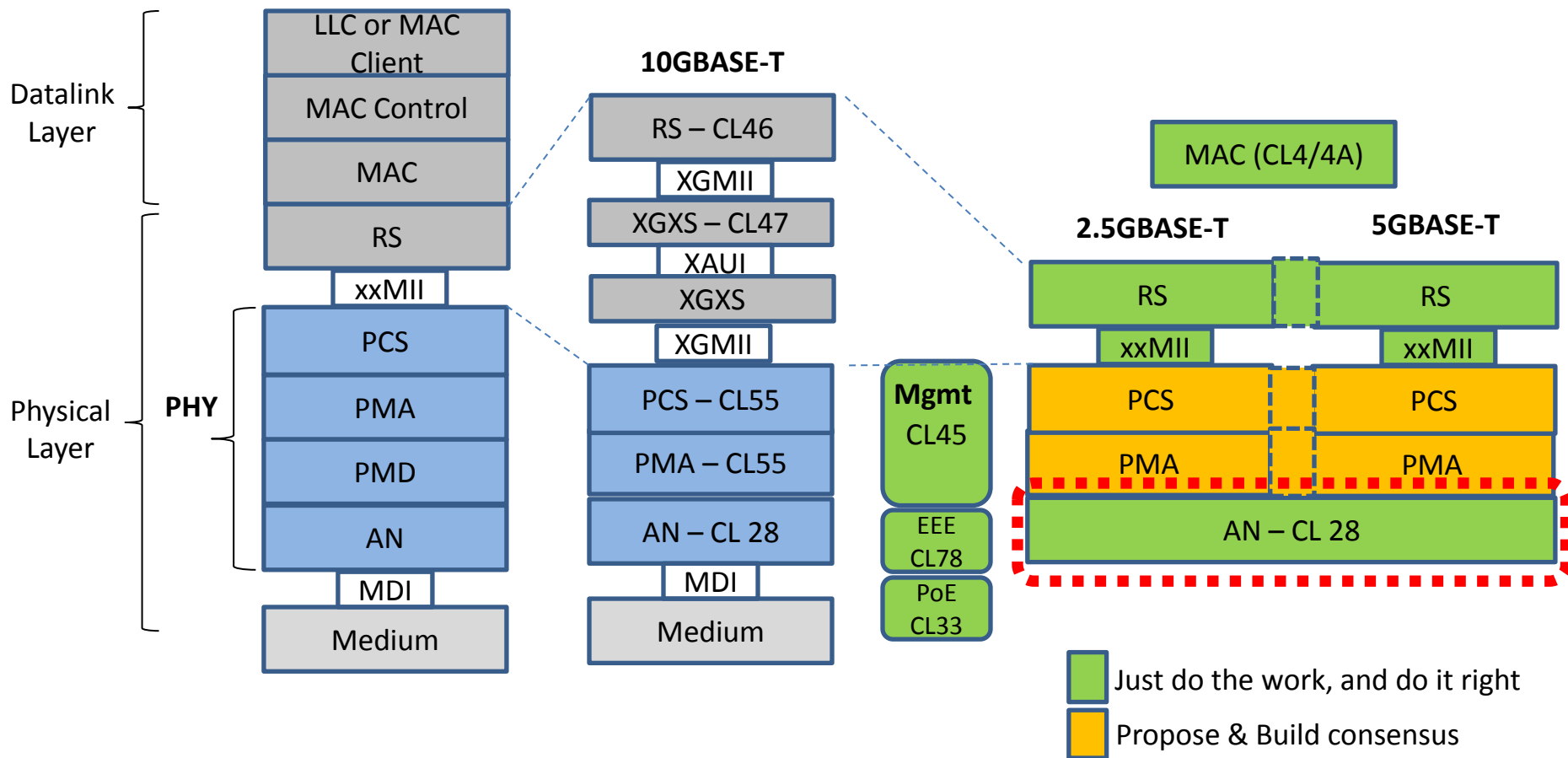# 802.3bz Layers – Auto-negotiation Proposal

(Presented on April 21st 2015 ad hoc call, Revised from April 14th, 2015,
Option 2 as the proposal based on .3bq direction.
Decouples .3bz from .3bq but still coordinated)

Yong Kim (ybkim at broadcom com), presenting
Tooraj Esmailian (ToorajE at broadcom com),
Brad Booth (BrBooth at microsoft com as the 10GBASE-T chair),
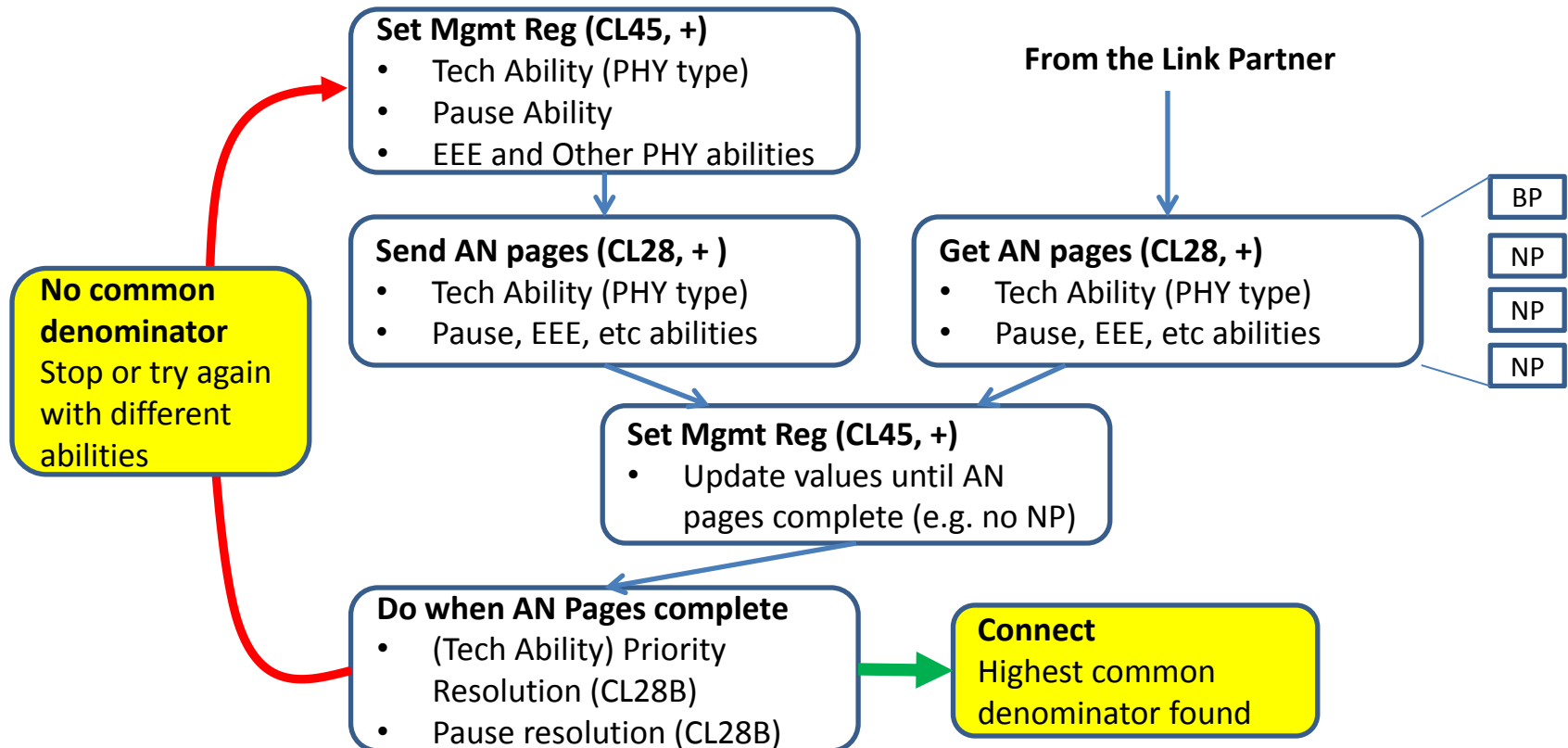
# Auto-Negotiation

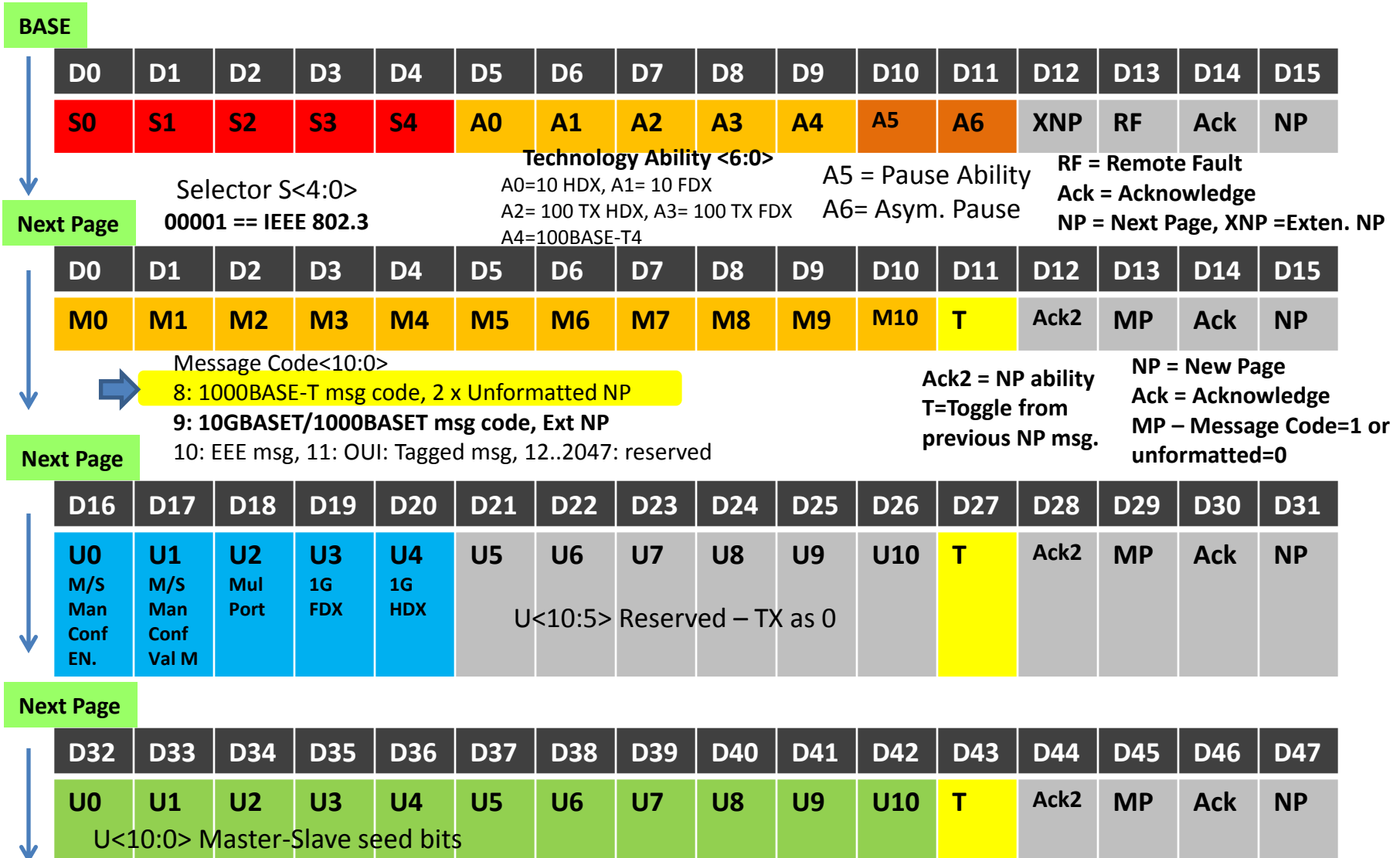## 2.5G and 5G BASE-T Layering considerations

# AUTO-NEGOTIATION BACKGROUND

# Auto-Negotiation System

- **Refer to Auto-Negotiation (AN) Overview and Read CL28 +**
  - http://www.ieee802.org/3/by/public/Mar15/booth_3by_01_0315.pdf
- AN is an <u>open loop advertisement</u> – not a stateful protocol, just "Ack"s.
  - **Qualitative** description below (**Mgmt** – **CL45 is optional**, but info required for AN resides locally regardless.

**Set Mgmt Reg (CL45, +)**
- Tech Ability (PHY type)
- Pause Ability
- EEE and Other PHY abilities

**From the Link Partner**

**No common denominator**
Stop or try again with different abilities

**Send AN pages (CL28, + )**
- Tech Ability (PHY type)
- Pause, EEE, etc abilities

**Get AN pages (CL28, +)**
- Tech Ability (PHY type)
- Pause, EEE, etc abilities

BP

NP

NP

NP

**Set Mgmt Reg (CL45, +)**
- Update values until AN pages complete (e.g. no NP)

**Do when AN Pages complete**
- (Tech Ability) Priority Resolution (CL28B)
- Pause resolution (CL28B)

**Connect**
Highest common denominator found

# Auto-Negotiation (CL28) <u>Review</u> – 1G (CL40.5)

**BASE**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| S0 | S1 | S2 | S3 | S4 | A0 | A1 | A2 | A3 | A4 | A5 | A6 | XNP | RF | Ack | NP |

Selector S<4:0>
00001 == IEEE 802.3

Technology Ability <6:0>
A0=10 HDX, A1= 10 FDX
A2= 100 TX HDX, A3= 100 TX FDX
A4=100BASE-T4

A5 = Pause Ability
A6= Asym. Pause

RF = Remote Fault
Ack = Acknowledge
NP = Next Page, XNP =Exten. NP

**Next Page**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| M0 | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | T | Ack2 | MP | Ack | NP |

Message Code<10:0>
8: 1000BASE-T msg code, 2 x Unformatted NP
**9: 10GBASET/1000BASET msg code, Ext NP**
10: EEE msg, 11: OUI: Tagged msg, 12..2047: reserved

Ack2 = NP ability
T=Toggle from
previous NP msg.

NP = New Page
Ack = Acknowledge
MP – Message Code=1 or
unformatted=0

**Next Page**

| D16 | D17 | D18 | D19 | D20 | D21 | D22 | D23 | D24 | D25 | D26 | D27 | D28 | D29 | D30 | D31 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U0<br>M/S<br>Man<br>Conf<br>EN. | U1<br>M/S<br>Man<br>Conf<br>Val M | U2<br>Mul<br>Port | U3<br>1G<br>FDX | U4<br>1G<br>HDX | U5 | U6 | U7 | U8 | U9 | U10 | T | Ack2 | MP | Ack | NP |

U<10:5> Reserved – TX as 0

**Next Page**

| D32 | D33 | D34 | D35 | D36 | D37 | D38 | D39 | D40 | D41 | D42 | D43 | D44 | D45 | D46 | D47 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U0 | U1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | U9 | U10 | T | Ack2 | MP | Ack | NP |

U<10:0> Master-Slave seed bits

# Auto-Negotiation (CL28) <u>Review</u> – EEE
## CL40.5 (1G) & CL45.2.7.13 (Mgmt.EEE)

For 1G **EEE** (CL40), do
1) This (1G auto-neg on the previous slide)
2) And then below (EEE)

Note: Msg code 9 based EEE is NOT referenced in CL40 (1G)
Maintenance item?

**Next Page**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|------|-----|-----|-----|
| M0 | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | T | Ack2 | MP | Ack | NP |

Message Code<10:0>
8: 1000BASE-T msg code, 2 x Unformatted NP
**9: 10GBASET/1000BASET msg code, Ext NP**
10: EEE msg, 11: OUI: Tagged msg, 12..2047: reserved

Ack2 = NP ability
T=Toggle from previous NP msg.

NP = New Page
Ack = Acknowledge
MP – Message Code=1 or unformatted=0

**Next Page**

| D16 | D17 | D18 | D19 | D20 | D21 | D22 | D23 | D24 | D25 | D26 | D27 | D28 | D29 | D30 | D31 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|-----|-----|-----|
| U0 | U1 100TX EEE | U2 1G EEE | U3 10G EEE | U4 1G KX EEE | U5 10G KX4 EEE | U6 10G KR EEE | U7 | U8 | U9 | U10 | T | Ack2 | MP | Ack | NP |

Note: U<10:0> is specified in 45.2.7.13 reference to copy bit by bit -- 28C.12. Bits 15:0 of register 7.60 – EEE Adv Register in CL45, Table 190

# Auto-Negotiation (CL28) <u>Review</u> – 10G (CL55.6)

**BASE**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| S0 | S1 | S2 | S3 | S4 | A0 | A1 | A2 | A3 | A4 | A5 | A6 | XNP | RF | Ack | NP |

Selector S<4:0>
**00001 == IEEE 802.3**

**Technology Ability <6:0>**
A0=10 HDX, A1= 10 FDX
A2= 100 TX HDX, A3= 100 TX FDX
A4=100BASE-T4

A5 = Pause Ability
A6= Asym. Pause

**RF = Remote Fault**
**Ack = Acknowledge**
**NP = Next Page, XNP =Exten. NP**

**Next Page**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| M0 | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | T | Ack2 | MP | Ack | NP |

Message Code<10:0>
8: 1000BASE-T msg code, 2 x Unformatted NP
**9: 10GBASET/1000BASET msg code, Ext NP**
10: EEE msg, 11: OUI: Tagged msg, 12..2047: reserved

Ack2 = NP ability
T=Toggle from
previous NP msg.

NP = New Page
Ack = Acknowledge
MP – Message Code=1 or
unformatted=0

**Ext NP W1**

| D16 | D17 | D18 | D19 | D20 | D21 | D22 | D23 | D24 | D25 | D26 | D27 | D28 | D29 | D30 | D31 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U0 | U1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | U9 | U10 | U11 10G M/S Man Conf EN. | U12 10G M/S Conf Val M | U13 (10G) Mul Port | U14 1G FDX | U15 1G HDX |

U<10:0> Master-Slave seed bits

**Ext NP W2**

| D32 | D33 | D34 | D35 | D36 | D37 | D38 | D39 | D40 | D41 | D42 | D43 | D44 | D45 | D46 | D47 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U16 10G | U17 (10G) LD Lp Time | U18 (10G) Short Reach | U19 (10G) Fast retrn | U20 (10G) LD train Rst rq | U21 | U22 100T X EEE | U23 1G EEE | U24 10G EEE | U25 | U26 | U27 | U28 | U29 | U30 | U31 |

# Auto-Negotiation (CL28) – .3bq D2.0

**BASE**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| S0 | S1 | S2 | S3 | S4 | A0 | A1 | A2 | A3 | A4 | A5 | A6 | XNP | RF | Ack | NP |

Selector S<4:0>
00001 == IEEE 802.3

Technology Ability <6:0>
A0=10 HDX, A1= 10 FDX
A2= 100 TX, HDX A3= 100 TX FDX
A4=100BASE-T4

A5 = Pause Ability
A6= Asym. Pause

RF = Remote Fault
Ack = Acknowledge
NP = Next Page, XNP =Exten. NP

**Next Page**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|------|-----|-----|-----|
| M0 | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | T | Ack2 | MP | Ack | NP |

Message Code<10:0>
8: 1000BASE-T msg code, 2 x Unformatted NP
➡ 9: ~~10GBASET/1000BASET~~ xGBASE-T msg code, Ext NP
10: EEE msg, 11: OUI: Tagged msg, 12..2047: reserved

Ack2 = NP ability
T=Toggle from previous NP msg.

NP = New Page
Ack = Acknowledge
MP – Message Code=1 or unformatted=0

**Ext NP W1**

| D16 | D17 | D18 | D19 | D20 | D21 | D22 | D23 | D24 | D25 | D26 | D27 | D28 | D29 | D30 | D31 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U0 | U1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | U9 | U10 | U11 ~~(10G)~~ M/S Man Conf EN. | U12 ~~(10G)~~ M/S Conf Val Mstr. | U13 ~~(10G)~~ Mul Port | U14 1G FDX | U15 1G HDX |

U<10:0> Master-Slave seed bits

**Ext NP W2**

| D32 | D33 | D34 | D35 | D36 | D37 | D38 | D39 | D40 | D41 | D42 | D43 | D44 | D45 | D46 | D47 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U16 10G | U17 (10G) LD Lp Time | U18 (10G) Short Reach | U19 (10G) Fast retrn | U20 (10G) LD train Rst rq | U21 40G | U22 100T X EEE | U23 1G EEE | U24 10G EEE | U25 40G EEE | U26 | U27 | U28 | U29 | U30 | U31 |

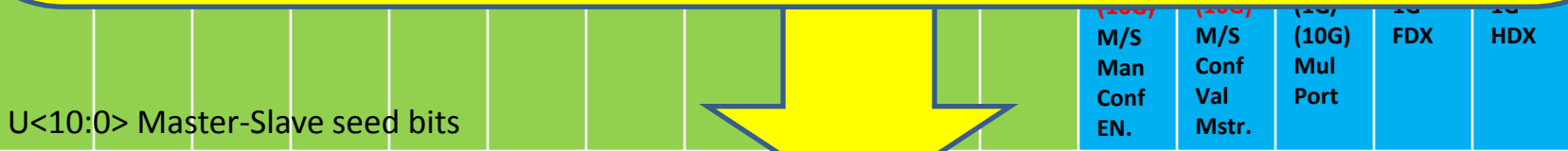# Auto-Negotiation (CL28) – .3bq ⬅ 25G

**BASE**

**Bits that needs to go in here for 25G and 40G.**

- **40G Fast Re-train (need separate bit from 10G)**
- **40G Repetitive Training Pattern mode (jul14/souvignier_3bq_01_0714.pdf)**

- **25G**
- **25G EEE**
- **25G Fast Re-train**
- **25G Repetitive Training Pattern Mode**

**6 Spare bts at present.**
**6 bits needed at present (no more bits left, if we find something later), and**
**no bits left for 2.5G and 5G.**
**Note: .3 bq reuse of 10G Master-Slave related fields (11 + 1 + 1 + 1)**

| | | | | | | | | | | | (10G) M/S Man Conf EN. | (10G) M/S Conf Val Mstr. | (1G) (10G) Mul Port | 1G FDX | 1G HDX |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

U<10:0> Master-Slave seed bits

**Ext NP W2**

| D32 | D33 | D34 | D35 | D36 | D37 | D38 | D39 | D40 | D41 | D42 | D43 | D44 | D45 | D46 | D47 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| U16 10G | U17 (10G) LD Lp Time | U18 (10G) Short Reach | U19 (10G) Fast retrn | U20 (10G) LD train Rst rq | U21 40G | U22 100T X EEE | U23 1G EEE | U24 10G EEE | U25 40G EEE | U26 | U27 | U28 | U29 | U30 | U31 |

# CL28 - Rules and Observations

- Reminder: Do not redefine bits.  Reusing bits without any functional changes may be ok (but need to be careful).

- Not very obvious and clean, especially inclusive  of 1G and 10G coding.
  - Attempted "cover all modern PHYs" in the message code 9 to serve 100TX/1G/10G & EEE versions is not working well for us now.

- Legacy replicated info (what's allowed in AN)
  - EEE message bits (EEE capabilities for 100TX, 1G, 10G BASE-T)  and Message code 9 (10G) already replicate CL45 mgmt register info on AN.   100TX, 1G FDX, and 1G HDX also replicated in M Code 8 and 9.

- 2.5G and 5G needs
  - 8 bits -- 4 bits for each <Speed, EEE, Fast Re-train, Repetitive Train?>
  - Master/Slave (14 bits = 11 + 1 + 1 + 1) -- could be common w/ 1G/10G/25G/40G.
  - Needs 22 bits of which 14 M/S bits *may* be reused w/ other BASE-T speeds.
    - Consequence of sharing of M/S bits -- Mixed speed multiport device where not all ports have the same (e.g. speed) capability – not likely as a product but possible and allowed by standard. E.g. what does a multi-port PHY that support 4 x 2.5 G and 1 x 2.5G/5G report?  Std is not clear.
    - Suggest NOT to dwell on this point.  Offered as an information for completeness.

# CL28 .3bz - So what are the options?

- 🚫 Not an option – fit into XNP msg code 9 (10G) in flight (.3bq)
- 🚫 Option 1 – go back to 1G method (msg code 8 & 10 (EEE)
  - BP + NP + NP + NP and add more NP (new) + NP (new) for 2.5G/5G  ….
- Option 2 – Define a new 2.5G/5G Extended NP
  - BP + XNP(1)-msg code 12 (new) + XNP (2) + XNP(3)
  - 22 bits out of 32 bits used.
  - Reflects 802.3bq D2.0 (current as of this PDF).
- Option 3 – Define 2.5G/5G/25G/40G Extended NP
  - BP + XNP(1)-msg code 12 (new) + XNP (2) + XNP(3)
  - 30 bits out of 32 bits used.  (or 16 bit out 32 used, if MC9 M/S re-used).
  - Coordinate w/ .3bq + 25G project
- Option 4 – Reuse 10G (MC 8) for 2.5G/5G, and "ask" 25G/40G to go to a new Extended NP (and let it replicate 10G bits perhaps).
  - 8 bits, Speed, EEE, Fast Re-train, Rep Train )*2, needed out of 8 bits available for 2.5G and 5G.  No spares.

# CL 28 .3bz Options and Consequences

| O | Description | 1G | 2.5G/5G | 10G | 25G | 40G |
|---|---|---|---|---|---|---|
| 2 | 2.5G & 5G gets its own new page | MC8 & (9 or 10) | MC12 | MC9 | MC9 | MC9 |
| | ← Optimizes well! | | SB =3(17) | MC9 - no Spare Bits, post .3bq work | | |
| 3 | 2.5G/5G/25G/40G to go to a new page | MC8 & (9 or 10) | MC12 | MC9 | MC12 | MC12 |
| | ← Not the .3bq direction | | SB =2(16) | SB = 8 | SB = 2(16) | |
| 4 | "Ask" 25G/40G to go to a new page. 2.5G/5G uses MC9 | MC8 & (9 or 10) | MC9 | MC9 | MC12 | MC12 |
| | ← Not the .3bq direction ← Not recommended for .3bz | | MC9 – No Spare Bits anticipated (4 x 2) | | SB=12(26) | |

Note: SB = Spare bits, (nn) denotes if MC9 Master/Slave related fields are re-used and common across 2.5G/5G/10G/25G/40G, regardless of MC12 use.

# 802.3BZ AUTO-NEGOTIATION PROPOSAL (based on the "Option 2")

**in**

**http://www.ieee802.org/3/NGEBASET/public/archadhoc/Kim_AutoNegotiation_v2_2015_4_14c.pdf**

Based on 802.3bq (early) indication to stay in Message Code 9.

Acting on "IF .3bq does not want to move, THEN …

# CL 28 .3bz AN Objectives &  A Proposal

**Objectives**

- 802.3bq and 802.3bz to be coordinated, i.e. 2.5G/5G/25G/40G.

- Got (early) .3bq Feedback – 802.3bq to stay in MC9, so acting on that preference,

- Define 2.5G/5G BASE-T PHY related AN bits in new MC.

**Proposal**

- Define a new message code 12 for 2.5G/5G ("Option 2")

- Design such a way that modern RJ-45 MDI PHYs only need to <u>support Base Page plus XNP MC12 (new) – optimize and help reduce AN duration.</u>
  - Does <u>NOT relieve PHY's support of</u> other message codes  (true, and has been true).

- Design such a way to recognize the following modern optimizations.
  - 10M/100M/1G to extend to support 10M/100M/1G/<span style="color:red">2.5G</span>
  - 1G/10G to extend (down) support for 1G/<span style="color:red">2.5G</span>/<span style="color:red">5G</span>/10G
  - Superset of the above two ranges – <u>does it fit?  - YES!</u> w/ spare bits.
  - No ability assignment for 1G HDX, no need to replicate.

- Master/Slave related fields are replicated (as done in MC8 and MC9).

- 3 spare bits, or 1 bit, left, if .3bz adopts "repeat train capability".

# Auto-Negotiation (CL28) – .3bz Proposal

**BASE**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| S0 | S1 | S2 | S3 | S4 | A0 | A1 | A2 | A3 | A4 | A5 | A6 | XNP | RF | Ack | NP |

**Technology Ability <6:0>**
A0=10 HDX, A1= 10 FDX
A2= 100 TX, HDX A3= 100 TX FDX
A4=100BASE-T4

Selector S<4:0>
00001 == IEEE 802.3

A5 = Pause Ability
A6= Asym. Pause

RF = Remote Fault
Ack = Acknowledge
NP = Next Page, XNP =Exten. NP

**Next Page**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| M0 | M1 | M2 | M3 | M4 | M5 | M6 | M7 | M8 | M9 | M10 | T | Ack2 | MP | Ack | NP |

Message Code<10:0>
8: 9: 10GBASET/1000BASET msg code, Ext NP
10: EEE msg, 11: OUI: Tagged msg,
12: 100M/1G/2.5G/5G/10G        13 ..2047: reserved

Ack2 = NP ability
T=Toggle from
previous NP msg.

NP = New Page
Ack = Acknowledge
MP – Message Code=1 or
unformatted=0

**Ext NP W1**

| D16 | D17 | D18 | D19 | D20 | D21 | D22 | D23 | D24 | D25 | D26 | D27 | D28 | D29 | D30 | D31 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U0 | U1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | U9 | U10 | U11 M/S Man Conf EN. | U12 M/S Conf Val Mstr. | U13 Mul Port | U14 1G FDX | U15 1G HDX |

U<10:0> Master-Slave seed bits

**Ext NP W2**

| D32 | D33 | D34 | D35 | D36 | D37 | D38 | D39 | D40 | D41 | D42 | D43 | D44 | D45 | D46 | D47 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U16 10G | U17 (10G) LD Lp Time | U18 (10G) Short Reach | U19 (10G) Fast retrn | U20 (10G) LD train Rst rq | U21 2.5G | U22 100T X EEE | U23 1G EEE | U24 10G EEE | U25 2.5G EEE | U26 2.5G Fast retrn | U27 5G | U28 5G EEE | U29 5G Fast retrn | U30 | U31 |

# Auto-Negotiation (CL28 + CL55.6.1 & .3bq.6.1)

- CL28 complete (no functional changes required, just revisions)
  - Table 28-9 – Timer min/max value
    - Link_Fail_inhibit_timer (10G**/25G?/40G**) – min 2000, max 2250 msec.
    - **Add** Link_Fail_inhibit_timer(s) for **2.5G/5G** – min **TBD**, max **TBD** msec.
- Annex 28B.3 – Priority resolution
  - Insert 2.5G and 5G above 1G and below 10G.
- Annex 28C – Next Page Msg Code field definitions
  - Table 28-C-1 (message code 9 (Ext NP, xGBASE-T) – code field value entry (or entries) for 2.5G and 5G, and corresponding message code text.
- Annex 28D – Description of extensions to CL 28 and assc. annexes.
  - 28D, insert as 28D.9(?), after 40G and 25GBASE-T
    - Auto-neg mandatory for 2.5G and 5GBASE-T, extended NP support, use of MASTER and SLAVE PHY operation, support of the priority resolution table (Annex 28B.3), and asymmetric pause (Annex 28B.2 "A6"), etc.
- And reflect the above changes to the PICS (28.5)

# Annex 28B

- **28B.3 Priority resolution**

Modify the priorities as:

   a) 40GBASE-T full duplex
   b) 25GBASE-T full duplex
   c) 10GBASE-T full duplex
   d) 5GBASE-T full duplex
   e) 2.5GBASE-T full duplex
   f) 1000BASE-T full duplex
   g) 1000BASE-T
   h) 100BASE-T2 full duplex
   i) 100BASE-TX full duplex
   j) 100BASE-T2
   k) 100BASE-T4
   l) 100BASE-TX
   m) 10BASE-T full
   n) 10BASE-T

# Annex 28D

**Add at the end of Annex 28D and replace <xx>, <yy>, <nn> as appropriate (provided as e.g.).**

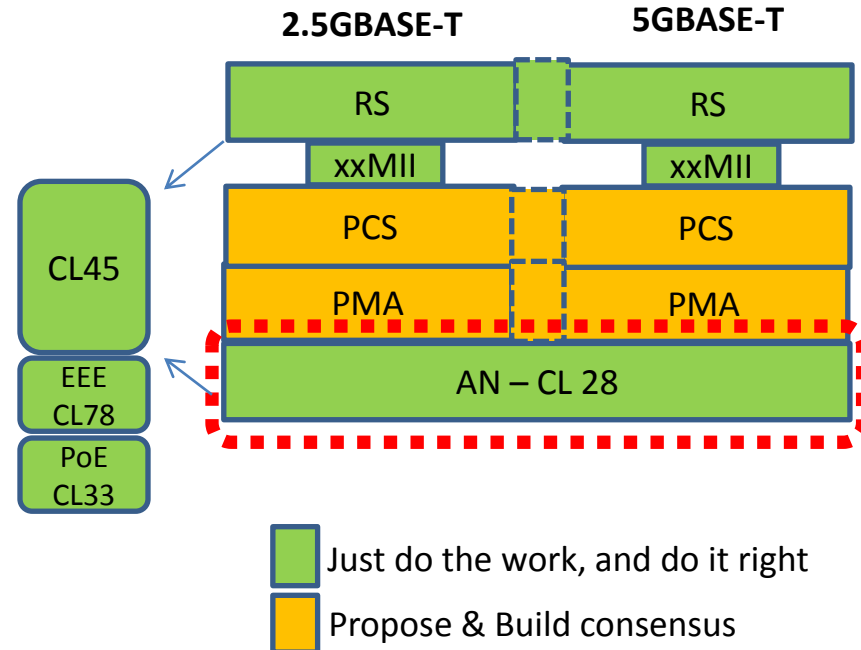**28D.xx  Extensions required for Clause yy (2.5GBASE-T and 5GBASE-T)**

Clause yy (2.5GBASE-T and 5GBASE-T) makes special use of Auto-Negotiation and requires additional MDIO registers. This use is summarized below.  Details are provided in <yy.nn>.

a)     Auto-Negotiation is mandatory for 2.5GBASE-T and 5GBASE-T.

b)     Extended Next Page support is mandatory for 2.5GBASE-T and 5GBASE-T

c)     2.5GBASE-T and 5GBASE-T requires an exchange of an Extended Next Page message.

d)     2.5GBASE-T and 5GBASE-T parameters are configured based on information provided by the exchange of an Extended Next Page message.

e)     2.5GBASE-T and 5GBASE-T uses MASTER and SLAVE to define PHY operations and to facilitate the timing of transmit and receive operations. Auto-Negotiation is used to provide information used to configure MASTER-SLAVE status.

f)     2.5GBASE-T and 5GBASE-T transmits and receives an Extended Next Page for exchange of information related to MASTER-SLAVE operation. The information is specified in 45.2.7.

g)     2.5GBASE-T and 5GBASE-T adds 2.5GBASE-T and 5GBASE-T full duplex capabilities to the priority resolution table (see 28B.3).

h)     2.5GBASE-T is defined as a valid value for "x" in 28.3.1 (e.g., link_status_2.5GigT.) 2.5GigT represents that the 2.5GBASE-T PMA is the signal source.

i)     5GBASE-T is defined as a valid value for "x" in 28.3.1 (e.g., link_status_5GigT.) 5GigT represents that the 5GBASE-T PMA is the signal source.

j)     2.5GBASE-T and 5GBASE-T supports Asymmetric Pause as defined in Annex 28B.

# Summary

- CL 28 auto-negotiation changes are [still] straight forward.
  - Defined MC12 and assign extended next page assignments (new)
  - Supports all modern PHYs between 10M ~ 10G, with some spare bits, and NO need to use BOTH MC9 and MC12.

- Next Steps
  - Consider this proposal for adoption in 802.3bz TF.

  Note: Feedbacks, objections, support,…, all welcome toward May Interim.

# Thank you!

# And for something completely different…..

Backup slides on new selector consideration

# How about the new base page?

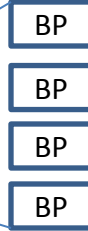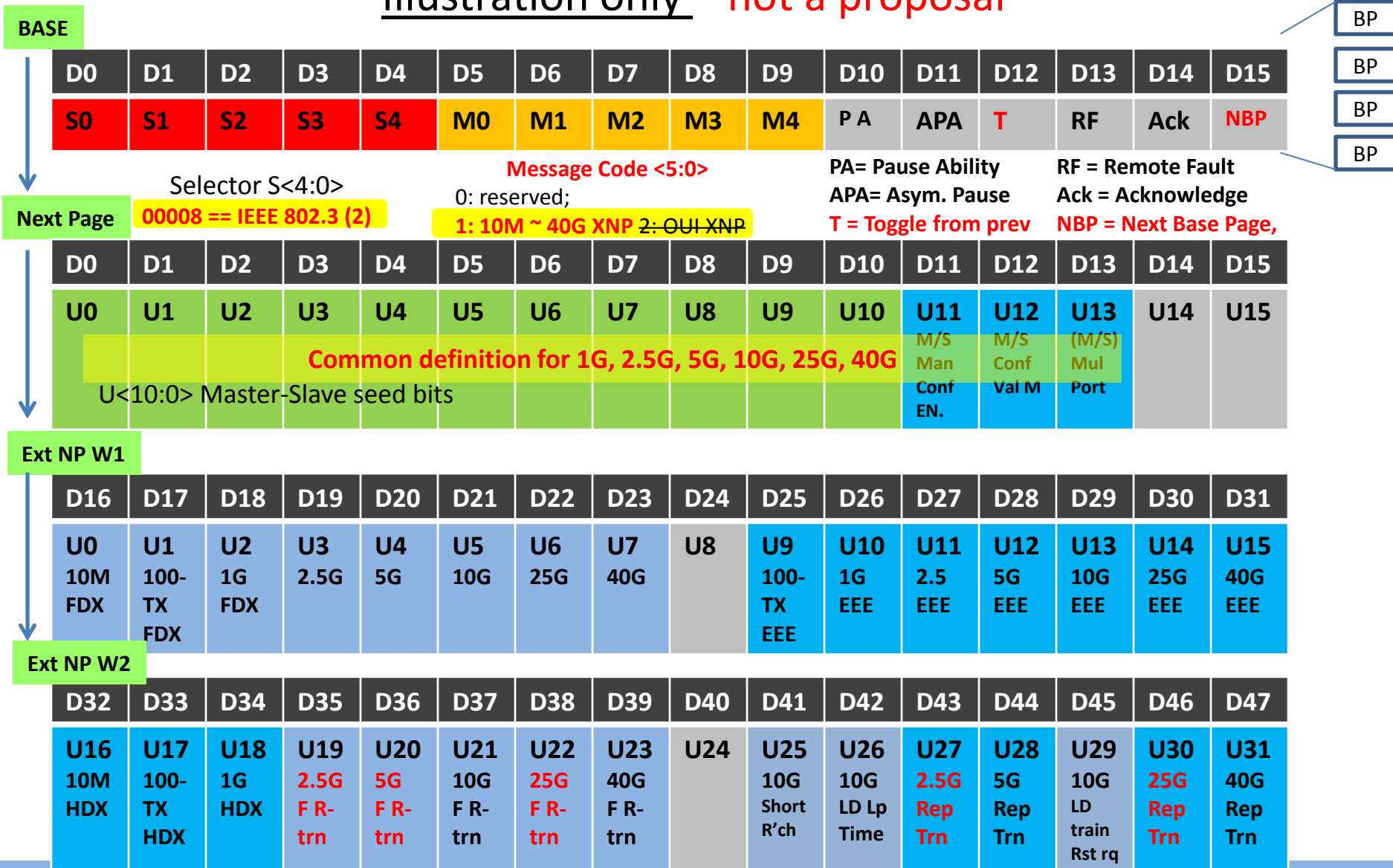| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| S0 | S1 | S2 | S3 | S4 | U0 | U1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | U9 | U10 |

Selector S<4:0>

1: IEEE 802.3; .. 6,7:R, 8== IEEE 802.3 (2)

- Selector for different network attachments that may share RJ45.
  - The value of sharing RJ45 is diminishing (TR, Firewire, little else in the pipeline).
  - Take few more values for 802.3 and define "cleaner" CL28 AN for modern PHYs (2015 and beyond).
  - **Consequence – likely a new AN Clause (or some manageable but substantial material to be added).**

- Observations
  - For each selector base page, we get 11 bits. -- not enough for multiple modern PHYs in BP.
    - Need 5 shared bits of [NP, ACK, RF, Pause, Asym Pause]
    - Need a field for message codes, e.g. OUI, future, etc (could be less than current 11 bits)
  - [Shared] Magnetic compatibility – Full range of 10M~40G AN [practically] irrelevant.
    - (Not encouraging this! **Scope**) Consider taking multiple selector fields around magnetic compatibility. -- some previous such optimization rendered not true w/ R&D.
  - Sample Extended NP format – 44 out of 48 bits used (next slide)
    - Need 30 bits = 2 for 10M; 3 each for 100M, 1G; 6 for 10G; 4 for 40G; 4 each for 2.5G/5G/25G.
    - Need Master/Slave bits 14 = Seed (11)+ Control (3)
    - Base Page conveys remote fault and pause abilities, plus AN related control bits.

# Auto-Negotiation (CL28)– New Selector
## Illustration only – not a proposal

**BASE**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| S0 | S1 | S2 | S3 | S4 | M0 | M1 | M2 | M3 | M4 | P A | APA | T | RF | Ack | NBP |

**Message Code <5:0>**

PA= Pause Ability
RF = Remote Fault

APA= Asym. Pause
Ack = Acknowledge

Selector S<4:0>

00008 == IEEE 802.3 (2)

0: reserved;

**1: 10M ~ 40G XNP** 2: OUI XNP

T = Toggle from prev
NBP = Next Base Page,

**Next Page**

| D0 | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 | D11 | D12 | D13 | D14 | D15 |
|----|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|
| U0 | U1 | U2 | U3 | U4 | U5 | U6 | U7 | U8 | U9 | U10 | U11 M/S Man Conf EN. | U12 M/S Conf Val M | U13 (M/S) Mul Port | U14 | U15 |

**Common definition for 1G, 2.5G, 5G, 10G, 25G, 40G**

U<10:0> Master-Slave seed bits

**Ext NP W1**

| D16 | D17 | D18 | D19 | D20 | D21 | D22 | D23 | D24 | D25 | D26 | D27 | D28 | D29 | D30 | D31 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U0 10M FDX | U1 100-TX FDX | U2 1G FDX | U3 2.5G | U4 5G | U5 10G | U6 25G | U7 40G | U8 | U9 100-TX EEE | U10 1G EEE | U11 2.5 EEE | U12 5G EEE | U13 10G EEE | U14 25G EEE | U15 40G EEE |

**Ext NP W2**

| D32 | D33 | D34 | D35 | D36 | D37 | D38 | D39 | D40 | D41 | D42 | D43 | D44 | D45 | D46 | D47 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| U16 10M HDX | U17 100-TX HDX | U18 1G HDX | U19 2.5G F R-trn | U20 5G F R-trn | U21 10G F R-trn | U22 25G F R-trn | U23 40G F R-trn | U24 | U25 10G Short R'ch | U26 10G LD Lp Time | U27 2.5G Rep Trn | U28 5G Rep Trn | U29 10G LD train Rst rq | U30 25G Rep Trn | U31 40G Rep Trn |

# Thank you!