# 800GbE PCS/FEC/PMA Baseline Proposal for PHYs using 8 x 100G PMD lanes

Kapil Shrikhande (Marvell), Eugene Opsasnick (Broadcom), Gary Nicholl (Cisco), Matt Brown (Huawei), Xinyuan Wang (Huawei), Mark Gustlin (Cisco), David Ofelt (Juniper), Eric Maniloff (Ciena), Shawn Nicholl (AMD), Jeff Slavick (Broadcom), Adee Ran (Cisco), Kent Lusted (Intel), Jerry Pepper (Keysight)

October 2022

IEEE 802.3df Task Force

1

# Supporters

- Rob Stone, Meta
- Lenin Patra, Marvell
- Arthur Marris, Cadence
- Howard Heck, Intel
- Venu Balasubramonian, Marvell
- Leon Bruckman, Huawei
- Yan Zhuangyan, Huawei
- Chris Cole, Quintessent
- Jeffery Maki, Juniper Networks
- Nathan Tracy, TE Connectivity
- Daniel Koehler, Synopsys
- David Malicoat, Malicoat Networking Solutions
- Frank Effenberger, Huawei
- Joshua Kim, Hirose Electric
- Kenneth Jackson, Sumitomo Electric
- Rick Rabinovich, Keysight
- Vipul Bhatt, II-VI
- Tom Palkert, MACOM and Samtec
- Mau-Lin Wu, Mediatek

- Ted Sprague, Infinera
- Dave Estes, Spirent
- Ed Nakamoto, Spirent
- Chris DiMinico, PHY-SI/SenTekse
- Ben Jones, AMD
- Ali Ghiasi, Ghiasi Quantum LLC
- Paul Brooks, Viavi Solutions
- Megha Shanbhag, TE Connectivity
- Pavel Zivny, Tektronix
- Mike Dudek, Marvell
- Shimon Muller, Enfabrica Corp.
- Xiang He, Huawei
- Viet Tran, Keysight
- Brian Welch, Cisco
- Phil Sun, Credo
- Lokesh Kabra, Synopsys
- Roberto Rodes, II-VI
- Brad Booth, Microsoft

# Outline

- **Introduction**
- PCS/FEC/PMA Baseline proposal
- Conclusions

# Goals

- Fast time to an 800GbE PCS/FEC/PMA specification for PMDs and AUIs using 100G/lane
  - Reuse 400GbE PCS/FEC (CL119) as much as possible
  - Support 800GbE with simple modification to the 400GbE PCS/FEC
  - Leverage 802.3bs CL120 PMA; leverage 802.3ck 100G/lane PMD and AUI specifications

- Maximize the reuse of existing logic sub-blocks used in 400GbE PCS/FEC
  - Leverage industry investment in 400GbE technology

- Enable systems using current 8-lane 800G connectors (OSFP / QSFP-DD) to also support 800GbE
  - E.g. 8-lane C2M AUIs used as : 8 x 100GAUI-1 / 4 x 200GAUI-2 / 2 x 400GAUI-4 and 1 x 800GAUI-8

# Scope

**802.3df Adopted PHY Objectives***

| Ethernet Rate | Assumed Signaling Rate | AUI | BP | Cu Cable | MMF 50m | MMF 100m | SMF 500m | SMF 2km | SMF 10km | SMF 40km |
|---|---|---|---|---|---|---|---|---|---|---|
| 200 Gb/s | 200 Gb/s | Over 1 lane | | Over 1 pair | | | Over 1 Pair | Over 1 Pair | | |
| 400 Gb/s | 200 Gb/s | Over 2 lanes | | Over 2 pairs | | | | Over 2 Pair | | |
| 800 Gb/s | 100 Gb/s | Over 8 lanes | Over 8 lanes | Over 8 pairs | Over 8 pairs | Over 8 pairs | Over 8 pairs | Over 8 pairs | | |
| | 200 Gb/s | Over 4 lanes | | Over 4 pairs | | | Over 4 pairs | 1) Over 4 pairs 2) Over 4 λ's | | |
| | TBD | | | | | | | | Over single SMF in each direction | Over single SMF in each direction |
| 1.6 Tb/s | 100 Gb/s | Over 16 lanes | | | | | | | | |
| | 200 Gb/s | Over 8 lanes | | Over 8 pairs | | | Over 8 pairs | Over 8 pairs | | |

**Making it all work together**

## Technology Reuse

- Leverage existing or work-in-progress 100 Gb/s per lane (e.g. 3cu, 3ck, 3db) to higher lane counts
- Develop 200 Gb/s per lane electrical signaling for 1/2/4/8 lane variants of AUIs and electrical PMDs
- Develop 200 Gb/s per optical fiber for 1/2/4/8 fiber based optical PMDs and 4 lambda WDM optical PMD
- Potential for either direct detect and / or coherent signaling technology

**Scope of this Baseline : 800GbE PCS/FEC/PMA for all PHY objectives that use 8 x 100G PMDs and AUIs**

*\* Table from https://www.ieee802.org/3/B400G/public/21_1028/B400G_overview_c_211028.pdf*
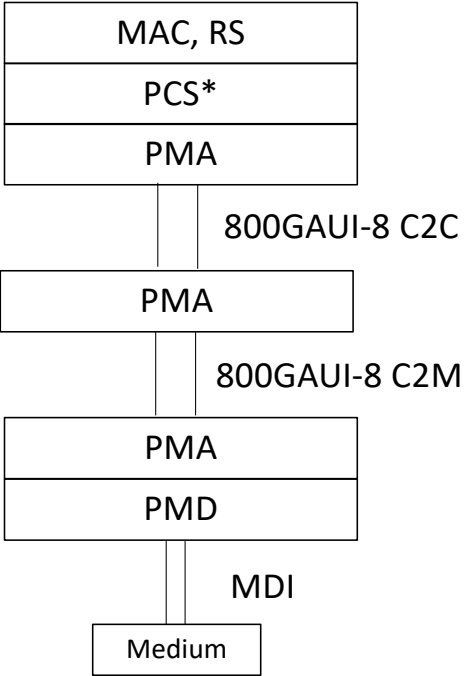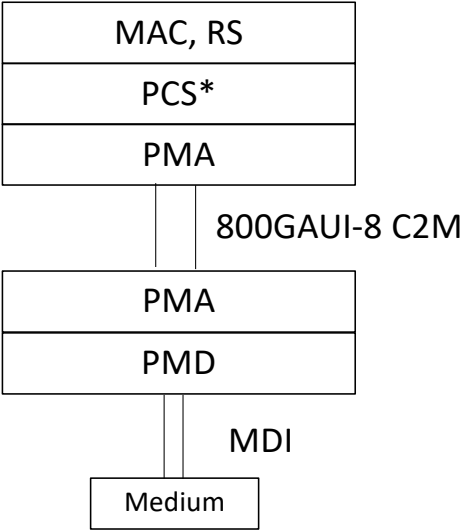
# AUI and PMD assumptions

- 802.3df Task Force has adopted 800GbE 8-lane AUI baseline proposals leveraging existing 100G/lane AUI specs, drafts
    - https://www.ieee802.org/3/df/public/22_03/lusted_3df_01a_220315.pdf

- 802.3df Task Force has adopted 800GbE 8-lane PMD baseline proposals leveraging existing 100G/lane PMD specs, drafts
    - https://www.ieee802.org/3/df/public/22_03/lusted_3df_01a_220315.pdf
    - https://www.ieee802.org/3/df/public/22_02/welch_3df_01a_220222.pdf
    - https://www.ieee802.org/3/df/public/22_03/murty_3df_01a_220315.pdf

- 802.3bs CL119 PCS works for all 100G/lane AUIs and PMDs for 400GbE

- Similarly, this PCS/FEC Baseline (leveraging CL119) works for all adopted 800GbE 8-lane AUIs and PMDs

# Outline

- Introduction
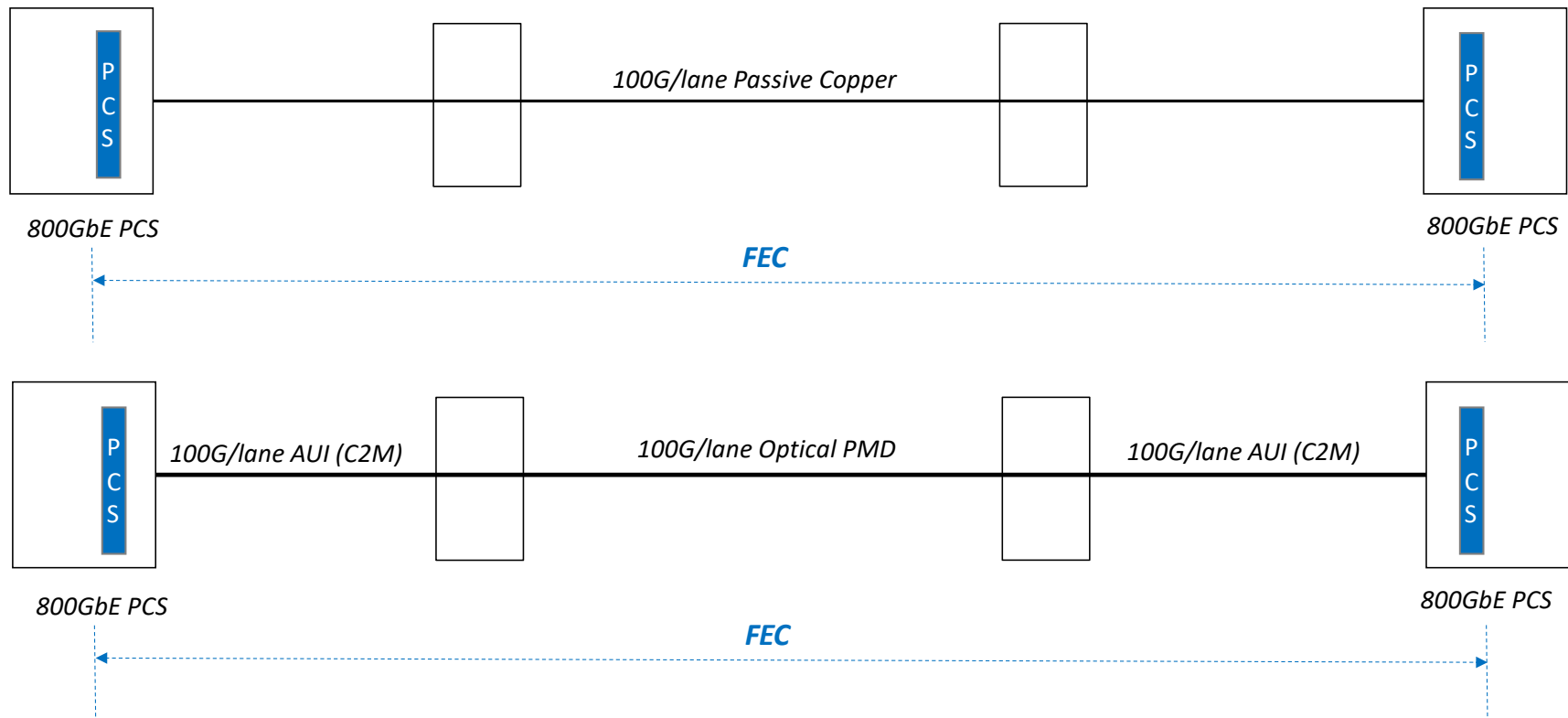- **PCS/FEC/PMA Baseline proposal**
- Conclusions

# Architecture

| MAC, RS |
|---------|
| PCS* |
| PMA |

800GAUI-8 C2M

| PMA |
|-----|
| PMD |

MDI

| Medium |
|--------|

| MAC, RS |
|---------|
| PCS* |
| PMA |

800GAUI-8 C2C

| PMA |
|-----|

800GAUI-8 C2M

| PMA |
|-----|
| PMD |

MDI

| Medium |
|--------|

*PCS and FEC are in the PCS sublayer (same as CL119)

Note : Not showing layering diagram for Cu PMD (will be same as other Cu PMD layering diagrams in 802.3)

# End-End PCS/FEC scheme for 800GbE (8 x 100G) PMDs



800GbE PCS

100G/lane Passive Copper

800GbE PCS

FEC

800GbE PCS

100G/lane AUI (C2M)

100G/lane Optical PMD

100G/lane AUI (C2M)

800GbE PCS

FEC

Note : This End-End PCS/FEC works with optional Chip to Chip AUIs and a combination of Chip to chip and Chip to module *(same as 400GAUI-4 in 802.3ck)*
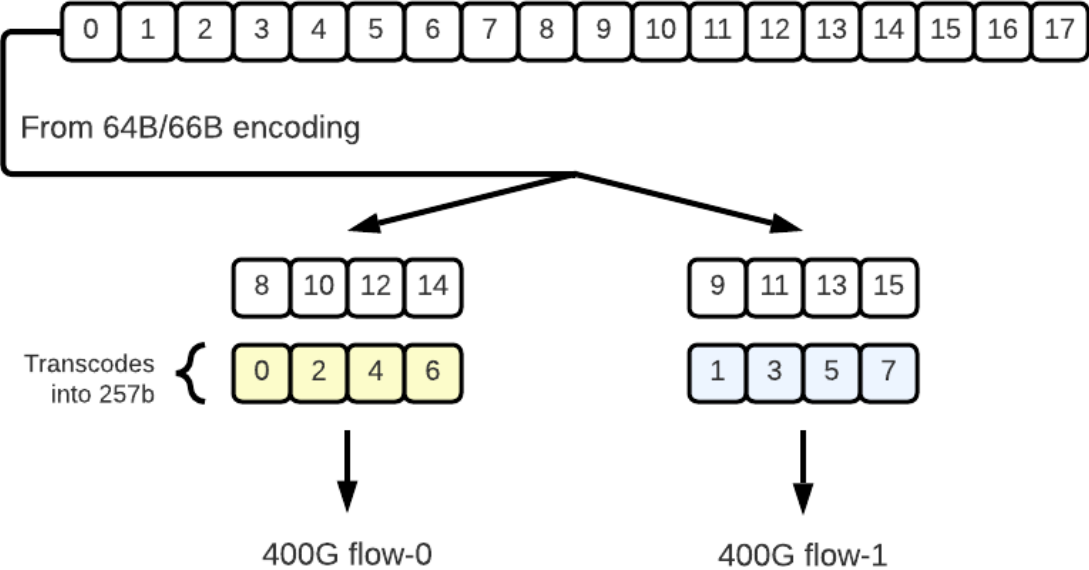
# Tx PCS/FEC Data Flow

- Based on two 802.3bs, CL119 sublayers in parallel
  - Two 400G FEC flows (flow-0 and flow-1)
- 66b round robin distribution into two 400G flows after 64B/66B encode
- Sub-blocks shown within each flow are identical to CL119, except :
  - AM values are made unique across the two flows
  - AM insertion is aligned across the two flows
- 32 Flow lanes per 800GbE PCS
  - 16 per 400G flow
- Specific Flow lanes mapped to a given PMA output lane
  - 4:1 bit-muxing
  - Lanes chosen so all 4 FEC codewords are equally represented on each PMA output lane
  - Bitmux can be specified to occur in either the PCS or PMA sublayer (TBD).



**800G MII**

TXD<63:0>
TXC<7:0>
TX_CLK

**800G PCS**

64B/66B Encode and rate match

66b Block Distribution

| 400G flow-0 | 400G flow-1 |
|---|---|
| 256B/257B Transcode | 256B/257B Transcode |
| Scramble | Scramble |

Alignment Insertion

| Pre-FEC distribution | Pre-FEC distribution |
| FEC Encode | FEC Encode |
| Distribution and Interleave | Distribution and Interleave |

802.3bs

FlowLane[0:15]      FlowLane[16:31]

PCS or PMA sublayer
(4:1 bitmuxing)

# Tx 66b Block Distribution

- Round Robin among two '400G Flows'

# Alignment Marker Insertion



Figure 119–8—400GBASE-R alignment marker insertion period

Source: IEEE Std 802.3-2018
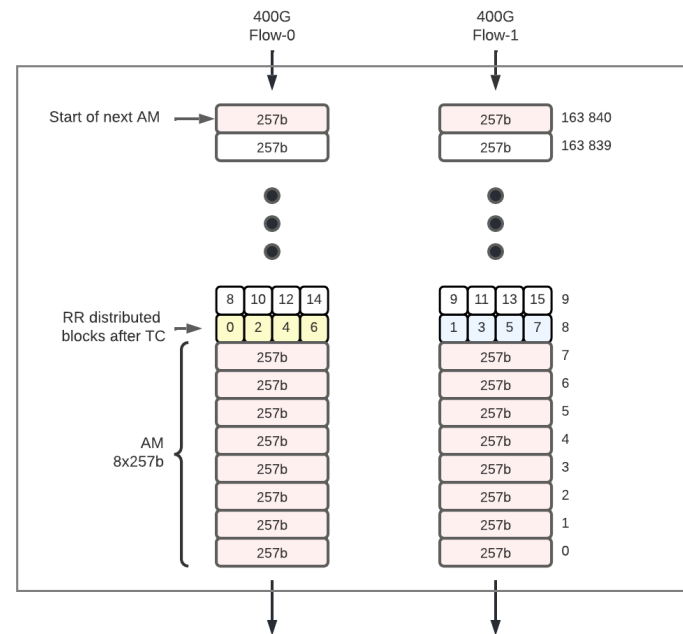
- 802.3bs 400G AM structure
  - AM size = 8 x 257b
  - Spacing = 160k x 257b = 8192 CWs
- AM total sizing for 800G = 2x400G
  - AM size = 16 x 257b
  - Spacing = 320k x 257b = 16384 CWs
- Markers inserted at consecutive 257b blocks across both 400G flows
  - Flow-0 is first in time carrying the even encoded 4x66b blocks
  - Flow-1 carries odd encoded 4x66b blocks
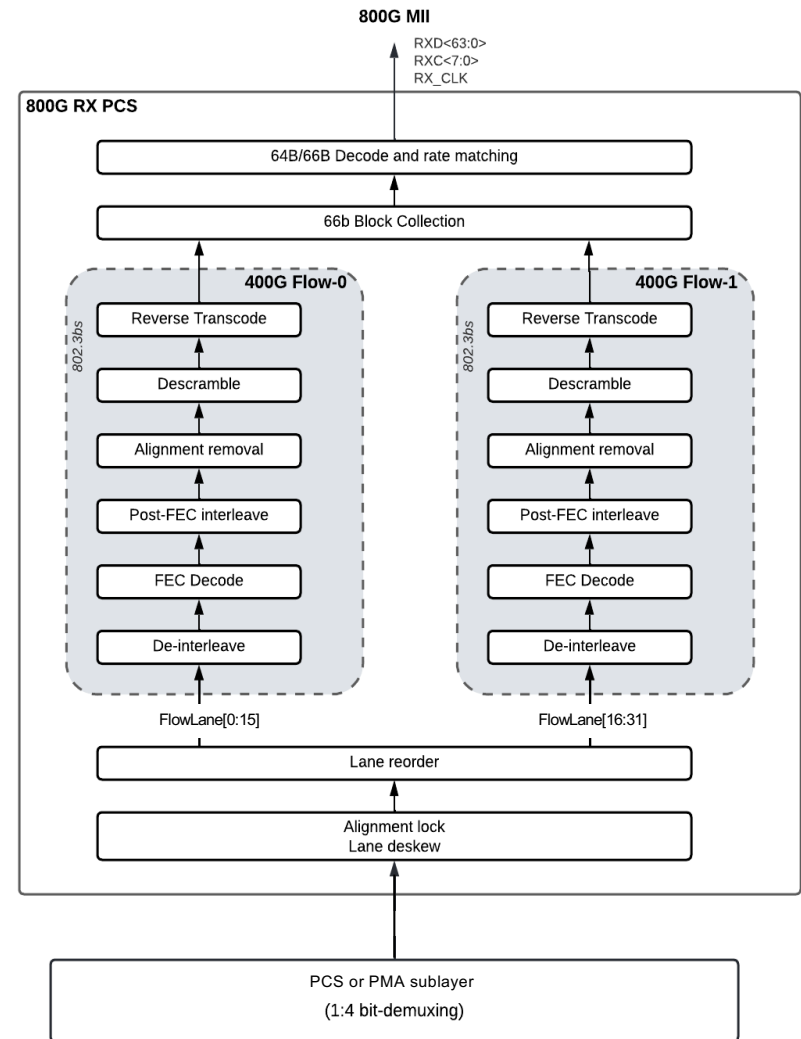
# AM Marker Encoding

- CM0-CM5 and UP0-UP2 are unchanged from 400GbE CL119

- UM0/UM3 for Flow lanes 0-15 are inverted from 400GbE

- UM1/UM2/UM4/UM5 for Flow lanes 16-31 are inverted from 400GbE

- Prevents lock with 400GbE ports

- Maintains DC balance

| Flow Lane # | Encoding | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CM0 | CM1 | CM2 | UP0 | CM3 | CM4 | CM5 | UP1 | UM0 | UM1 | UM2 | UP2 | UM3 | UM4 | UM5 |
| 0 | 0x9A | 0x4A | 0x26 | 0xB6 | 0x65 | 0xB5 | 0xD9 | 0xD9 | **0xFE** | 0x71 | 0xF3 | 0x26 | **0x01** | 0x8E | 0x0C |
| 1 | 0x9A | 0x4A | 0x26 | 0x04 | 0x65 | 0xB5 | 0xD9 | 0x67 | **0xA5** | 0xDE | 0x7E | 0x98 | **0x5A** | 0x21 | 0x81 |
| 2 | 0x9A | 0x4A | 0x26 | 0x46 | 0x65 | 0xB5 | 0xD9 | 0xFE | **0xC1** | 0xF3 | 0x56 | 0x01 | **0x3E** | 0x0C | 0xA9 |
| 3 | 0x9A | 0x4A | 0x26 | 0x5A | 0x65 | 0xB5 | 0xD9 | 0x84 | **0x79** | 0x80 | 0xD0 | 0x7B | **0x86** | 0x7F | 0x2F |
| 4 | 0x9A | 0x4A | 0x26 | 0xE1 | 0x65 | 0xB5 | 0xD9 | 0x19 | **0xD5** | 0x51 | 0xF2 | 0xE6 | **0x2A** | 0xAE | 0x0D |
| 5 | 0x9A | 0x4A | 0x26 | 0xF2 | 0x65 | 0xB5 | 0xD9 | 0x4E | **0xED** | 0x4F | 0xD1 | 0xB1 | **0x12** | 0xB0 | 0x2E |
| 6 | 0x9A | 0x4A | 0x26 | 0x3D | 0x65 | 0xB5 | 0xD9 | 0xEE | **0xBD** | 0x9C | 0xA1 | 0x11 | **0x42** | 0x63 | 0x5E |
| 7 | 0x9A | 0x4A | 0x26 | 0x22 | 0x65 | 0xB5 | 0xD9 | 0x32 | **0x29** | 0x76 | 0x5B | 0xCD | **0xD6** | 0x89 | 0xA4 |
| 8 | 0x9A | 0x4A | 0x26 | 0x60 | 0x65 | 0xB5 | 0xD9 | 0x9F | **0x1E** | 0x73 | 0x75 | 0x60 | **0xE1** | 0x8C | 0x8A |
| 9 | 0x9A | 0x4A | 0x26 | 0x6B | 0x65 | 0xB5 | 0xD9 | 0xA2 | **0x8E** | 0xC4 | 0x3C | 0x5D | **0x71** | 0x3B | 0xC3 |
| 10 | 0x9A | 0x4A | 0x26 | 0xFA | 0x65 | 0xB5 | 0xD9 | 0x04 | **0x6A** | 0xEB | 0xD8 | 0xFB | **0x95** | 0x14 | 0x27 |
| 11 | 0x9A | 0x4A | 0x26 | 0x6C | 0x65 | 0xB5 | 0xD9 | 0x71 | **0xDD** | 0x66 | 0x38 | 0x8E | **0x22** | 0x99 | 0xC7 |
| 12 | 0x9A | 0x4A | 0x26 | 0x18 | 0x65 | 0xB5 | 0xD9 | 0x5B | **0x5D** | 0xF6 | 0x95 | 0xA4 | **0xA2** | 0x09 | 0x6A |
| 13 | 0x9A | 0x4A | 0x26 | 0x14 | 0x65 | 0xB5 | 0xD9 | 0xCC | **0xCE** | 0x97 | 0xC3 | 0x33 | **0x31** | 0x68 | 0x3C |
| 14 | 0x9A | 0x4A | 0x26 | 0xD0 | 0x65 | 0xB5 | 0xD9 | 0xB1 | **0x35** | 0xFB | 0xA6 | 0x4E | **0xCA** | 0x04 | 0x59 |
| 15 | 0x9A | 0x4A | 0x26 | 0xB4 | 0x65 | 0xB5 | 0xD9 | 0x56 | **0x59** | 0xBA | 0x79 | 0xA9 | **0xA6** | 0x45 | 0x86 |
| 16 | 0x9A | 0x4A | 0x26 | 0xB6 | 0x65 | 0xB5 | 0xD9 | 0xD9 | 0x01 | **0x8E** | **0x0C** | 0x26 | 0xFE | **0x71** | **0xF3** |
| 17 | 0x9A | 0x4A | 0x26 | 0x04 | 0x65 | 0xB5 | 0xD9 | 0x67 | 0x5A | **0x21** | **0x81** | 0x98 | 0xA5 | **0xDE** | **0x7E** |
| 18 | 0x9A | 0x4A | 0x26 | 0x46 | 0x65 | 0xB5 | 0xD9 | 0xFE | 0x3E | **0x0C** | **0xA9** | 0x01 | 0xC1 | **0xF3** | **0x56** |
| 19 | 0x9A | 0x4A | 0x26 | 0x5A | 0x65 | 0xB5 | 0xD9 | 0x84 | 0x86 | **0x7F** | **0x2F** | 0x7B | 0x79 | **0x80** | **0xD0** |
| 20 | 0x9A | 0x4A | 0x26 | 0xE1 | 0x65 | 0xB5 | 0xD9 | 0x19 | 0x2A | **0xAE** | **0x0D** | 0xE6 | 0xD5 | **0x51** | **0xF2** |
| 21 | 0x9A | 0x4A | 0x26 | 0xF2 | 0x65 | 0xB5 | 0xD9 | 0x4E | 0x12 | **0xB0** | **0x2E** | 0xB1 | 0xED | **0x4F** | **0xD1** |
| 22 | 0x9A | 0x4A | 0x26 | 0x3D | 0x65 | 0xB5 | 0xD9 | 0xEE | 0x42 | **0x63** | **0x5E** | 0x11 | 0xBD | **0x9C** | **0xA1** |
| 23 | 0x9A | 0x4A | 0x26 | 0x22 | 0x65 | 0xB5 | 0xD9 | 0x32 | 0xD6 | **0x89** | **0xA4** | 0xCD | 0x29 | **0x76** | **0x5B** |
| 24 | 0x9A | 0x4A | 0x26 | 0x60 | 0x65 | 0xB5 | 0xD9 | 0x9F | 0xE1 | **0x8C** | **0x8A** | 0x60 | 0x1E | **0x73** | **0x75** |
| 25 | 0x9A | 0x4A | 0x26 | 0x6B | 0x65 | 0xB5 | 0xD9 | 0xA2 | 0x71 | **0x3B** | **0xC3** | 0x5D | 0x8E | **0xC4** | **0x3C** |
| 26 | 0x9A | 0x4A | 0x26 | 0xFA | 0x65 | 0xB5 | 0xD9 | 0x04 | 0x95 | **0x14** | **0x27** | 0xFB | 0x6A | **0xEB** | **0xD8** |
| 27 | 0x9A | 0x4A | 0x26 | 0x6C | 0x65 | 0xB5 | 0xD9 | 0x71 | 0x22 | **0x99** | **0xC7** | 0x8E | 0xDD | **0x66** | **0x38** |
| 28 | 0x9A | 0x4A | 0x26 | 0x18 | 0x65 | 0xB5 | 0xD9 | 0x5B | 0xA2 | **0x09** | **0x6A** | 0xA4 | 0x5D | **0xF6** | **0x95** |
| 29 | 0x9A | 0x4A | 0x26 | 0x14 | 0x65 | 0xB5 | 0xD9 | 0xCC | 0x31 | **0x68** | **0x3C** | 0x33 | 0xCE | **0x97** | **0xC3** |
| 30 | 0x9A | 0x4A | 0x26 | 0xD0 | 0x65 | 0xB5 | 0xD9 | 0xB1 | 0xCA | **0x04** | **0x59** | 0x4E | 0x35 | **0xFB** | **0xA6** |
| 31 | 0x9A | 0x4A | 0x26 | 0xB4 | 0x65 | 0xB5 | 0xD9 | 0x56 | 0xA6 | **0x45** | **0x86** | 0xA9 | 0x59 | **0xBA** | **0x79** |

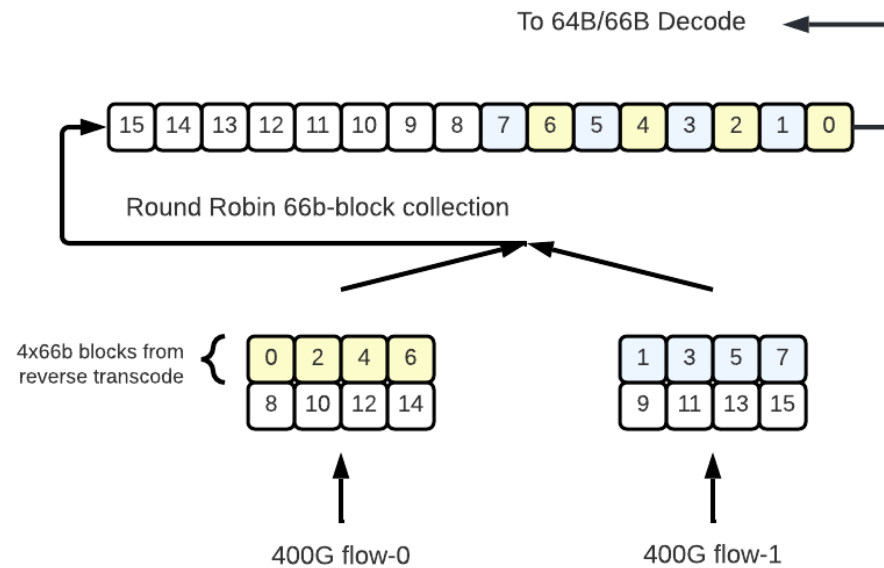Note: in table above, bolded text indicates changes from CL 119 AM values

# Rx PCS/FEC Data Flow

- Alignment Lock and Deskew
  - AM lock : per lane, same as CL119
  - De-skew : across 32 PCS lanes

- Lane reorder (and split)
  - Reorder and split 32 PCS lanes into 2 groups of 16
    - Lanes 0-15 : Flow-0
    - Lanes 16-31 : Flow-1

- FEC decode, de-scramble, transcode decode – same as CL119

- Round robin block collection must be aligned across Flow-0/1 based on Alignment Marker location

# Rx 66b Block Collection

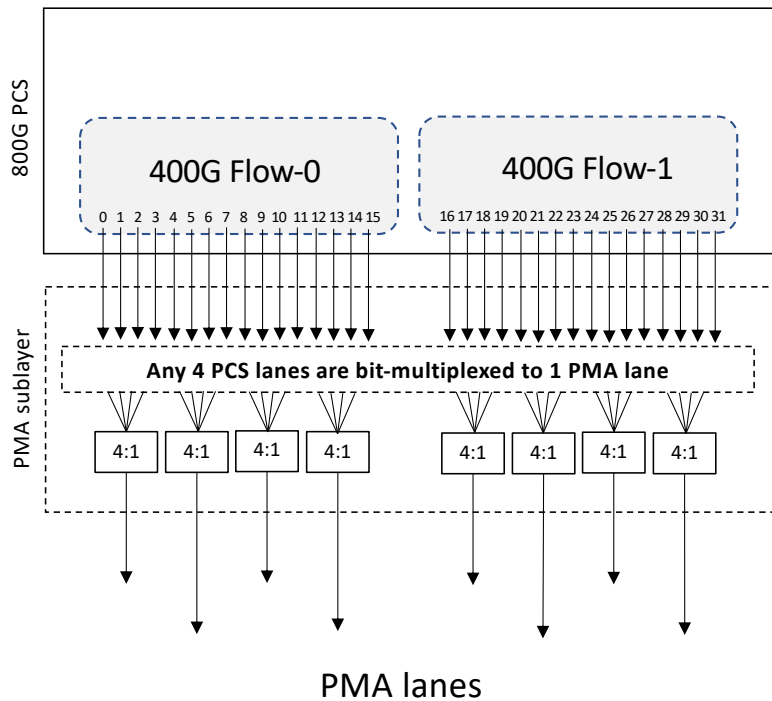- Round Robin 66b Block Collection is opposite of Tx Block Distribution

# Reuse CL119 State Diagrams

- Reuse the following
  - Figure 119–12—Alignment marker lock state diagram
  - Figure 119–13—PCS synchronization state diagram
  - Figure 119–14—Transmit state diagram
  - Figure 119–15—Receive state diagram

- Minor modification to the following
  - Add restart_lock<y> variable per 400G flow
    - restart_lock =  restart_lock<0> OR restart_lock<1>
  - Add hi_ser<y> variable per 400G flow
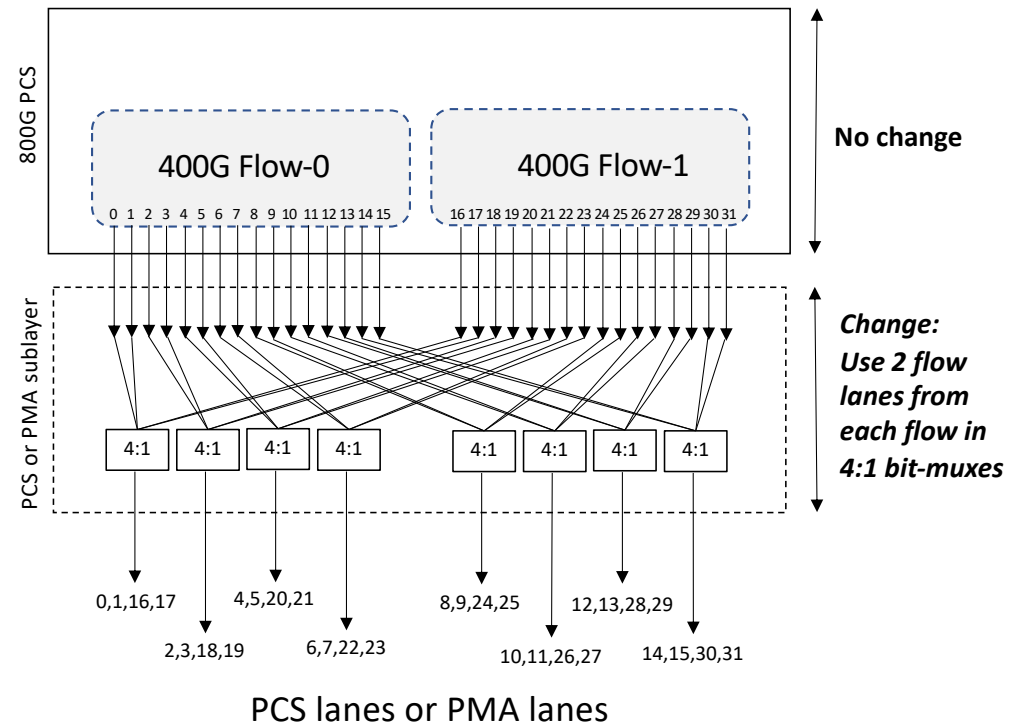    - hi_ser = hi_ser<0> OR hi_ser<1>

# Flow lane Muxing

- 32 Flow Lanes to 8 PMA Lanes such that

  - Each PMA lane is a result of bitmux of 2 flow lanes from Flow 0 and 2 flow lanes from Flow 1
    - This applies to all PMAs in the PHY

  - The PCS receiver includes full 32 lane reorder and deskew block so that
    - Any PMA output lane can connect to any PMA input lane
    - There can be non-zero skew between the 32 lanes (same skew limits as CL120)

# Flow lane Muxing (Tx) : Example illustrating proposed change



**Baseline proposal from July 2022**

**Example bitmuxing that meets new proposal**

# PMA

- PMA functions as defined in CL120, with latest 802.3ck updates for 100G/lane
  - Bit-multiplexing (4:1) [if PMA Service interface below the PCS is 32 lanes]
  - Modulation (PAM4)
  - AUI Physical lane instantiation (8 lane)
  - Signaling lane rate (106.25Gb/s)
  - Coding (Gray, precoding)
  - Clock and data recovery
  - Loopbacks
  - Test patterns


- Per lane AUI specifications from 802.3ck

## Outline

- Introduction
- PCS/FEC/PMA Baseline proposal
- **Conclusions**

# Conclusions

- This Baseline: 800GbE PCS, FEC and PMA for 8 x 100G PMDs and 8 x 100G AUIs

- Supports all adopted 802.3df copper and optical PMDs baselines using 100G/lane

- Fits into an overall 800GbE Logic Architecture
  - Does not constrain future PCS/FEC/PMA schemes using 200G/lane AUIs, PMDs and/or Coherent PMDs

- 1.6TbE PCS/FEC can be chosen independently of 800GbE
  - Decisions made in this baseline will not restrict options / choices for 1.6TbE

**Thanks**