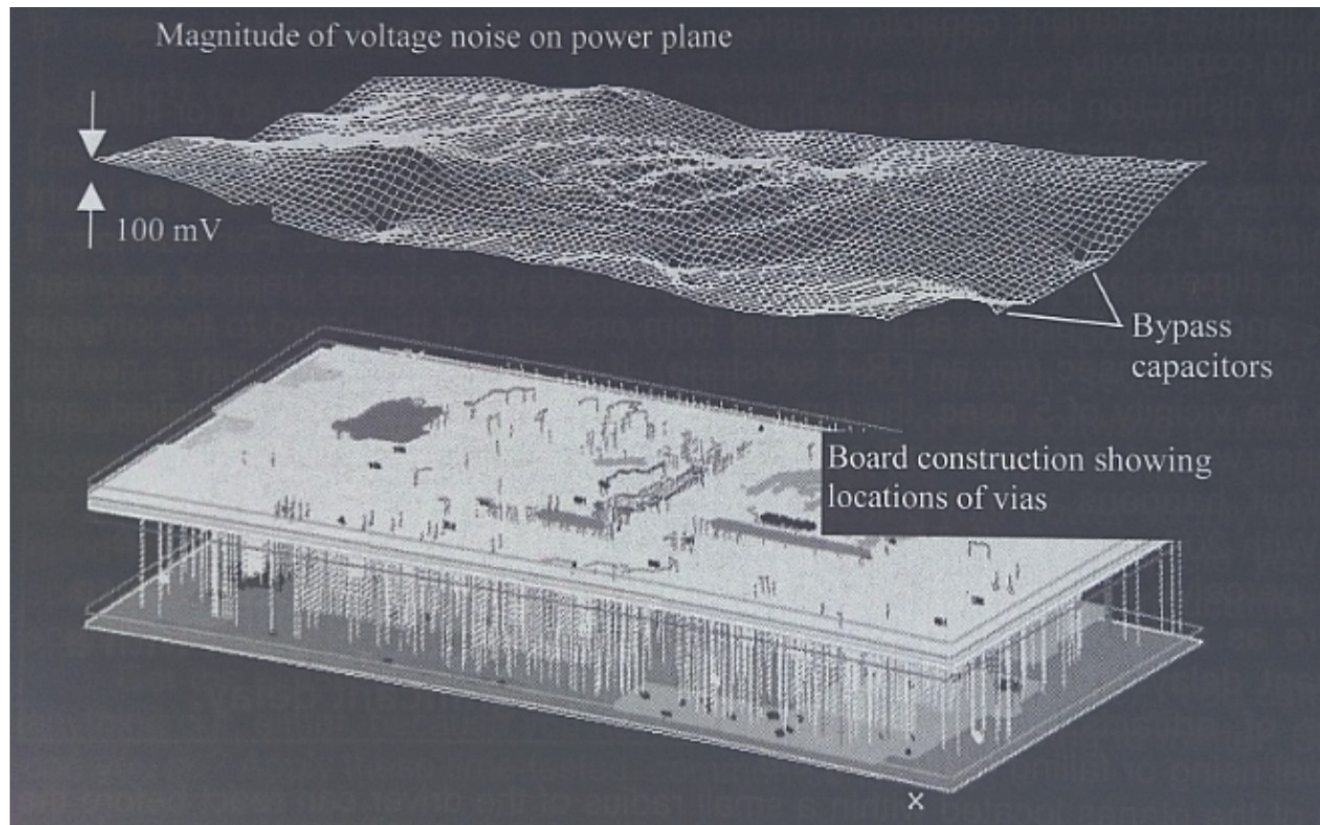
Frequency	Time	Average temperature	Average current	Errors found
5	100	41	5278	216
10	100	43	5929	0
15	100	44	6562	0
20	100	45	7180	2462
25	100	46	7701	18503123109715
30	100	47	8419	37593657128
35	100	48	9162	0
40	100	49	9753	0
80	100	54	15173	0
120	100	53	20619	0
160	100	58	26114	0
200	100	62	31746	0
240	100	67	37486	0
280	100	72	43351	0
320	100	78	49683	0

Limited measurement
because of current drawn

- Real toggle rate of 50%
- More current drawn
 - ➡ Maximum reached at 320 MHz

Power plane resonance ?

- Peaks at ~50 MHz, errors at ~25 MHz
- Resonance could happen if injected frequency is reflected on board edges

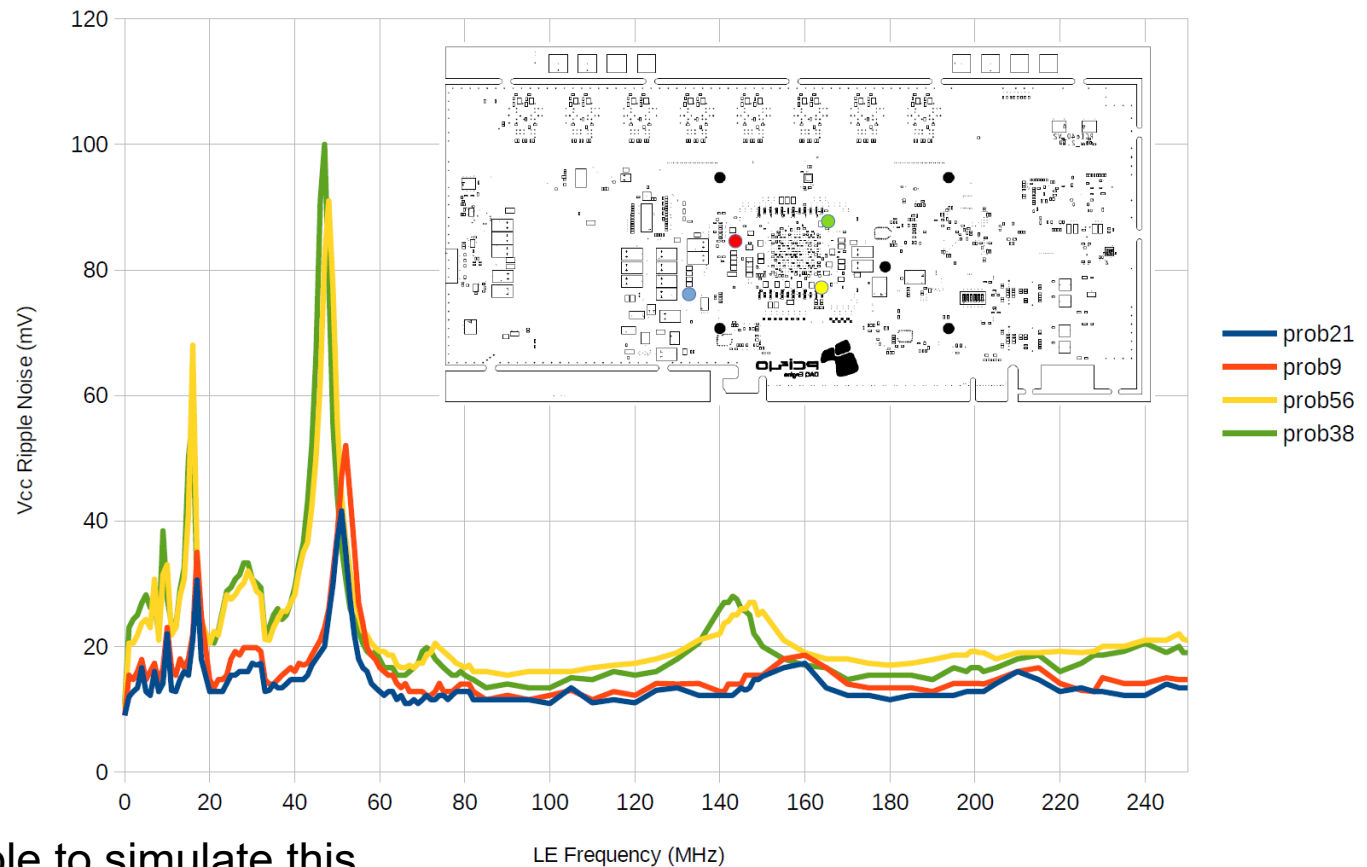


Power plane resonance ?

Geographical measurements to check this hypothesis

- 4 measurement points

Vcc Ripple Noise vs LE Frequency
for different probe locations on board



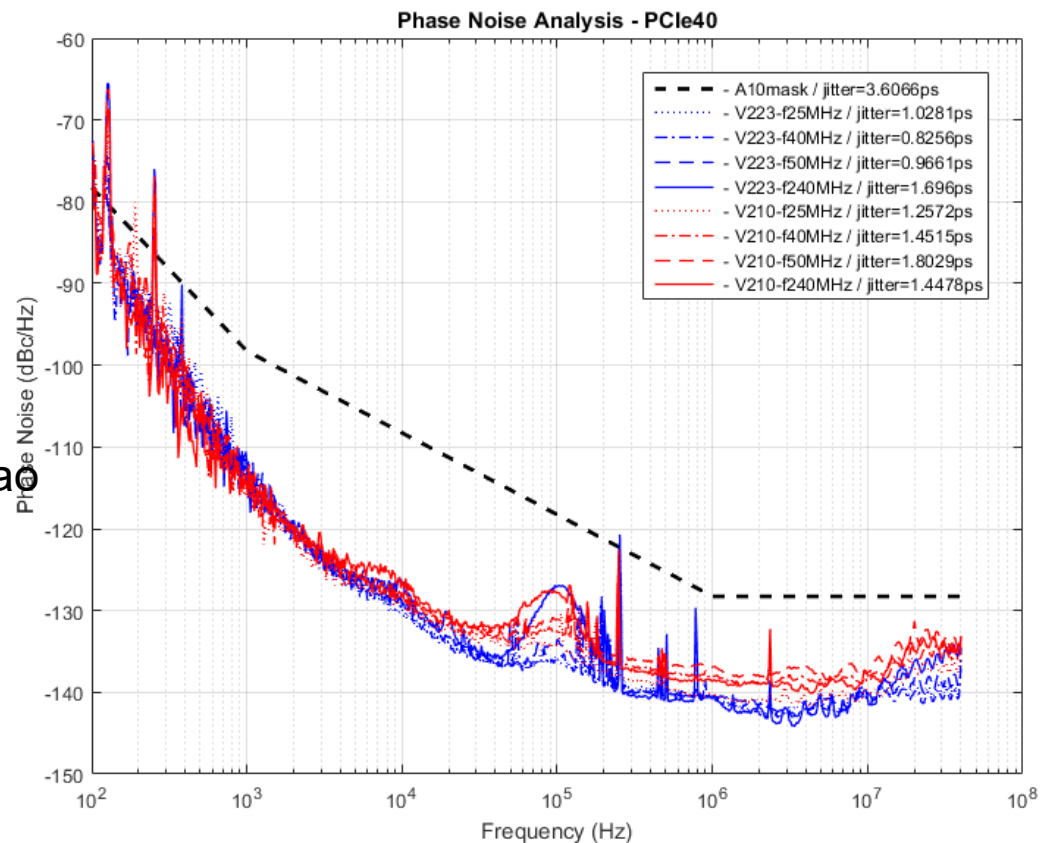
- No tool able to simulate this

PLL phase noise ?

Errors detected mainly on the emitting side: PLLs suspected

► *Phase noise on external PLLs could generate jitter on TX lines*

- To check this we measured the phase noise of PLLs feeding the refclks
- 4 frequencies injected:
 - 25 MHz
 - 50 MHz
 - 40 MHz
 - 240 MHz
- Two cards tested:
 - V210
 - V223
- Measured with Eduardo Brandao on an Agilent E5052B Signal Source Analyzer

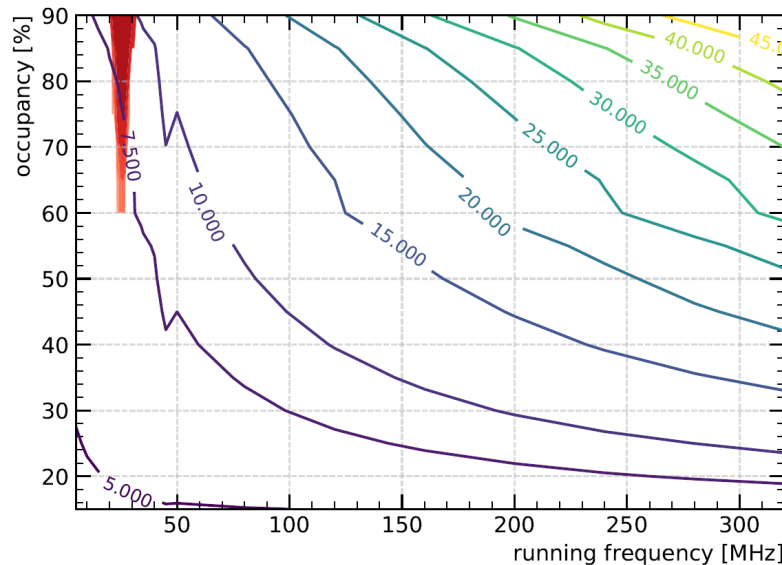


Results

- Phase noise within spec

Decision to go in production

No visible impact on frequencies used by LHCb and Alice

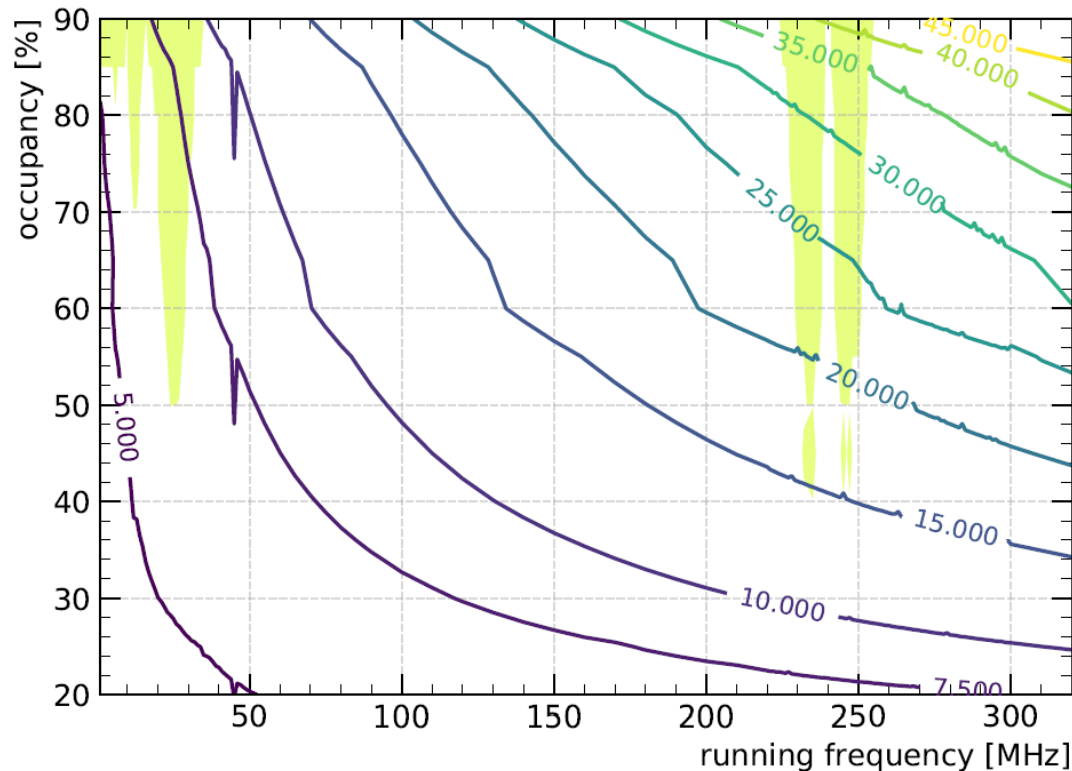


Possible alternatives in case we would work in the critical frequencies

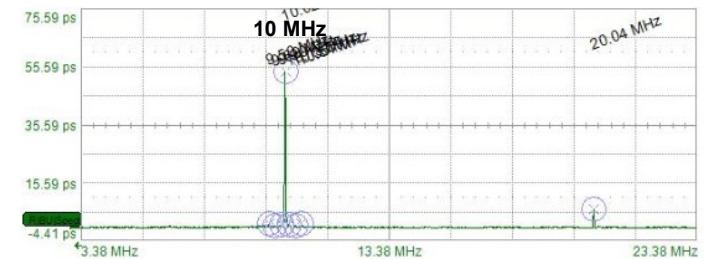
- Spreading the toggle rate over several shifted-phase clock domains
- Increase clock frequency

Refining the error space

Measurements with a finer granularity: new type of errors !



- Rightmost peaks clearly identified as beats between the 240 MHz refclk and the core logic injected clock



Jitter noise spectrum
when injecting a 230 MHz clock

FPGA internal PLLs testing

Change many PLL related parameters to see if errors disappear or increase

- Internal vs external feed back
- Bandwidth

Check other possible sources of beat

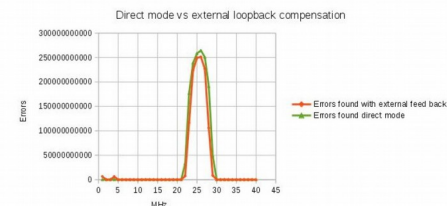
- PCIe DMA at 250 MHz

► *No significant change until ...*

Change external feed back by internal one in fPLL

Replaced loopback compensation mode by direct mode in fPLL

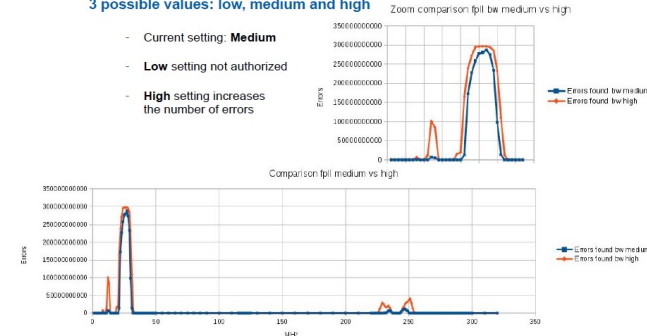
- External loopback circulates in the fabric
→ Could be subject to noise
- No visible effect.



FPLL loop bandwidth influence

3 possible values: low, medium and high

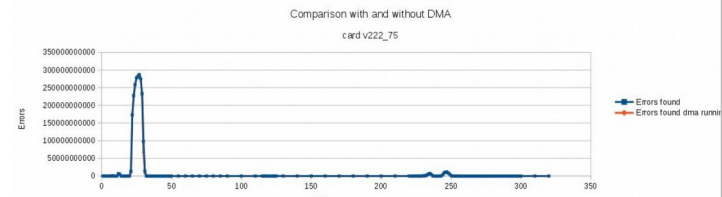
- Current setting: **Medium**
- **Low** setting not authorized
- **High** setting increases the number of errors



Run full DMA while running BER

Check that 250 MHz of PCIe does not interfere with 240 MHz of refclks

- Run BER when injecting all frequencies from 1 to 320 MHz
- Run DMA full speed
- Compare results with and without DMA
→ No difference

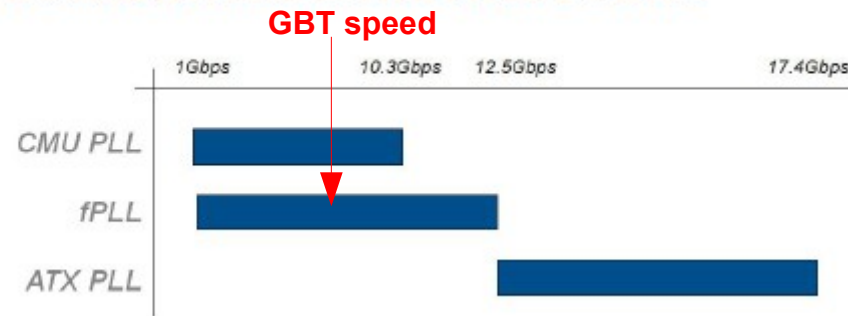


Replacement of fPLL by ATX PLLs

Although not recommended because not compatible with two rules given by Intel

- Use of PLL type:

Figure 171. Transmit PLL Recommendation Based on Data Rates



- Spacing

Arria 10 Transceiver ATX PLL
altera_xcvr_atx_pll_a10

Parameter Names	Parameter Values
K counter (valid in fractio...	1
L counter (valid in non-ca...	4
M counter	20
N counter	1
L cascade predivider/VC...	select_vco_output
L cascade counter (valid ...	1
PLL output frequency	2400.0 MHz
vco_freq	9600.0 MHz
datarate	4800.0 Mbps

3.1.1. Transmit PLLs Spacing Guideline when using ATX PLLs and fPLLs

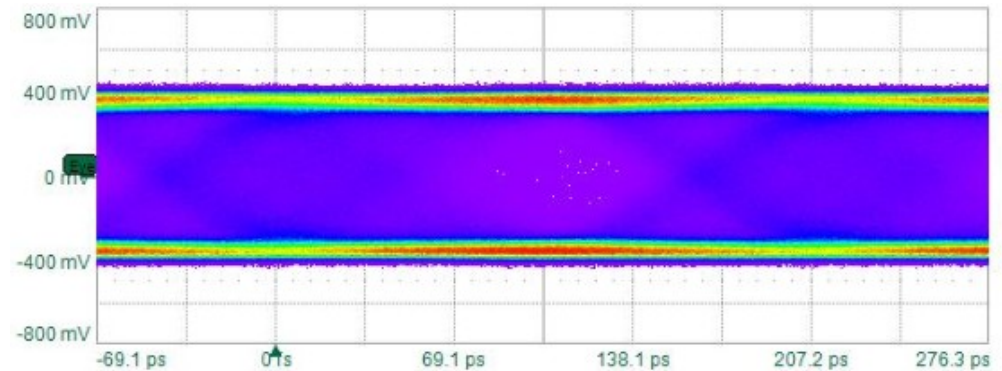
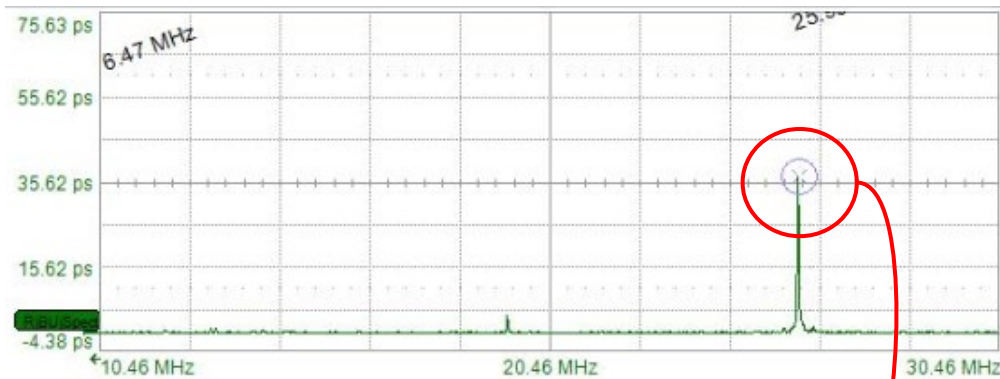
ATX PLL-to-ATX PLL Spacing Guidelines

For ATX PLL VCO frequencies between 7.2 GHz and 11.4 GHz, when two ATX PLLs operate at the same VCO frequency (within 100 MHz), they must be placed 7 ATX PLLs apart (skip 6).

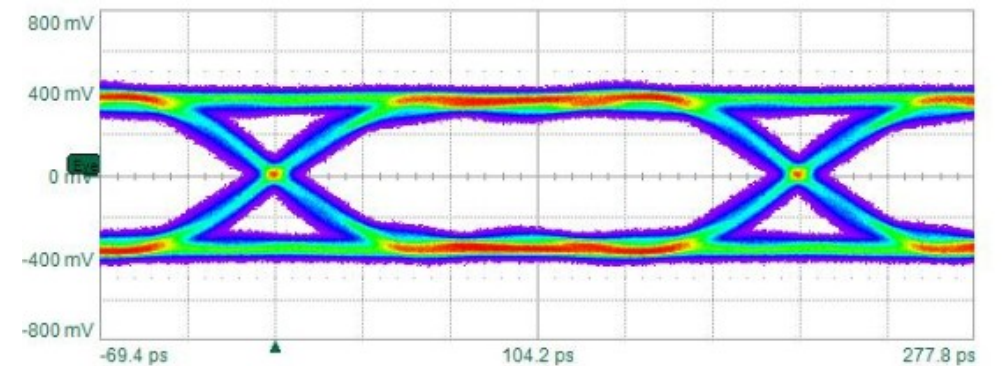
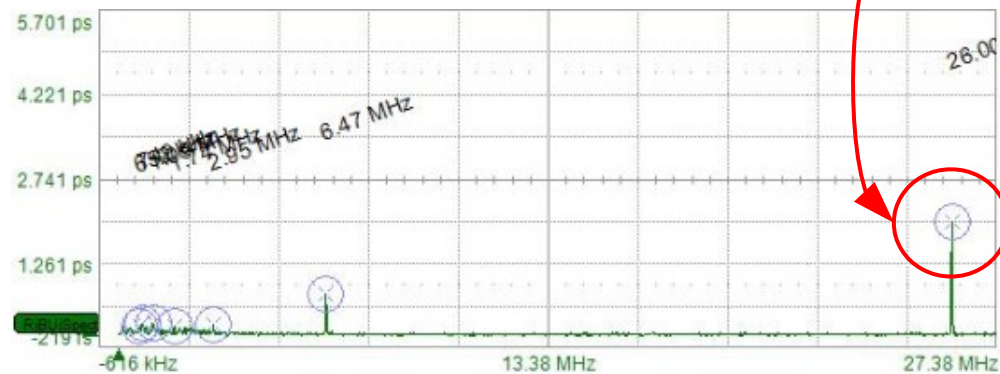
No more errors !

Worst case 1: 26 MHz injection

- **fPLL design**, jitter contribution of the 26 MHz core clock = **35 ps**



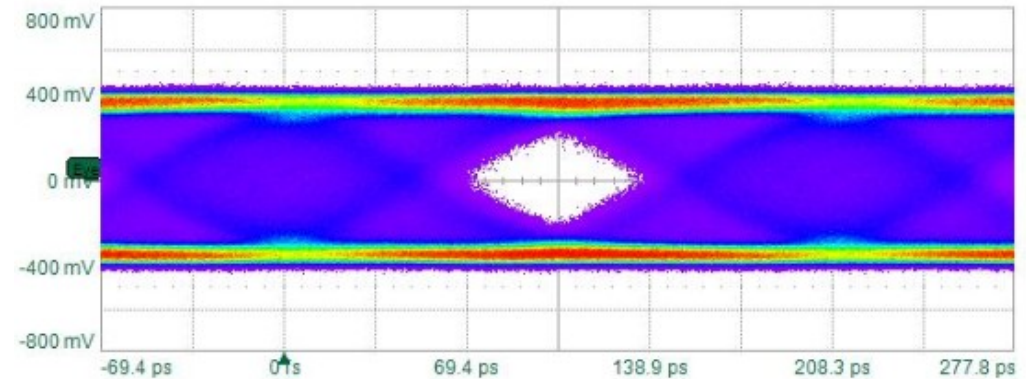
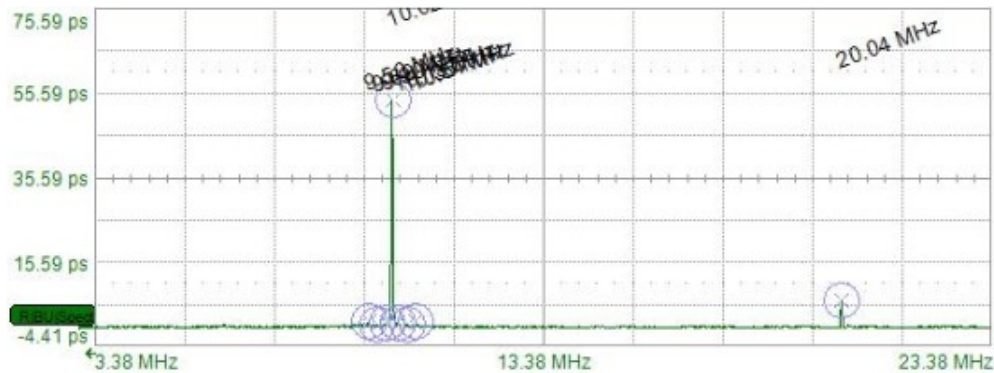
- **ATX PLL design**, jitter contribution of the 26 MHz core clock = **2.2 ps**



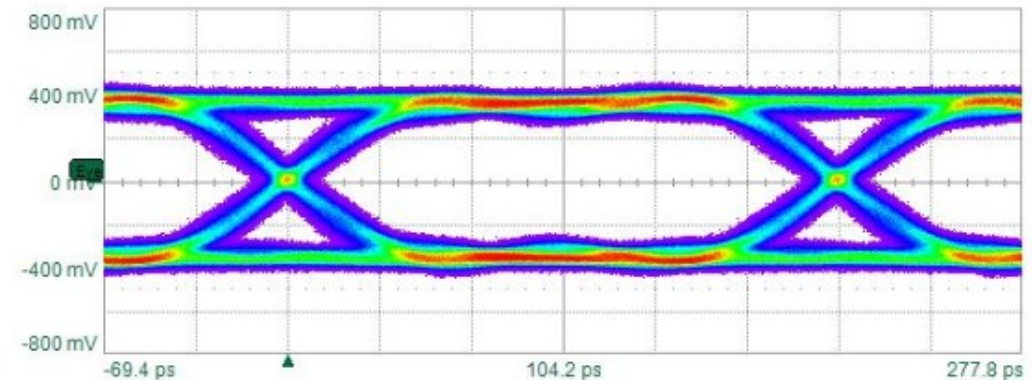
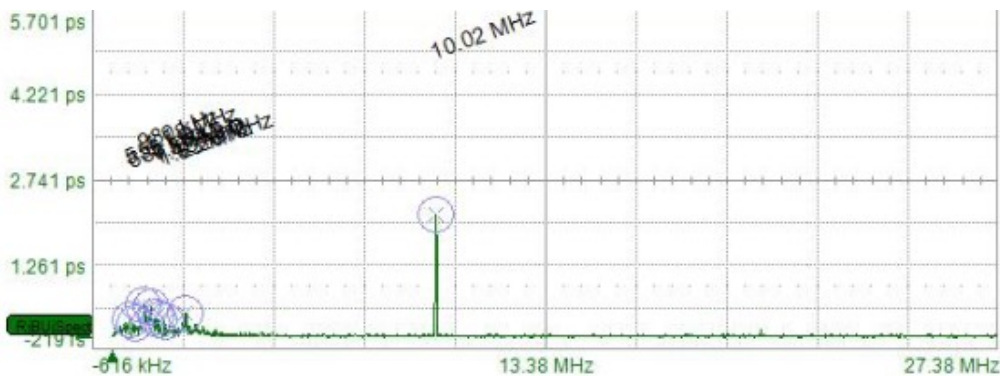
No more errors !

Worst case 2: 230 MHz injection (beats at 10 MHz)

- **fPLL design**, jitter contribution of the 26 MHz core clock = **56 ps**



- **ATX PLL design**, jitter contribution of the 26 MHz core clock = **2.2 ps**



Why does it works with ATX PLLs ?

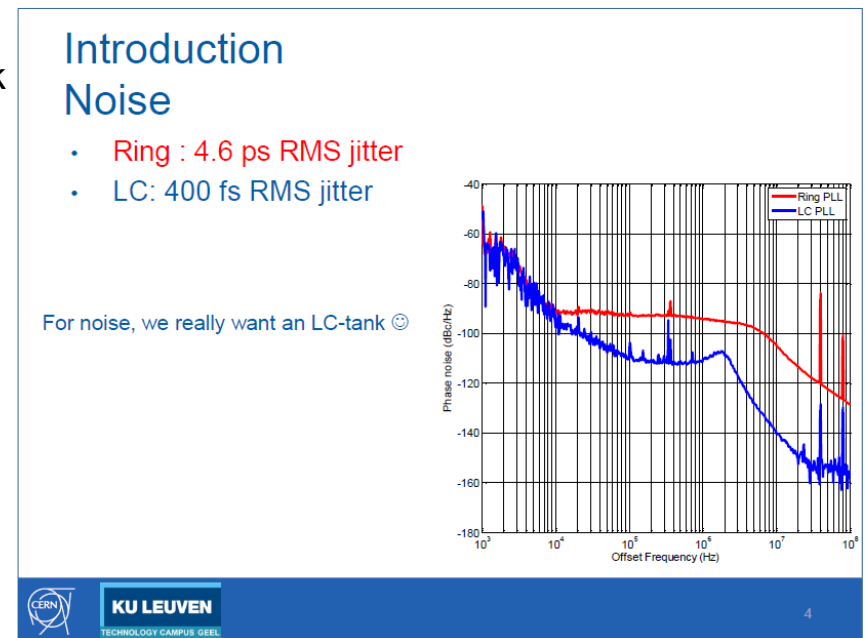
Not the same technology

- fPLLs are ring oscillator based VCO PLLs
- ATX PLLs are LC tank based

➔ The second type is much more robust to noise

see :

- 2.56-GHz SEU Radiation Hard LC-Tank VCO for High-Speed Communication Links in 65-nm CMOS Technology
(Jeffrey Prinzie , Student Member, IEEE, Jorgen Christiansen, Paulo Moreira, Michiel Steyaert, Fellow, IEEE, and Paul Leroux, Senior Member, IEE)
- Phase Noise and Jitter in CMOS Ring Oscillators
(Asad A. Abidi, Fellow, IEEE)

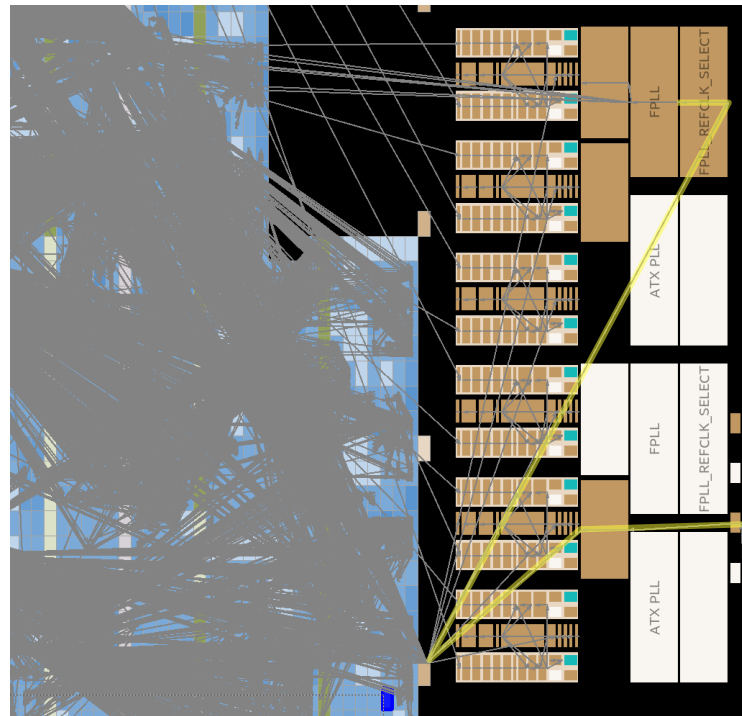


Ring oscillator vs LC tank PLLs noise comparison
Jeffrey Prinzie (CERN)

Further investigation with Intel

Finally find out to be a firmware issue !

- Refclk went into the matrix before going on fPLL !



- Due to residual constraint placed for allowing first versions of Quartus to converge (formerly crashed with a clock tree congestion message)
- Same topology when ATX PLL is instantiated but ATX PLL more robust

fPLL design results after modification

Frequency In MHz	Total jitter In ps	Random RMS	Deterministic p-p	Core current In A
10	37.23	2.167	6.762	5.953
20	37.23	2.169	6.728	7.209
25	37.47	2.170	6.963	7.710
40	37.81	2.170	6.728	9.735
80	37.72	2.170	7.200	15.072
120	37.62	2.170	7.109	20.367
160	37.59	2.169	7.076	25.814
200	37.55	2.170	7.031	31.206
240	37.35	2.168	6.865	36.499
280	37.45	2.169	6.938	41.956
320	37.37	2.168	6.877	47.462

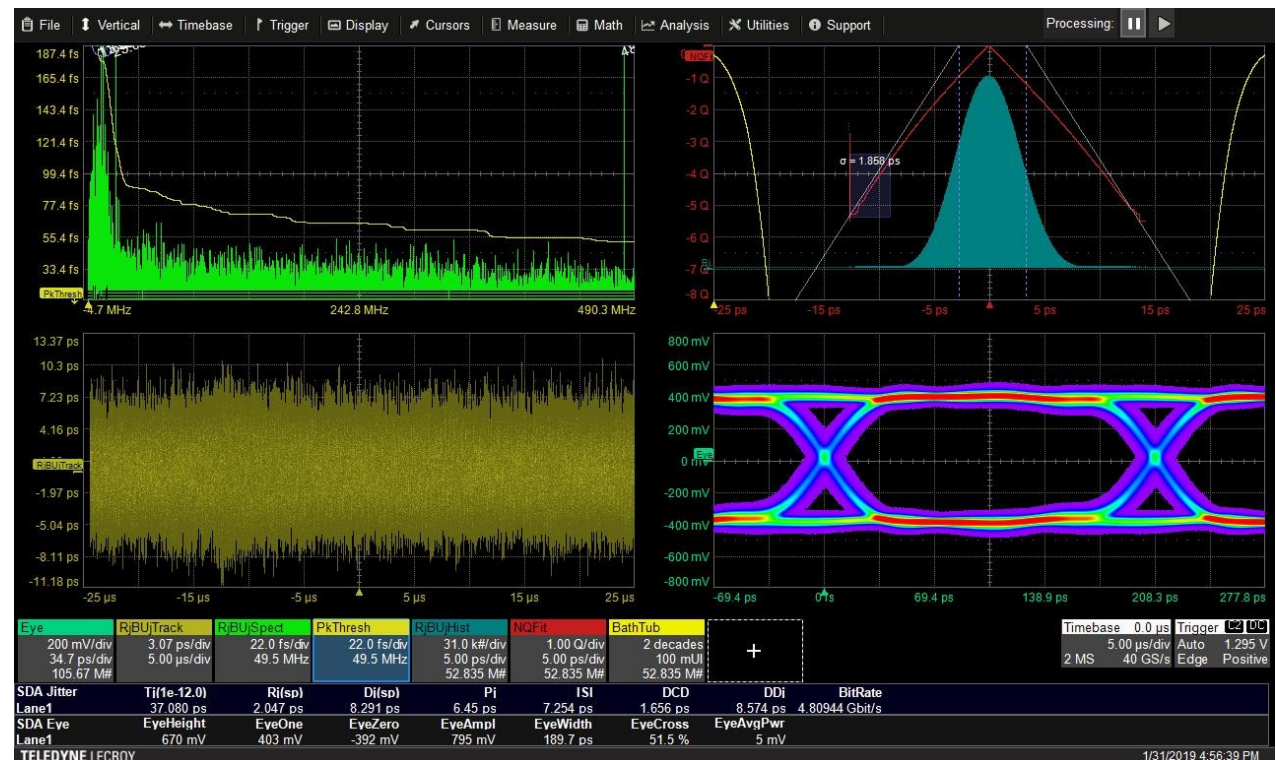
Worst case: 90% occupancy, 25 MHz core clock

- Very stable over the full spectrum:
 - o Total jitter 37 ps
 - o Random: 2 ps RMS
 - o Deterministic: 8.2 ps p-p

Test run during 2 days

- Over 48 links
- On 16 cards
- 300 MHz injected clock (45 A on the core)
- 10 m optical loopback

➡ **No errors**

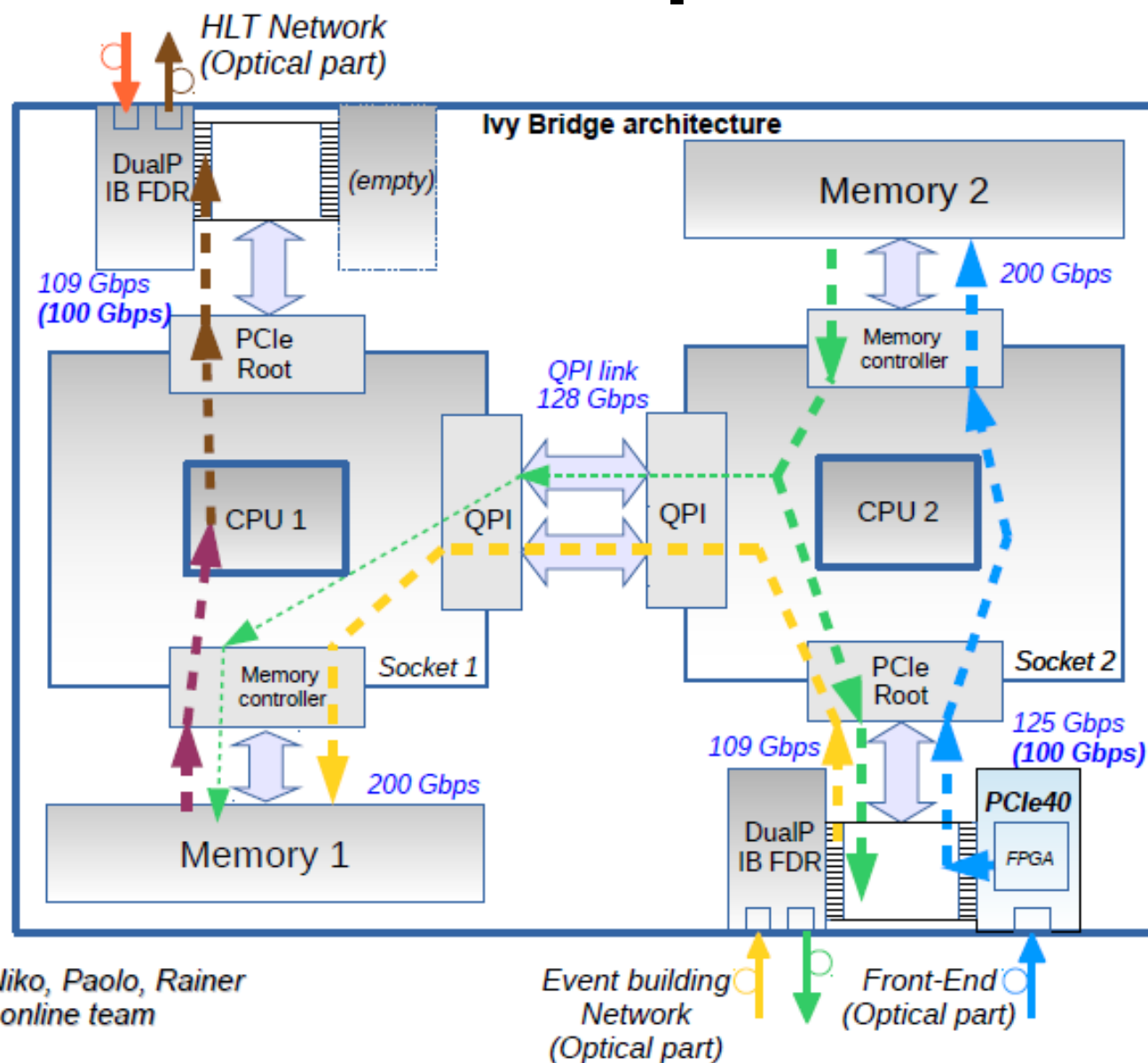


Conclusion

- Cards addressing many needs in our community
 - o Large acquisition capability, high processing power
 - o Powerful interface between dedicated Front-Ends and commercial computer CPUs
 - o Flexible enough to be used in many ways (3 functions in LHCb: DAQ, ECS and TFC, can fit ALICE needs as well, also selected for the readout of the $\mu 3E$ experiment)
 - o Lots of effort spent for optimizing the card for production (automatic testing, long time acceptance testing, automatic recording)
- Importance of testing a card to its limits
 - o Most of debugs are led with minimum firmwares
 - o High currents in high ends FPGAs raise new problems
- The PCIe40 card has been exhaustively tested
- Many lessons learned
 - o Better understanding of power plane geometry effects
 - o Better understanding of decoupling
 - o Limits of simulation tools
- Cards are now in production
 - o 700 cards for LHCb, 550 cards for ALICE
 - o Available end of this year

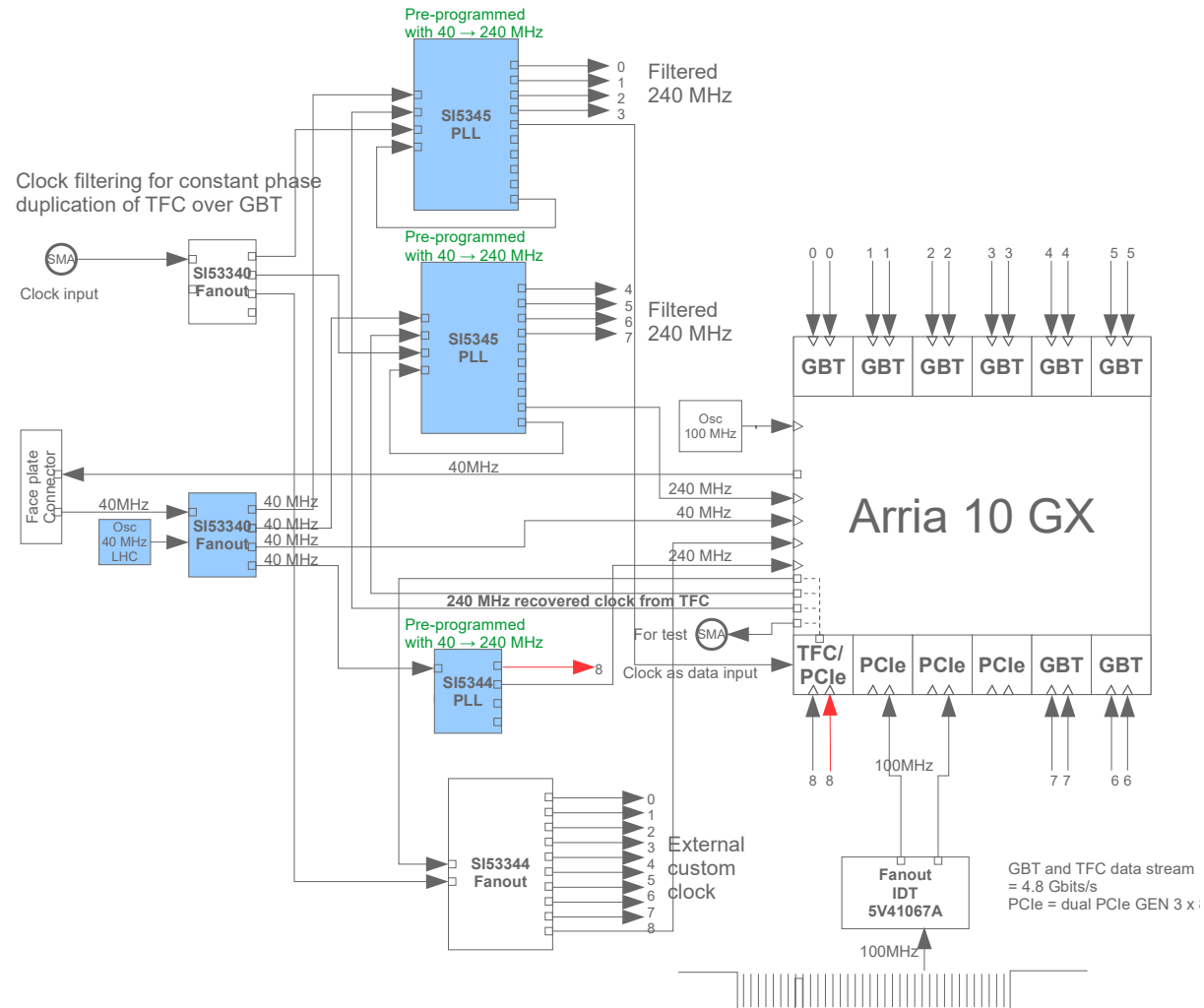
More information

Data path in the computer

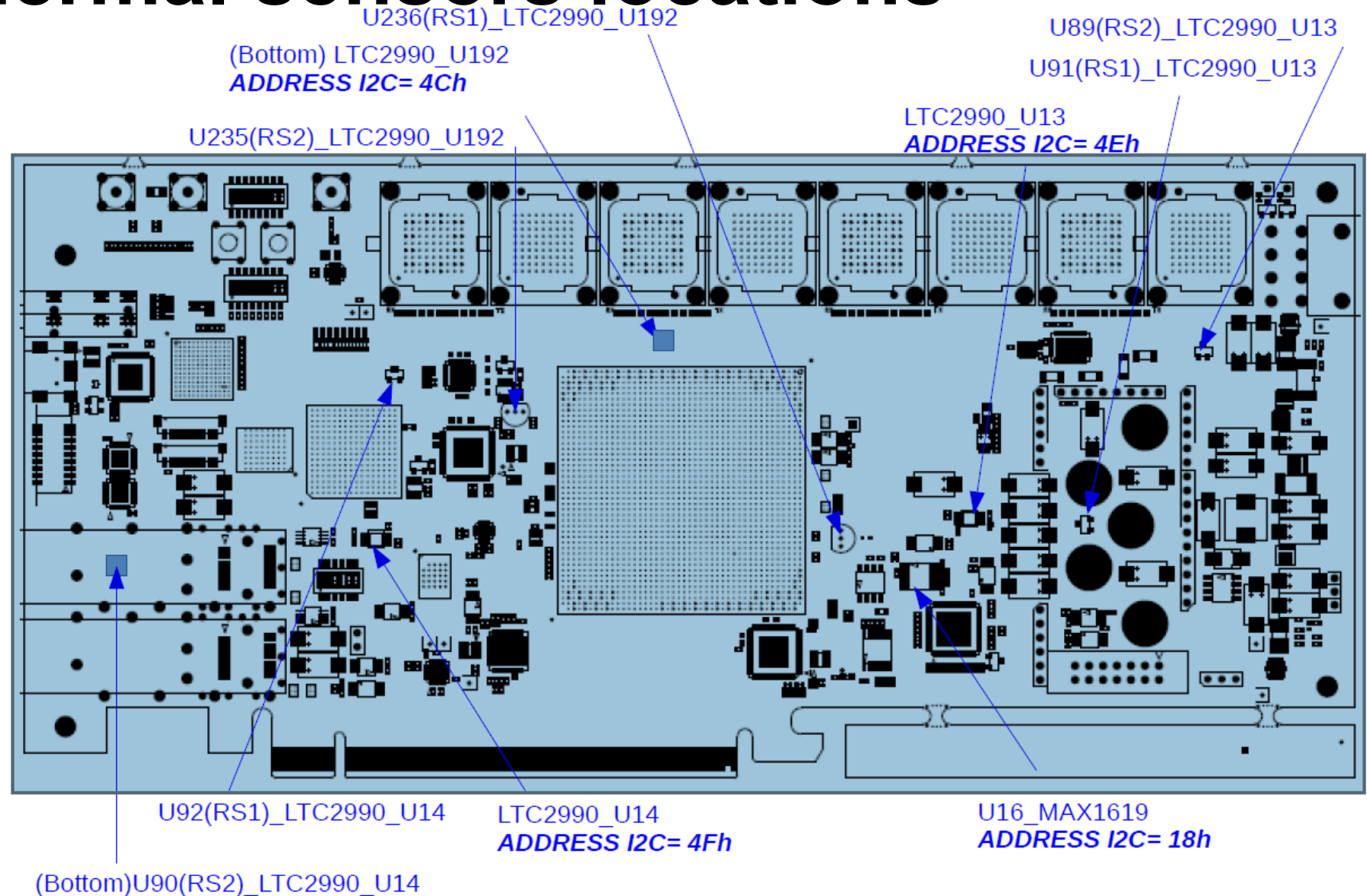


Clock distribution

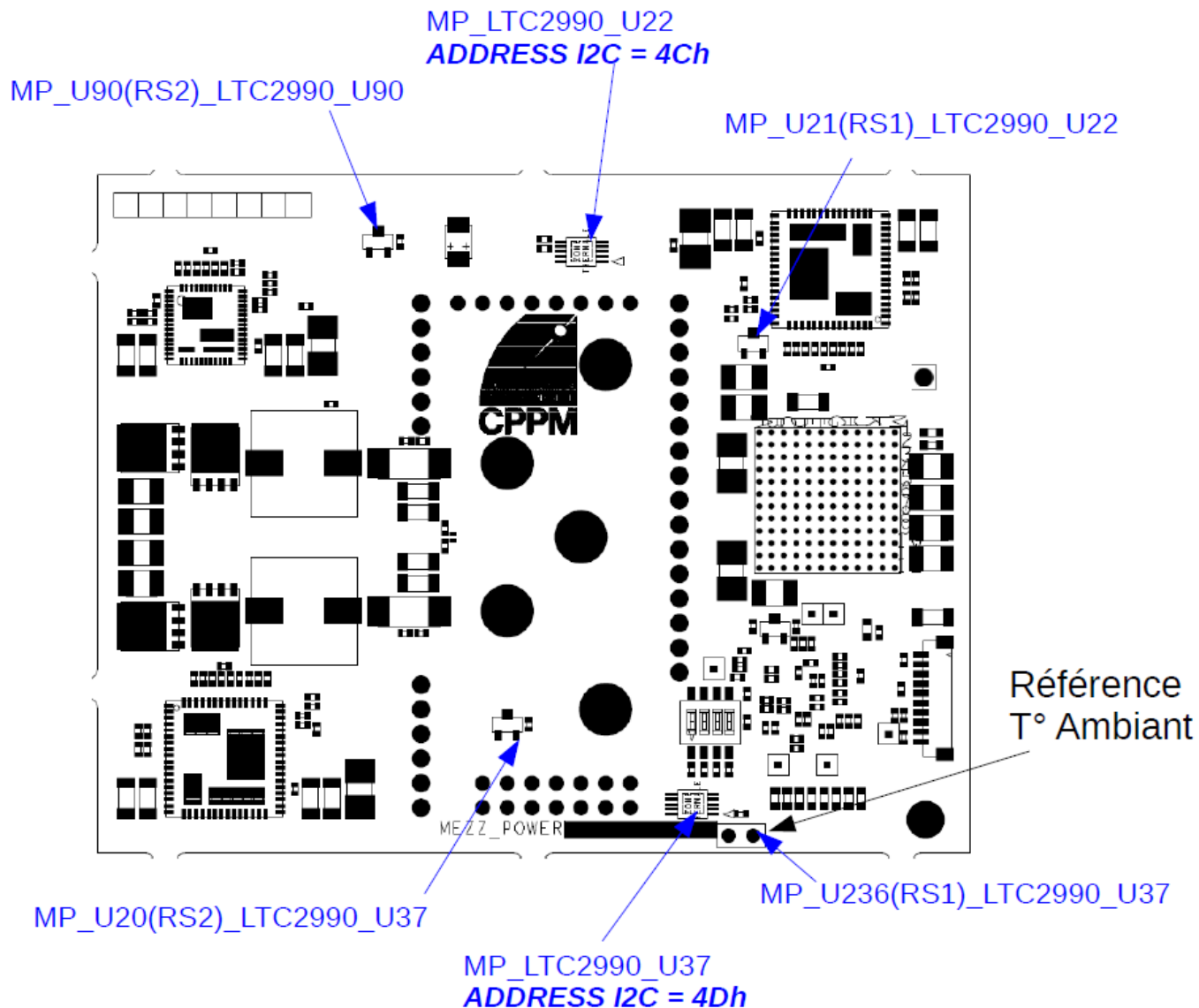
Clock Tree PCIe40V2



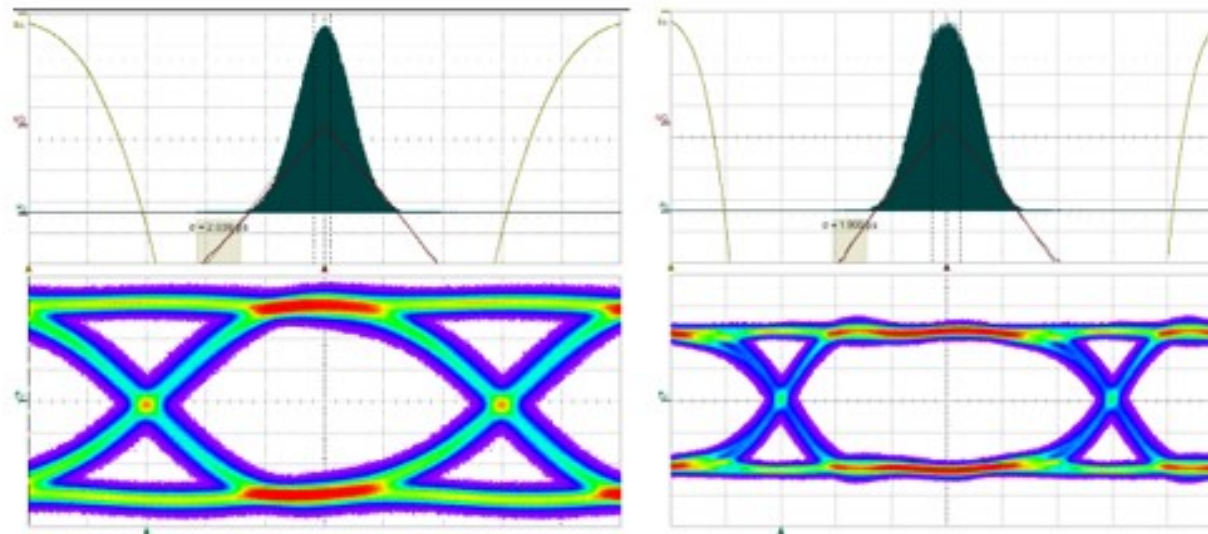
Thermal sensors locations



Thermal sensors locations



Eye diagrams



Measurements at 10.0Gbit/s
Total jitter ~ 36.82ps
Random Jitter ~ 2.22ps
Deterministic Jitter ~ 5.6ps

Measurements at 5.0Gbit/s
Total jitter ~ 37ps
Random Jitter ~ 2.1 ps
Deterministic Jitter ~ 3.44 ps



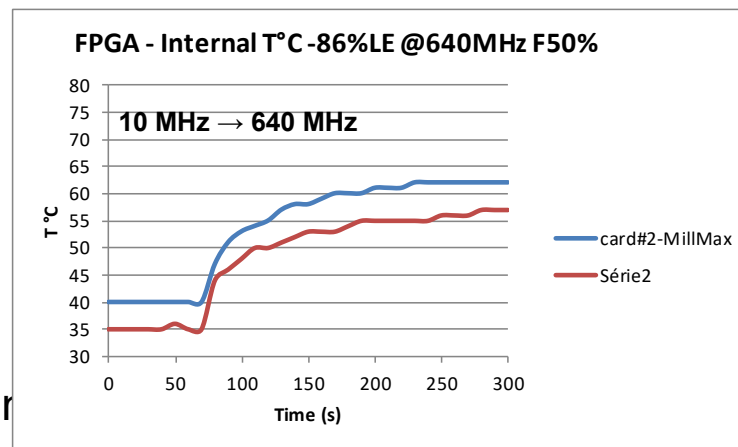
Mezzanine connector

Two choices : Samtec or Millmax

- Samtec : classical « full » connectors
- Millmax « transparent » connectors to let the air flow under the mezzanine

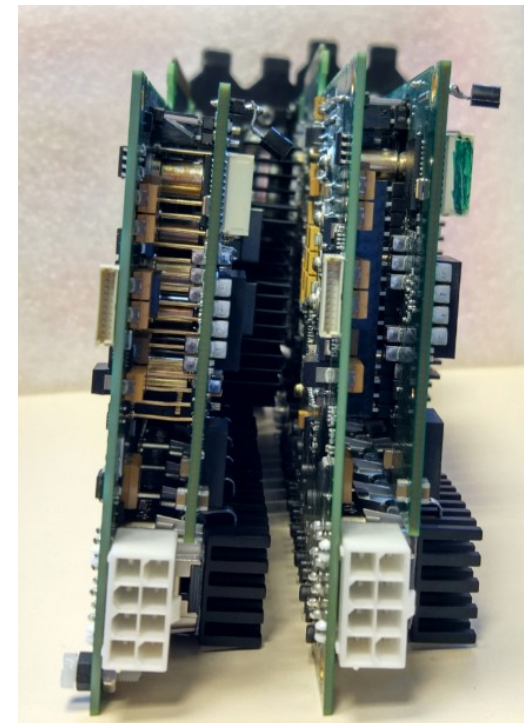
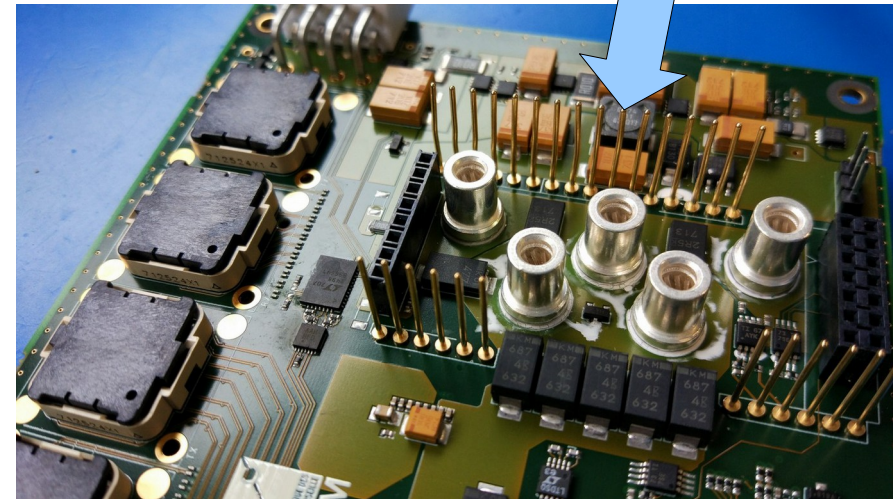
Cooling tests made with both solutions

- Counter intuitive results : Millmax card hotter than Samtec one (~5 to 6°C)
 - ➡ Venturi effect ?



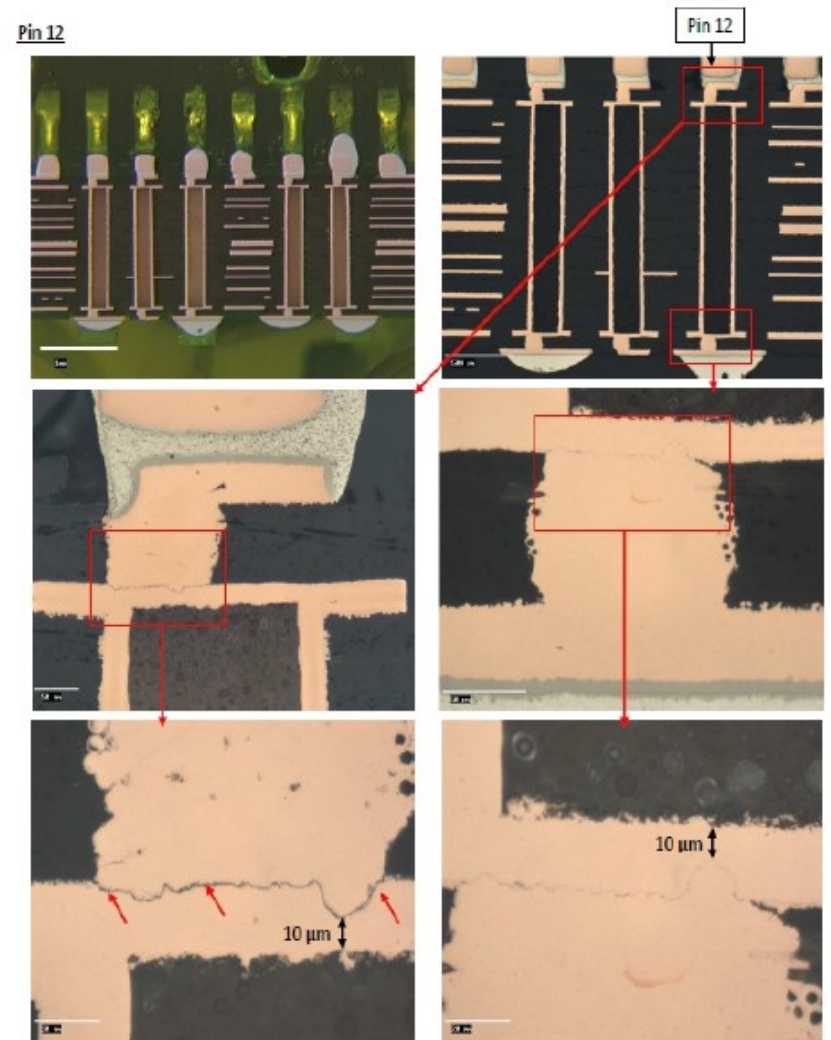
- Fir

➡ Much easier to mount



The PCB episode

- First batch of 6 MiniDAQ2 almost failed. Three boards survived but would die soon.
- After a long investigation, the issue was localized on the PCB. It was due to micro-cracks in the so-called stacked vias.
- A new board with a PCB from a different manufacturer was delivered Feb 15, 2017.
- After an extensive campaign of tests we concluded that the board is fully functional.

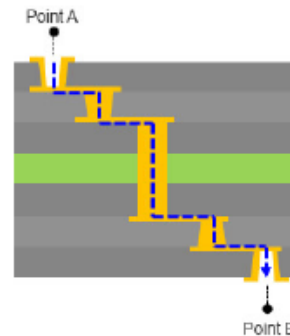


Routing

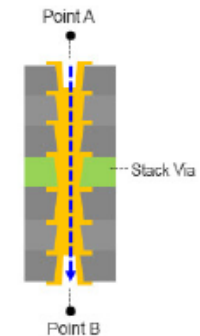
Use of staggered vias instead of stacked vias

- Slight degradation of signal integrity
- But more subcontractors able to manufacture the card

Staggered vias



Stacked vias



Stackup

- 14 layers
- 70μ thick planes for power
- HR408 high speed PCB
- More than 10000 vias among which 67% are microvias
- ~ 1750 components

PCB stackup & manufacturing requirements

SM above CU: 20μ +/-15μ registration 50μ
 CF: 12μ + plating = 55μ +/-10μ
 1 x 1086 RC 62% = 70μ +/-10μ
 CF: 12μ + plating = 40μ +/-10μ
 2 x 106 RC 70% = 90μ +/- 13μ
 CU: 17μ +/-5μ
 CCL: 1 x1086 RC 58% = 75μ +/- 13μ
 CU: 35μ +/-10μ
 2 x 106 RC 70% = 85μ +/- 13μ
 CU: 17μ +/-5μ
 CCL: 1 x1086 RC 58% = 75μ +/- 13μ
 CU: 35μ +/-10μ
 2 x 106 RC 70% = 80μ +/-13μ
 CU: 70μ +/-16μ
 CCL: 1 x1086 RC 58% = 75μ +/- 13μ
 CU: 70μ +/-16μ
 2 x 106 RC 70% = 80μ +/-13μ
 CU: 35μ +/-10μ
 CCL: 1 x1086 RC 58% = 75μ +/- 13μ
 CU: 17μ +/-5μ
 2 x 106 RC 70% = 85μ +/- 13μ
 CU: 35μ +/-10μ
 CCL: 1 x1086 RC 58% = 75μ +/- 13μ
 CU: 17μ +/-5μ
 2 x 106 RC 70% = 90μ +/- 13μ
 CF: 12μ + plating = 40μ +/-10μ
 1 x 1086 RC 62% = 70μ +/-10μ
 CF: 12μ + plating = 55μ +/-10μ
 SM above CU: 20μ +/-15μ registration 50μ



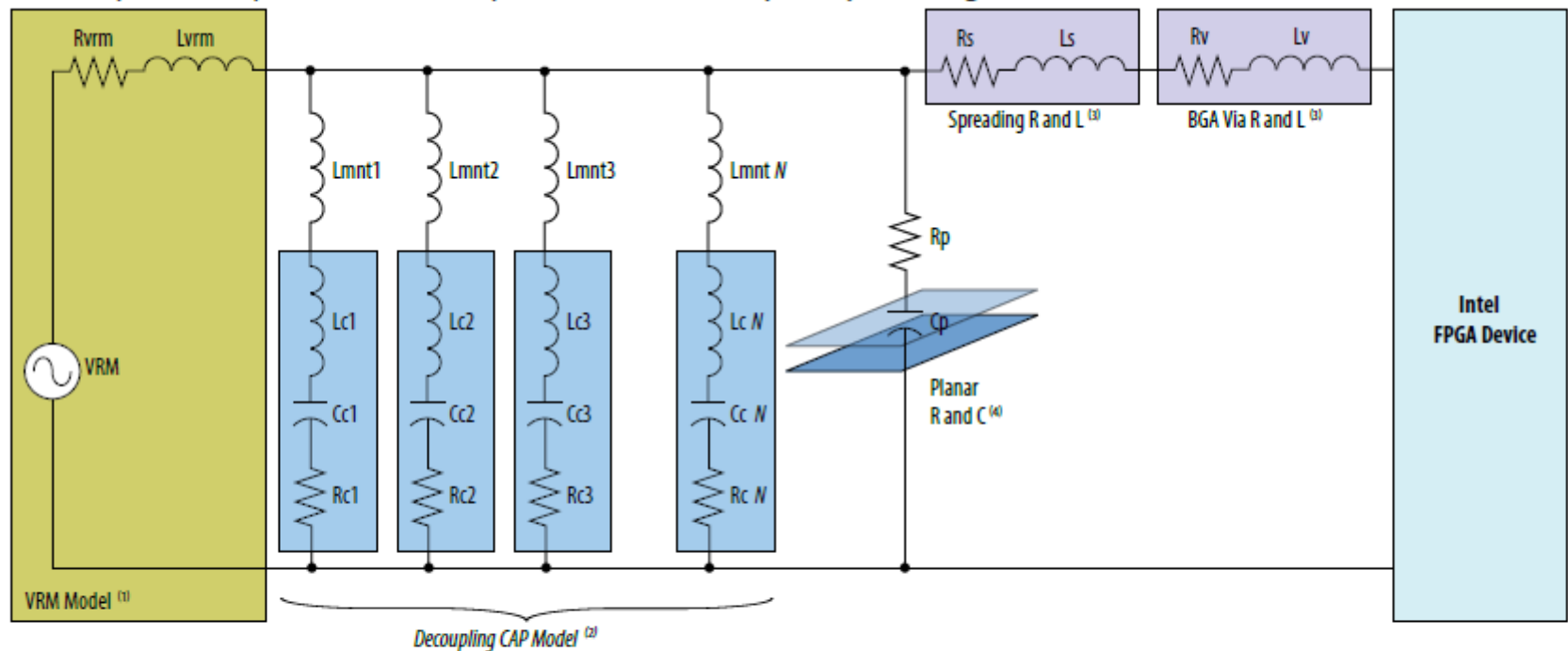
CAD Design rules & Copper balancing requirements

Bridge mini 100μ, Registration : QFN 50μ, BGA 50μ
 Top + W/Gmini 110μ + GND = Layout 60% = L1
 Full GND Smini 100μ = Layout 90% = L2
 W/Gmini 90μ + GND = Layout 50% = L3
 Full GND Smini 100μ = Layout 90% = L4
 Wmini 80μ, Smini 90μ + GND = Layout 50% = L5
 VCC1/2 Smini 120μ = layout 80% = L6
 VCC 0.9V + VCC 12V, Smini 240μ = layout 80% = L7
 Full GND S mini 240μ = Layout 90% = L8
 VCC5/6 Smini 120μ = layout 80% = L9
 Wmini 80μ, Smini 90μ + GND = Layout 60% = L10
 Full GND Smini 100μ = Layout 90% = L11
 W/Gmini 90μ + GND = Layout 50% = L12
 Full GND Smini 100μ = Layout 90% = L13
 Bottom + W/Gmini 110μ + GND = Layout 70% = L14
 Bridge mini 100μ, Registration : QFN 50μ, BGA 50μ

Decoupling

Principle

The PDN impedance profile is the impedance-over-frequency looking outward from the device.



Notes:

1. You can define or change VRM parameters in the Library sheet of the PDN tool.
2. You can define or change Decoupling Caps parameters in the Cap Mount, XYZ Mount, and Library sheets of the PDN tool.
3. R* and L* are parasitic resistances and inductances from BGA balls and PCB traces and connections.
4. Represents PCB layers dedicated to power and ground planes.

F_{EFFECTIVE} in PCB Decoupling

The PCB PDN cutoff frequency ($F_{EFFECTIVE}$) calculated by the PDN tool depends on the design trade-offs made on the PCB. The role of $F_{EFFECTIVE}$ is analyzed for both OPD and non-OPD packages.

Non-OPD Scenario

Figure 9 shows a simple topology for a rail without on-package decoupling.

Figure 9. Non-OPD Topology

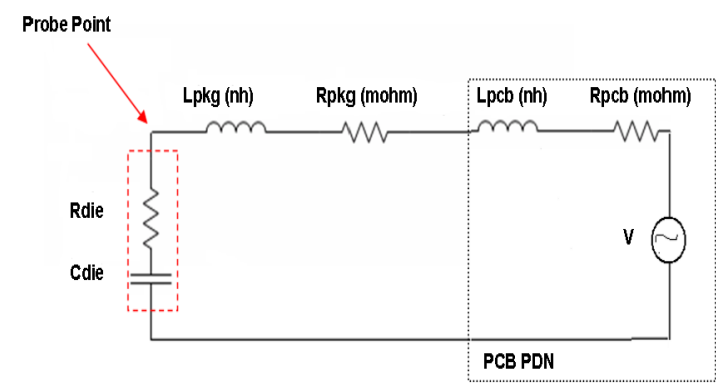
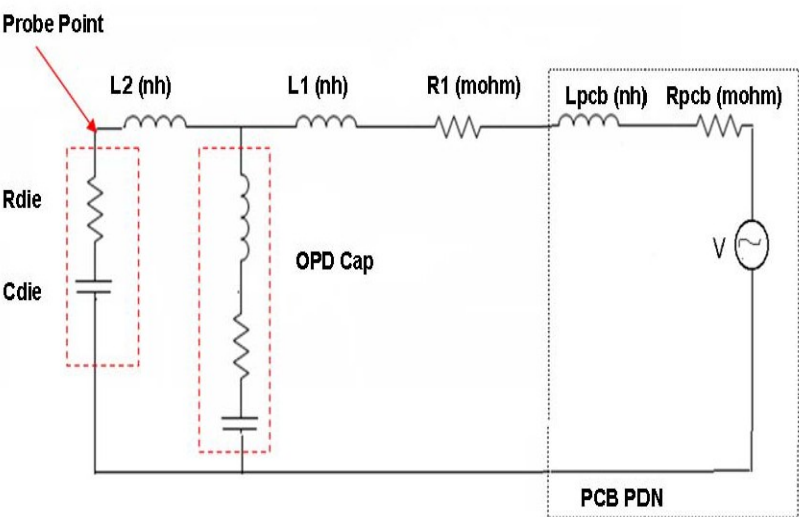


Figure 11. OPD Topology



$$F = 1 / 2\pi \sqrt{(Lpkg + Lpcb) * Cdie}$$

Figure 10. Non-OPD Topology Frequency Response

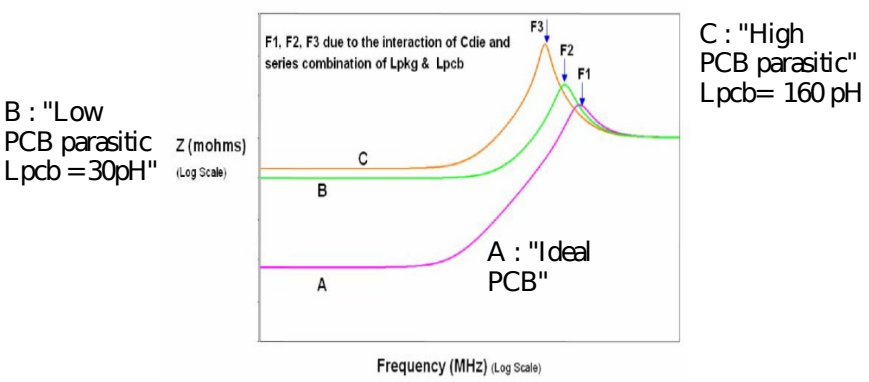


Figure 12. OPD Frequency Response

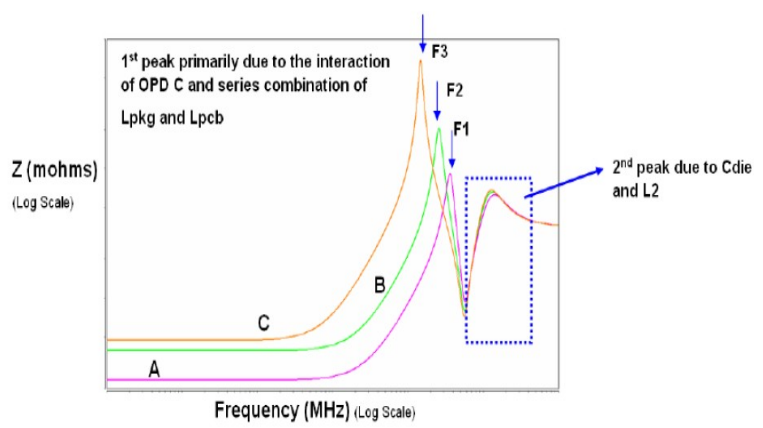
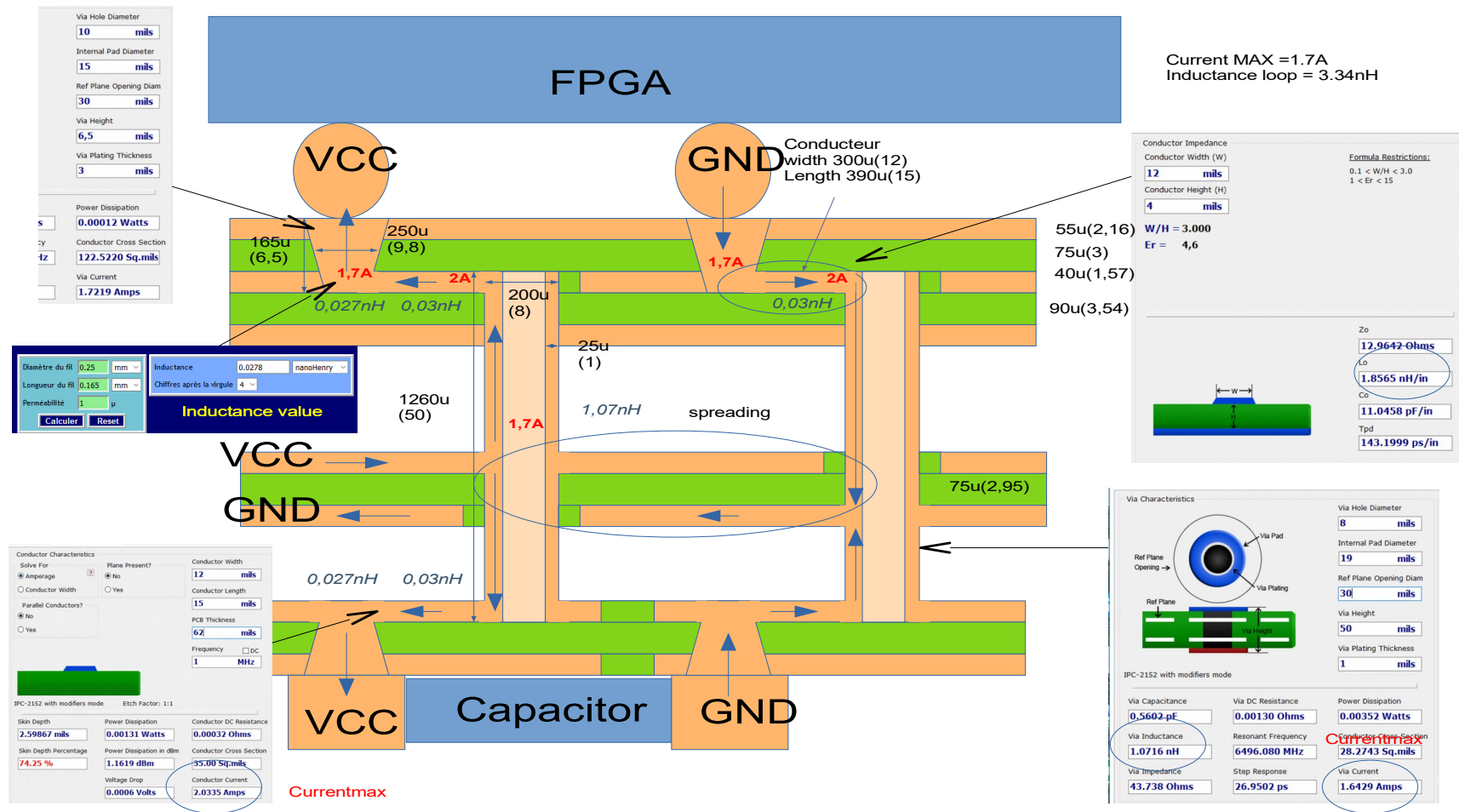


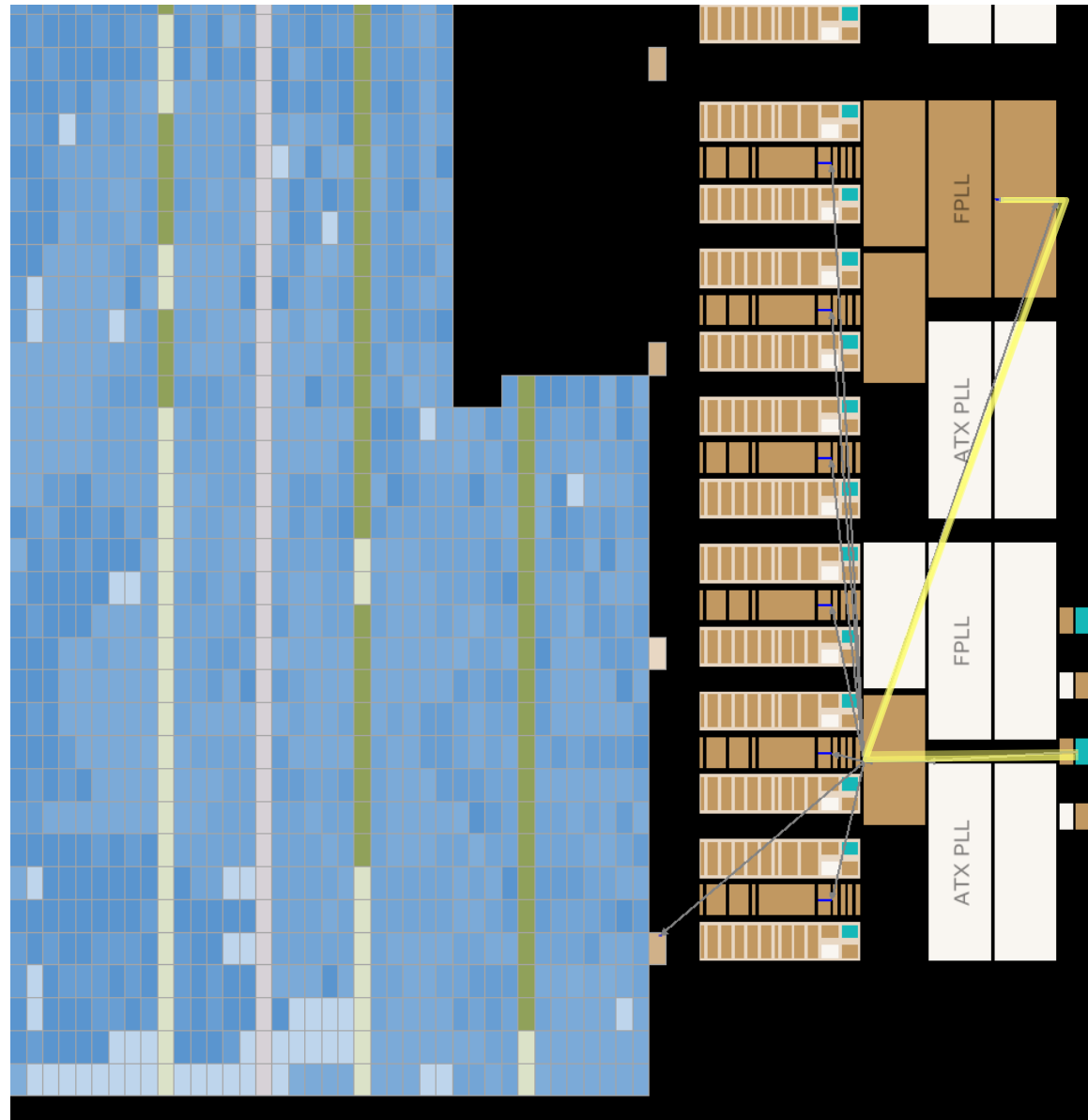
Figure 12 shows the simulated waveforms for three scenarios:

- Z-profile for "Ideal PCB"—The purple waveform (A) with resonance frequency F1
- Z-profile for "Low PCB parasitics"—The green waveform (B) with resonance frequency F2
- Z-profile for "High PCB parasitics"—The orange waveform (C) with resonance frequency F3

PDN parameters estimation



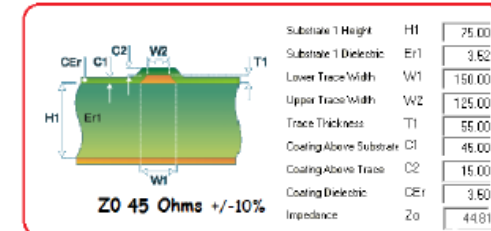
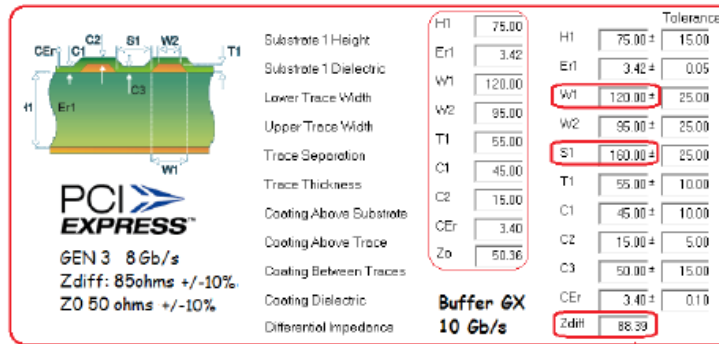
Correct clock routing



Stackup

PCIe40 V2-1 cards for LHCb
CERN market survey:
IT-4080/PH/LHCB
CPPM/IN2P3/CNRS

PCIe40 V 2-1 : ISOLA FR 408 HR Stackup

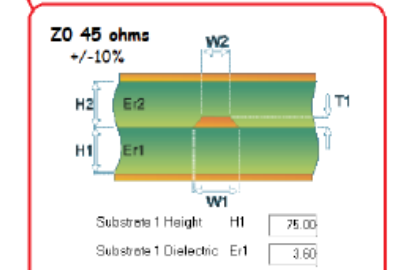
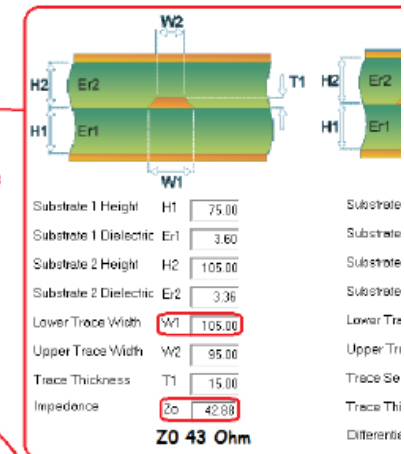


PCB stackup & manufacturing requirements

SM above CU: 20μ +/-15μ registration 50μ
CF: 12μ + plating = 55μ +/-10μ
1 x 1086 RC 65% = 75μ +/-15μ
CF: 12μ + plating = 40μ +/-10μ
2 x 106 RC 70% = 90μ +/- 13μ
CU: 17μ +/-5μ
CCL: 1 x1086 RC 58% = 75μ +/- 13μ
CU: 35μ +/-10μ
2 x 106 RC 70% = 85μ +/- 13μ
CU: 17μ +/-5μ
CCL: 1 x1086 RC 58% = 75μ +/- 13μ
CU: 35μ +/-10μ
2 x 106 RC 70% = 80μ +/-13μ
CU: 70μ +/-16μ
CCL: 1 x1086 RC 58% = 75μ +/- 13μ
CU: 70μ +/-16μ
2 x 106 RC 70% = 80μ +/-13μ
CU: 35μ +/-10μ
CCL: 1 x1086 RC 58% = 75μ +/- 13μ
CU: 17μ +/-5μ
2 x 106 RC 70% = 85μ +/- 13μ
CU: 35μ +/-10μ
CCL: 1 x1086 RC 58% = 75μ +/- 13μ
CU: 17μ +/-5μ
2 x 106 RC 70% = 90μ +/- 13μ
CF: 12μ + plating = 40μ +/-10μ
1 x 1086 RC 65% = 75μ +/-15μ
CF: 12μ + plating = 55μ +/-10μ
SM above CU: 20μ +/-15μ registration 50μ

CAD Design rules & Copper balancing requirements

Bridge mini 100μ, Registration : QFN 50μ, BGA 60μ
Top + W/Smini 110μ + GND = Layout 60% = L1
Full GND Smini 100μ = Layout 90% = L2
W/Smini 90μ + GND = Layout 50% = L3
Full GND Smini 100μ = Layout 90% = L4
Wmini 80μ, Smini 90μ + GND = Layout 60% = L5
VCC1/2 Smini 120μ = layout 80% = L6
VCC 0.9V + VCC 12V, Smini 240μ = layout 80% = L7
Full GND S mini 240μ = Layout 90% = L8
VCC5/6 Smini 120μ = layout 80% = L9
Wmini 80μ, Smini 90μ + GND = Layout 60% = L10
Full GND Smini 100μ = Layout 90% = L11
W/Smini 90μ + GND = Layout 50% = L12
Full GND Smini 100μ = Layout 90% = L13
Bottom + W/Smini 110μ + GND = Layout 70% = L14
Bridge mini 100μ, Registration : QFN 50μ, BGA 60μ



3D model

