

400GbE PCS Architectural Options

IEEE 400 Gb/s Ethernet Study Group

July 2013 Geneva

Mark Gustlin - Xilinx

Gary Nicholl - Cisco

Dave Ofelt - Juniper

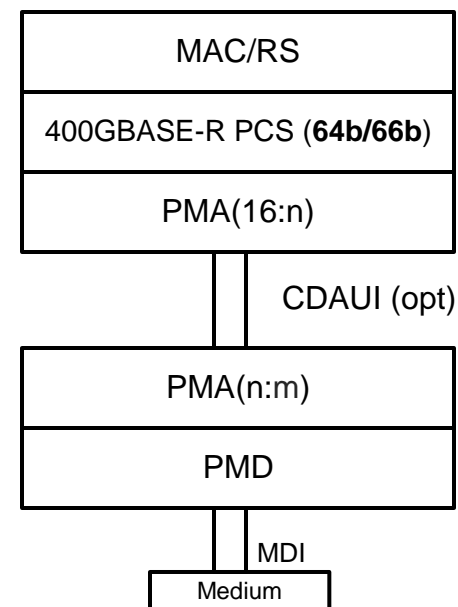
Introduction

- The following slides explore the feasibility of a 400GbE PCS
- Many feasible PCS architecture options are shown at 400GbE, building on the 802.3ba PCS and the work that has been done within P802.3bj so far

400GbE Possible Architecture #1

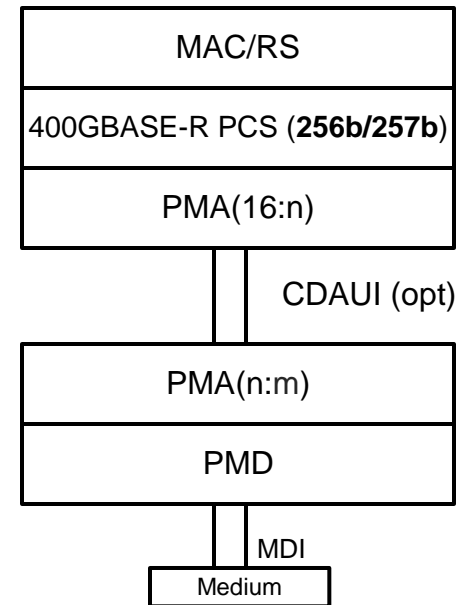
- Based on a 16 Lane PCS with 64B/66B encoding (25 Gb/s per PCS Lane)
- Data is striped to PCS lanes 66-bit blocks at a time
- Alignment Markers are periodically added to all PCS lanes to enable alignment in the RX PCS
- PMAs do simple bit multiplexing to change lane widths
- Lane widths of 16, 8, 4, 2, 1 can all be easily supported

- Pros of this architecture
 - Very flexible, can support future lane widths without a PCS change
 - Most of the complexity is in the PCS, PMAs are very simple bit multiplexers
 - Low latency solution
- Cons of this architecture
 - No low latency FEC
 - If FEC is added on then it likely requires transcoding, similar to 802.3bj
 - Susceptible to burst error MTTFFPA issues, when bit muxing PCS



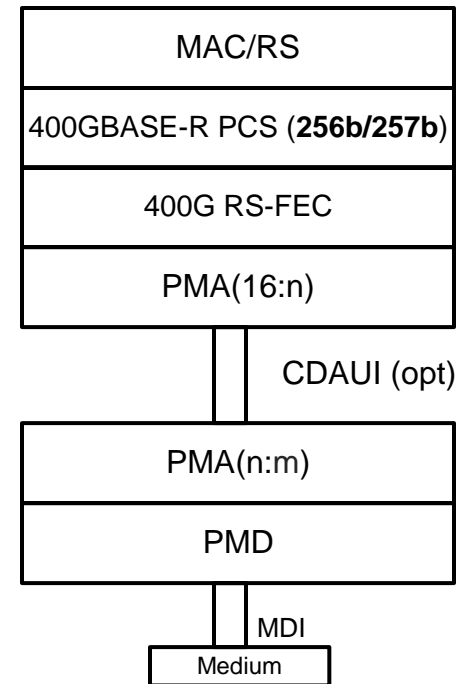
400GbE Possible Architecture #2

- Based on a 16 Lane PCS with 256B/257B encoding (25 Gb/s per PCS Lane)
- Data is striped to PCS lanes 257-bit blocks at a time
- Alignment Markers are periodically added to all PCS lanes to enable alignment in the RX PCS
- PMAs do simple bit multiplexing to change lane widths
- Lane widths of 16, 8, 4, 2, 1 can all be easily supported
- No FEC, but you can use the extra bits to add a robust checksum
 - 7bits per 257b block are available if running at 25.78125G per lane
 - Native rate without an additional checksum is 25.09765625G
- Pros of this architecture
 - Very flexible, can support future lane widths without a PCS change
 - Most of the complexity is in the PCS, PMAs are very simple bit multiplexers
 - Low latency solution
 - Robust error detection in the face of errors
- Cons of this architecture
 - No low latency FEC, but one can be added without transcoding
 - Different bit rate from today unless you add a checksum or other fill
 - Weak encoding is susceptible to MTTFPA issues unless robust checksum is added



400GbE Possible Architecture #3

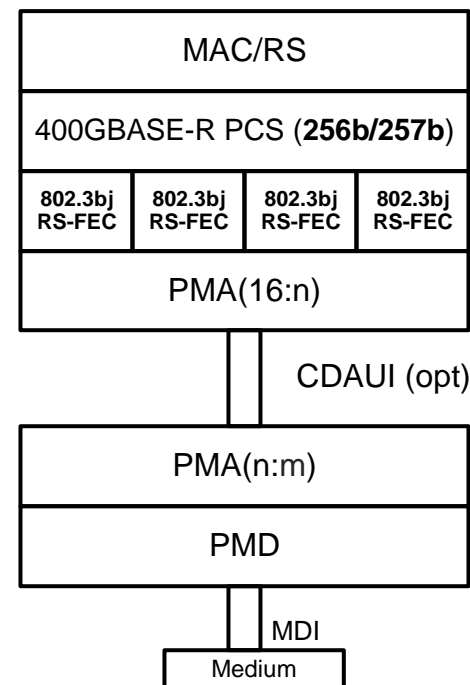
- Based on a 16 Lane PCS with 256B/257B encoding (25 Gb/s per PCS Lane)
- Data is striped to PCS lanes 257-bit blocks at a time?
 - Or distributed in RS symbol boundaries (10b for instance)?
- Alignment Markers are periodically added to all PCS lanes to enable alignment in the RX FEC block and PCS
- 400 Gb/s RS-FEC is added, no transcoding
 - Other FEC options should be explored also, what error signatures do we expect for electrical and optical lanes?
- PMAs do block multiplexing to change lane widths to preserve the error detection capability in the face of burst errors
 - You can do bit multiplexing if burst errors is not a concern
- Lane widths of 16, 8, 4, 2, 1 can all be supported
- Pros of this architecture
 - Very flexible, can support future lane widths without a PCS change
 - A lot of the complexity is in the PCS, PMAs though do have to find AM lock before muxing to preserve error correction capability
 - Pretty low latency solution (~25ns of added latency due to FEC, depends on block size though)
 - Robust error detection correction
- Cons of this architecture
 - Limited re-use from 100GbE



400GbE Possible Architecture #4

- Based on a 16 Lane PCS with 256B/257B encoding (25 Gb/s per PCS Lane)
- Data is striped to PCS lanes 257-bit blocks at a time
 - Or distributed in RS symbol boundaries (10b for instance)?
- Alignment Markers are periodically added to all PCS lanes to enable alignment in the RX FEC block and PCS
 - Need 16 unique AMs, unlike 802.3bj
- A portion of the 100 Gb/s RS-FEC x 4 is added, no transcoding
- PMAs do simple block multiplexing to change lane widths to preserve the error detection capability in the face of burst errors
 - You can do bit multiplexing if burst errors is not a concern
- Lane widths of 16, 8, 4, 2, 1 can all be supported

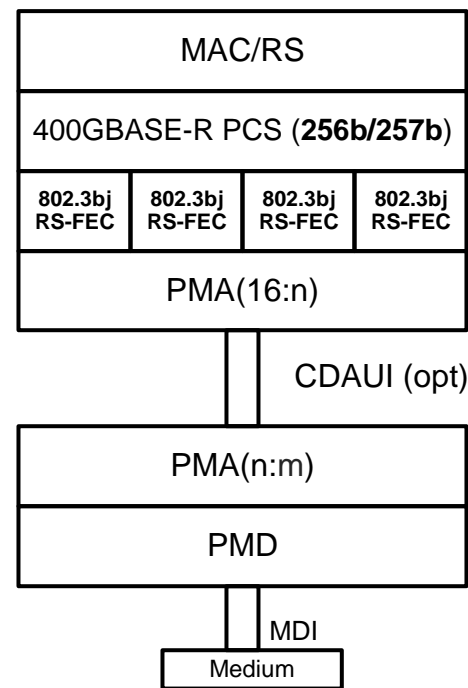
- Pros of this architecture
 - Very flexible, can support future lane widths without a PCS change
 - A lot of the complexity is in the PCS, PMAs though do have to find AM lock before muxing to preserve error correction capability
 - Re-use of some of the 802.3bj RS-FEC
 - Robust error detection correction
- Cons of this architecture
 - Higher latency than the other options, ~100ns with correction



400GbE Possible Architecture #5

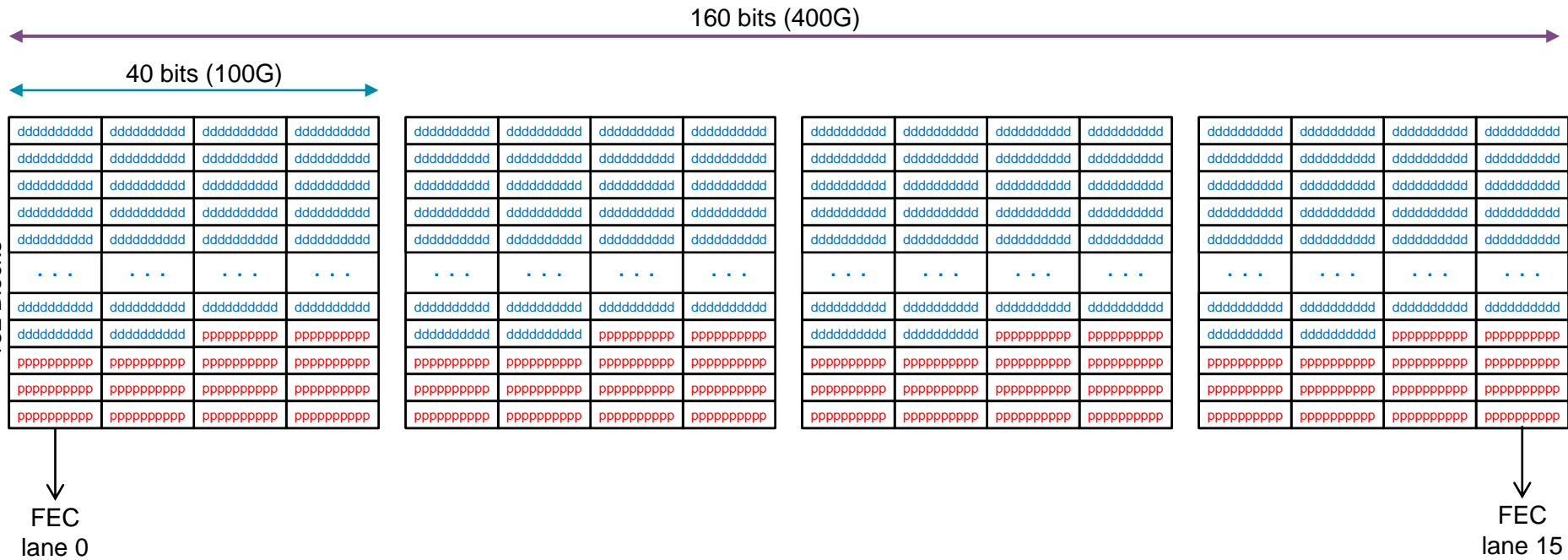
- Based on a 80 Lane PCS with 256B/257B encoding (25 Gb/s per PCS Lane)
- Data is distributed on RS symbol boundaries (10b in this case)
- Alignment Markers are periodically added to all PCS lanes to enable alignment in the RX FEC block and PCS
 - Need 80 unique AMs, unlike 802.3bj
- A portion of the 100 Gb/s RS-FEC x 4 is added, no transcoding
- Because 802.3bj FEC has only 4 FEC lanes towards the PMA, not sure how you would distribute/create 80 lanes below the FEC sublayer when there are 16 FEC lanes?
- Lane widths of ? can all be supported?

- Pros of this architecture
 - Consistent with today's 100GbE architecture
- Cons of this architecture
 - Not sure how 80 PCS lanes helps you when you only have 16 FEC lanes?



400GbE Possible Architecture #5 - Continued

- You could have 80 PCS lanes above the FEC sublayer, but does that help at all?
- Below the FEC sublayer, with using 4x802.3bj FEC, you would naturally have 16 FEC lanes
- If you want to distribute to 80 FEC lanes how would you divide up the data?
- If you broke each 10b block into 5 pieces then you could distribute to 80 FEC lanes, but you would need to ensure the AM gets mapped and broken up correctly, also in the face of burst errors you would degrade error correcting performance depending how the bits are later multiplexed



Stronger FEC

- With option 2,3 or 4, if a stronger FEC is needed than the base FEC (or no FEC in the case of option 2), you simply add the FEC on top of what is already there, no transcoding is needed
- You can also strip off the current FEC and then add a stronger FEC to the PCS encoded data, again without having to do transcoding

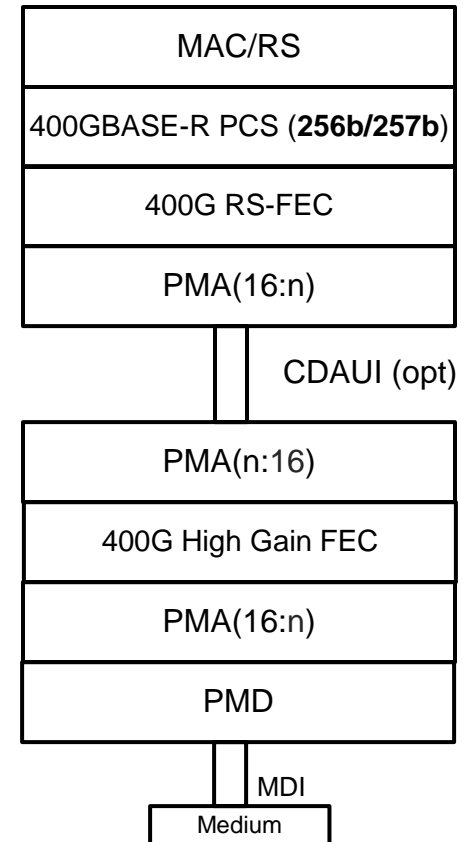


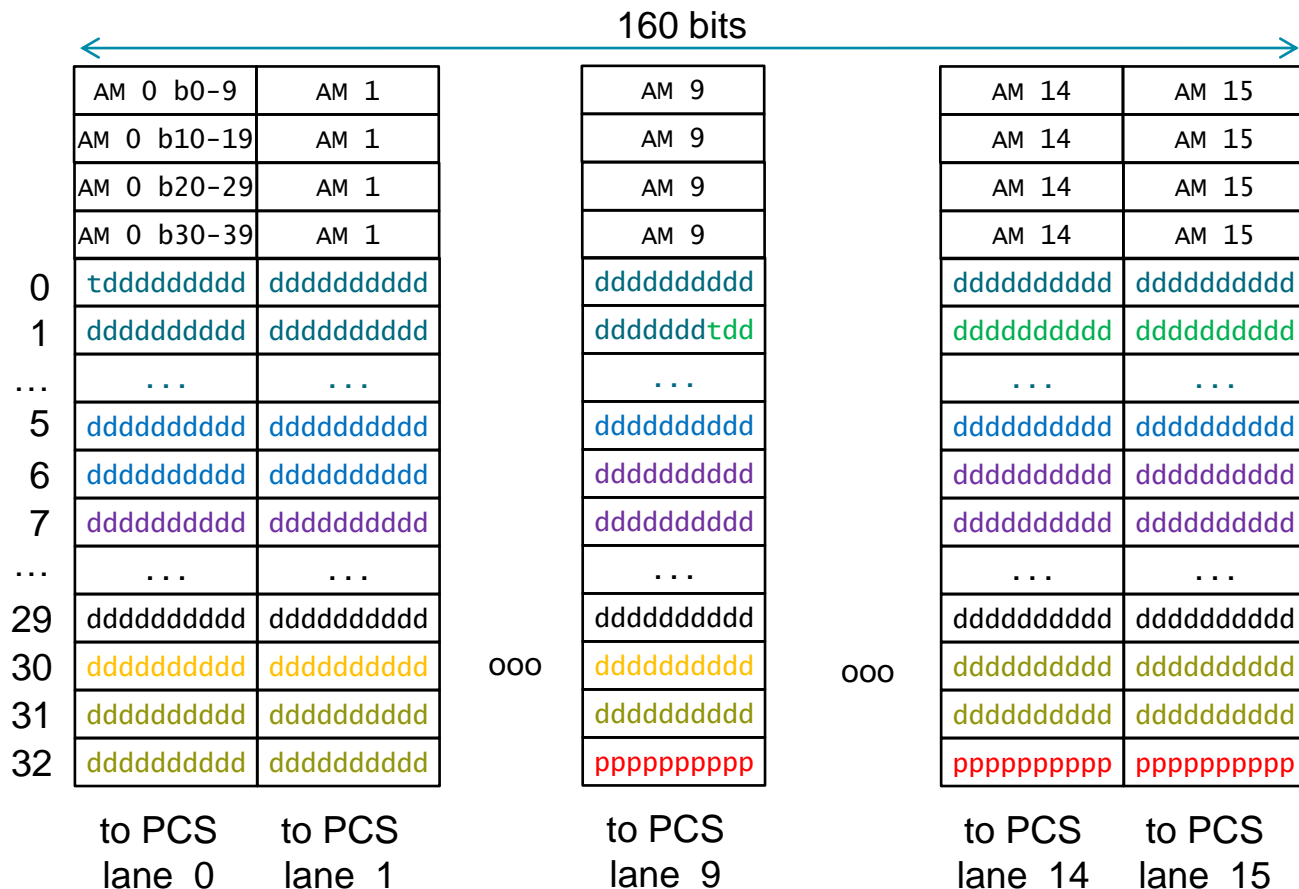
Table of Options

Option	Encoding	# PCS Lanes	FEC?	PMA Width Change	Added Latency	MTTFPA concerns?
1	64B/66B	16	None, adding FEC likely requires transcoding	Bit muxing	0	Yes when bit muxing + burst errors
2	256B/257B	16	None, but ready for FEC	Bit muxing	~5ns	Yes, especially when bit muxing + burst errors
3	256B/257B	16	400G FEC	Block muxing	~25-50ns?	No
4	256B/257B	16	4xRS-FEC	Block muxing	~100ns	No
4	256B/257B	80	4xRS-FEC	Block muxing	~100ns	No

Note: At this point supporting 10x40G lanes are not addressed with these options

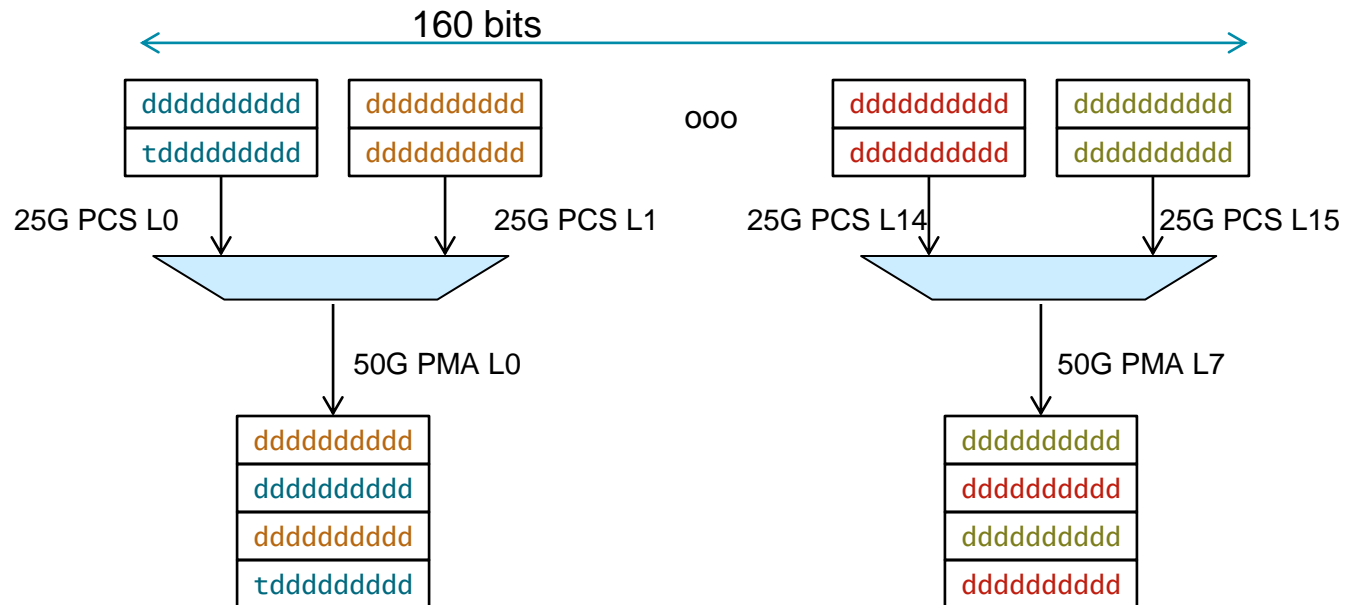
Alignment Markers

- Add an alignment marker to each PCS lane periodically, it does not need to be part of the FEC blocks, and it seems to make it easier if they are not part of the FEC block (so you don't have Alignment issues)
- Below the AMs are 40 bits each, but this is flexible, just must be nx10-bit



Multiplexing

- With 16 PCS lanes, you can multiplex down to 8, 4, 2, or 1 lane(s)
- Multiplexing is typically done on RS boundaries (10-bit in the case shown)
 - To preserve error correction capability in the face of burst errors
 - If you are running across a medium that only has uncorrelated errors, then you have the option of doing bit multiplexing
- First you must find alignment marker lock to find 10-bit boundaries, then you multiplex on RS boundaries
 - No need to deskew the various lanes
- Below shows muxing from 16 lanes down to 8 lanes



Things to Look Into

- Does one sized FEC fit for all applications, if not how to handle multiple FECs?
- Should the architecture allow a given FEC to be physically separated from the PCS?
- Will FEC be required for electrical interfaces?
 - If yes, how is the coding gain partitioned between the electrical interface and optical?
- What is the FEC equivalent of an FEC high BER?
- Future complex modulation PMDs will require high gain FEC, how does that fit into the architecture?
- EEE interactions
- Support for OTN, what does that mean to the various options
- How to do Alignment Marker mapping for some of the options
- How to distribute data for some of the options
- How to address 10x40G lanes? Is it needed?
- How does this fit in with the MLG protocol?

Summary

- There are many possible solutions for a 400GbE PCS, this paper shows a couple of options that are feasible with today's technology (either ASIC or FPGA)
- One simple option is scaling the 802.3ba PCS up in speed
- But if there will be interfaces that require FEC, and low latency is important, then a PCS could be defined that incorporates a low latency FEC from the start
 - This applies to both electrical and optical interfaces

Thanks!