



EEE Support for 100 Gb/s



Mark Gustlin, Hugh Barrass

IEEE P802.3bj

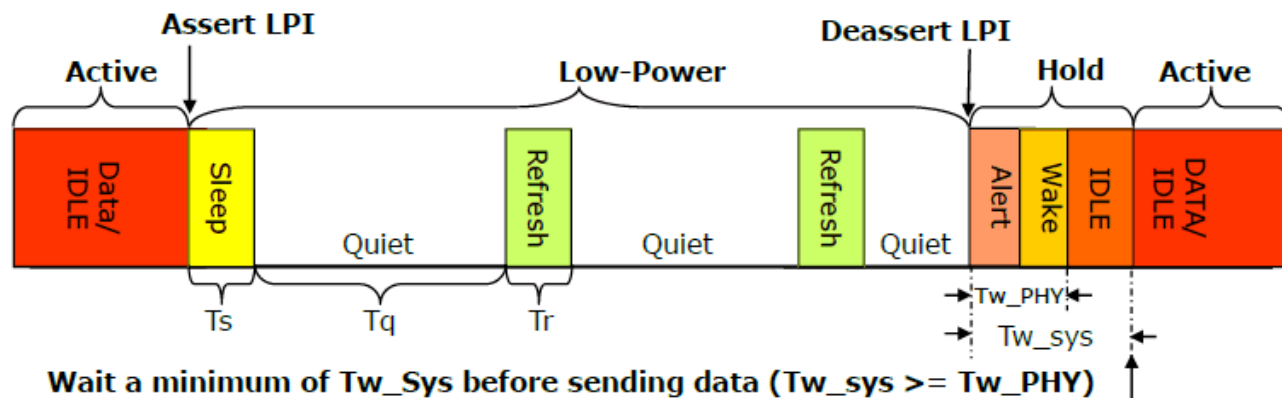
Atlanta November 2011

EEE for 100 Gb/s Overview

- This presentation will review the technical issues that need to be addressed in order to add EEE to 100 Gb/s copper interfaces
- The main issue has to do with the Alignment Marker lock, a proposal to address this concern is described
- Details and examples of other considerations for EEE at 100 Gb/s are also explored
- This presentation assumes re-using Low Power Idle, it does not investigate the complexities involved in any type of modular EEE

EEE Overview

LPI Overview



Wait a minimum of T_{w_Sys} before sending data ($T_{w_sys} \geq T_{w_PHY}$)

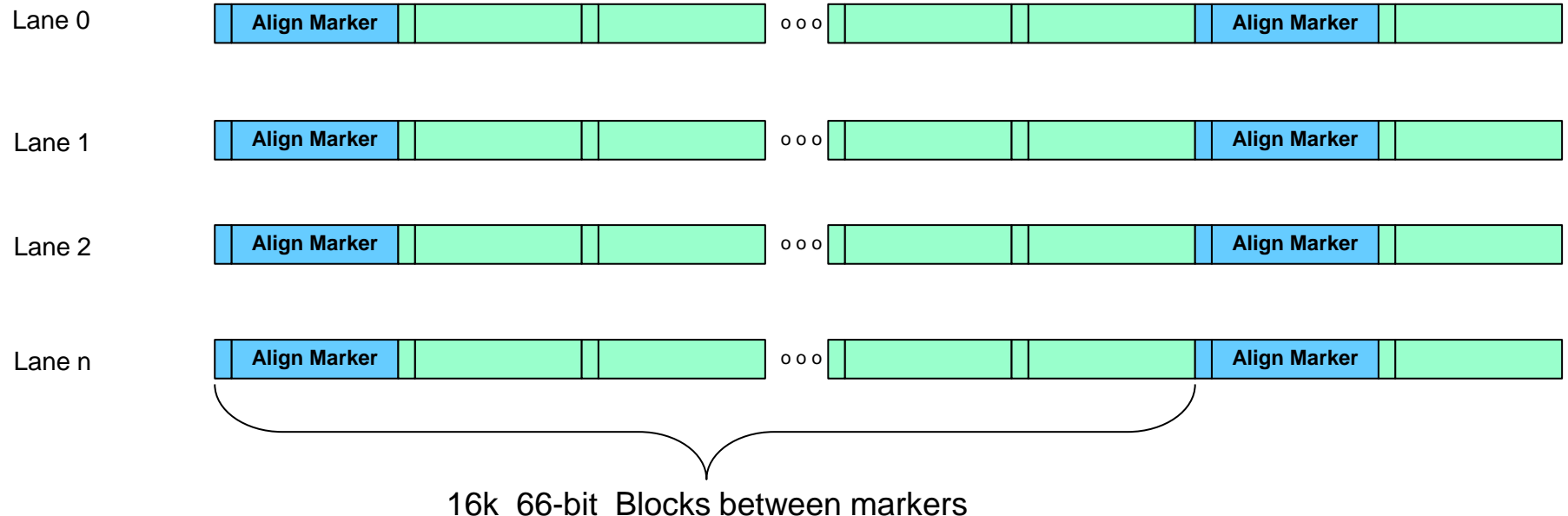
- LPI – PHY non-essential circuits shut down during idle periods
- During power-down, maintain coefficients and sync to allow rapid return to Active state
- Wake times for the respective backplane PHYs:
 - 1000BASE-KX: $T_{w_PHY(min)}$ = 11.25 usec
 - 10GBASE-KX4 $T_{w_PHY(min)}$ = 9.25 usec
 - 10GBASE-KR: $T_{w_PHY(min\ w/o\ FEC)}$ = 12.25 usec
 - 10GBASE-KR: $T_{w_PHY(min\ w/FEC)}$ = 14.25 usec

EEE Overview

- Wake time range is 9 to 14usec for existing EEE PHYs
- Note that the wake time does not scale down with speed even though data accumulates faster at higher interface speeds
- So for 100 Gb/s should we shoot for a wakeup time of < 5usec?
 - Note that in 5 μ s, 0.5Mb of data accumulates, per port
- Are there any concerns in the 100 Gb/s PCS that would prevent us from supporting a 5 usec or faster wakeup?
 - Alignment marker lock is >> 5 μ s, the next few slides look at this issue

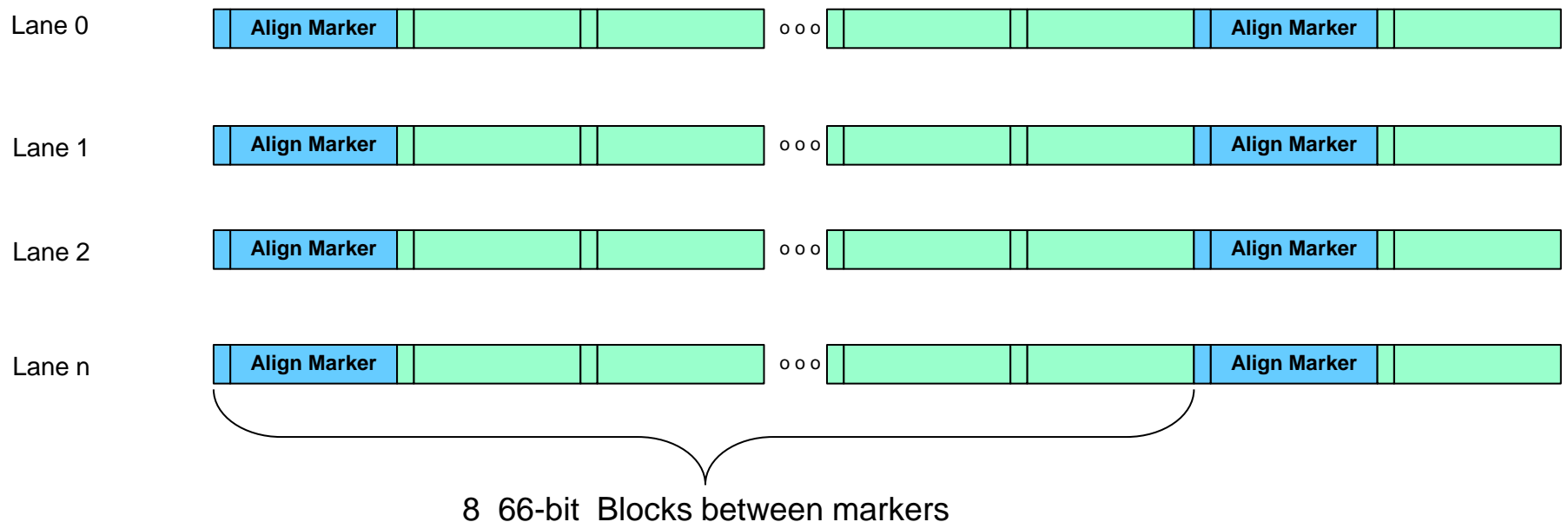
100/40GE Standard Alignment Marker Distance

- The alignment markers are widely spaced for 100 Gb/s and 40 Gb/s, 16k blocks apart on each PCS lane
- The alignment marker lock SM looks for two that match in a row before declaring lock and allowing alignment, so that is $16384 * 2 * 66 * 194ps = 419\mu s$ (for 100GE, $\frac{1}{2}$ that for 40GE)
- This would mean that startup would take $> 400\mu s$ today!
 - Ok for 802.3ba, not ok for a EEE interface



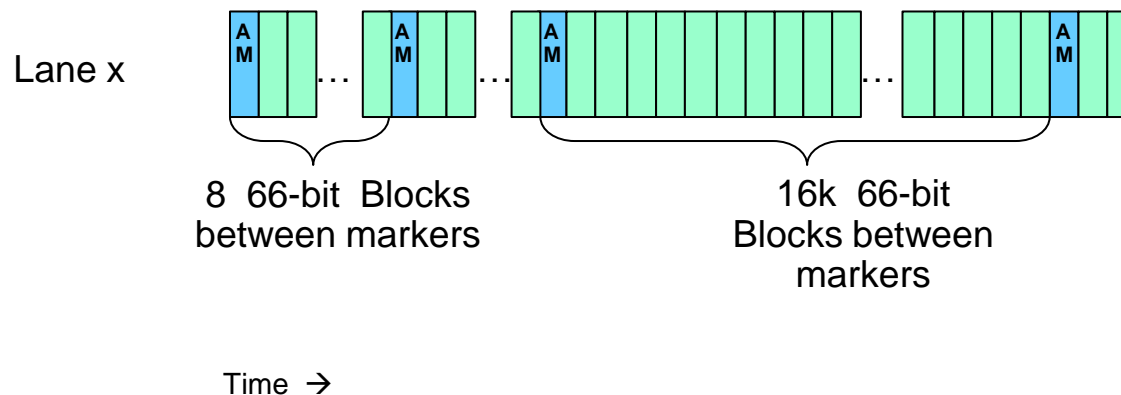
Start-up AM Distance Reduction

- When the lanes are starting up, reduce the distance between AMs temporarily
- The allowed minimum distance could be dependent on the total skew that is allowed (14 66-bit words for 100G at the RX PCS), or you could require that the receiver maintain last known alignment and only worry about the skew variation (44 bits total at the receiver, less than 1-66b word)
- Another option is to take advantage of the count down field which is located within the Rapid Alignment Marker (RAM) word, this allows the RAMs to be placed as close as we desire
- So let's say every 8 words there is an alignment marker until startup is finished, then revert to normal distance
- Alignment Marker lock would now take at least $8 * 2 * 66 * 194\text{ps} = 205\text{ns}$ (for 100GE), it can take longer with errors



Start-up Alignment Marker Distance

- When the lanes are starting back up, then Alignment Marker spacing is small, 8 for instance
- Then after a fixed time, the spacing goes back to 16k
- The transmitter has to signal to the far end when this change will take place, this can be done through repeated signals such as a count down field in the Alignment Marker, so that the receiver will know when the transition will take place even in the face of errors on the link
- The receiver should lock to the sync field in the AM in order to get 66b alignment, instead of taking 64 or more words to do that



How to Signal a Change in AM Distance?

- **Standard Alignment Marker format:**



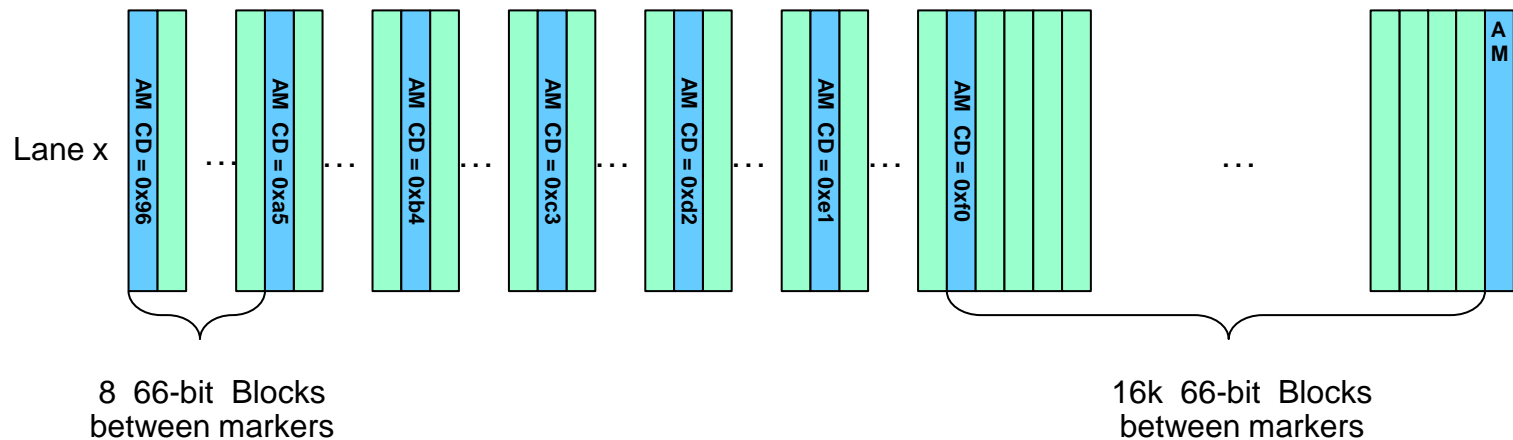
- **Proposed Rapid Alignment Marker format:**



- Add a count down field to the alignment marker, the receiver can use this to predict when the transmitter will switch the distance of the alignment markers
 - But how to know when the rx is listening, we might send AMs for a while before sending the last one? How many times are the fast AMs sent?
 - Send a fixed value in the CD field, and then Tx can start changing the countdown field with 16 AMs left before reverting to the standard distance
- Note that the AM is not scrambled, so anything we add should not negatively impact the baseline wander or clock content
- Assume that we want to be resilient in the face of up to 3 individual errors
- Convert the !BIP field to a Count Down field for the rapid AMs
- For the last 16 rapid AMs, set the !BIP field to 0x0F, 0x1E, 0x2D, 0x3C, 0x4B, 0x5A, 0x69, 0x78, 0x87, ...0xF0.
 - One nibble is the inverse of other to maintain DC balance and create at least one clock transition
 - How many do we send, if we send > 16 what do we do with this field? Field starts as a fixed value of 0xCC (DC balanced and average clock transitions)

Start-up Alignment Marker Distance

- The figure shows the short spacing of the AMs on link power up, the last 16 have a count down pattern, which allows the rx to tell when the AMs will go back to their normal spacing
- This occurs on all PCS lanes at the same time
- The receiver will look for multiple AMs at the short spacing, and look at the count down fields, compare and lock in the count down so that it can predict when the spacing returns to 16k words
- Once AM spacing returns to 16k, the !BIP field returns to its normal .ba function
- BIP values are still calculated the same way as is standard, just over shorter distances when the link is coming up, and the !BIP7 is not populated



Link Start-up

- When a device first powers up, it will power up sending 802.3ba AMs, 16k apart
- What if a fault causes the receiver to miss startup, and it first sees AMs on a power bring up from LPI?
 - State Machines already handle these kind of issues for 802.3az EEE, so the same thing would apply, the SM would have timeouts to handle the error cases, no need to do anything else.

BIP Coverage

- With the very short AM distance on LPI startup, does the BIP field stay as is? And just covers a lot less data?
- Proposal is yes, that BIP values are still sent as is, but instead of covering 131k bits or more, each BIP bit covers only 64 bits or so (8×8).
- Note that the receiver should not signal BIP errors until after the rapid AMs are locked to (prevents spurious BIPs from the initial LPI startup phase)

Data 'Randomness'

- There was a lot of work done during the 802.3ba project on baseline wander and clock content of the 100 Gb/s data stream and various sub 100 Gb/s lanes (10G CAUI lanes, 25G PMD lanes etc.) in order to ensure that the characteristics of a given serial data stream are good, see:

http://www.ieee802.org/3/ba/public/jan08/anslow_01_0108.pdf

http://www.ieee802.org/3/ba/public/nov08/anslow_06_1108.pdf

- When we send RAMs, does that impact the baseline wander or clock content negatively?
 - Simulations need to be run determine this
- There is some concern with when the lanes are being powered up, will the randomness of the data being sent at that time be sufficient to quickly train the receivers?
 - Simulations need to be run determine this
- The count down field is a new concept to this protocol, if you bit multiplex multiple streams together with an 802.3ba PMA, how will this impact the clock content (given that the count down field is not scrambled)?
 - Simulations need to be run determine this

Start-up timing

- How to know when to start sending and then stop sending rapid AMs?
- Proposal is to send Rapid AMs in TX_ALERT and TX_WAKE states, for 10GE that is 5usec total time
- Only Idles or LPIs are sent at this point

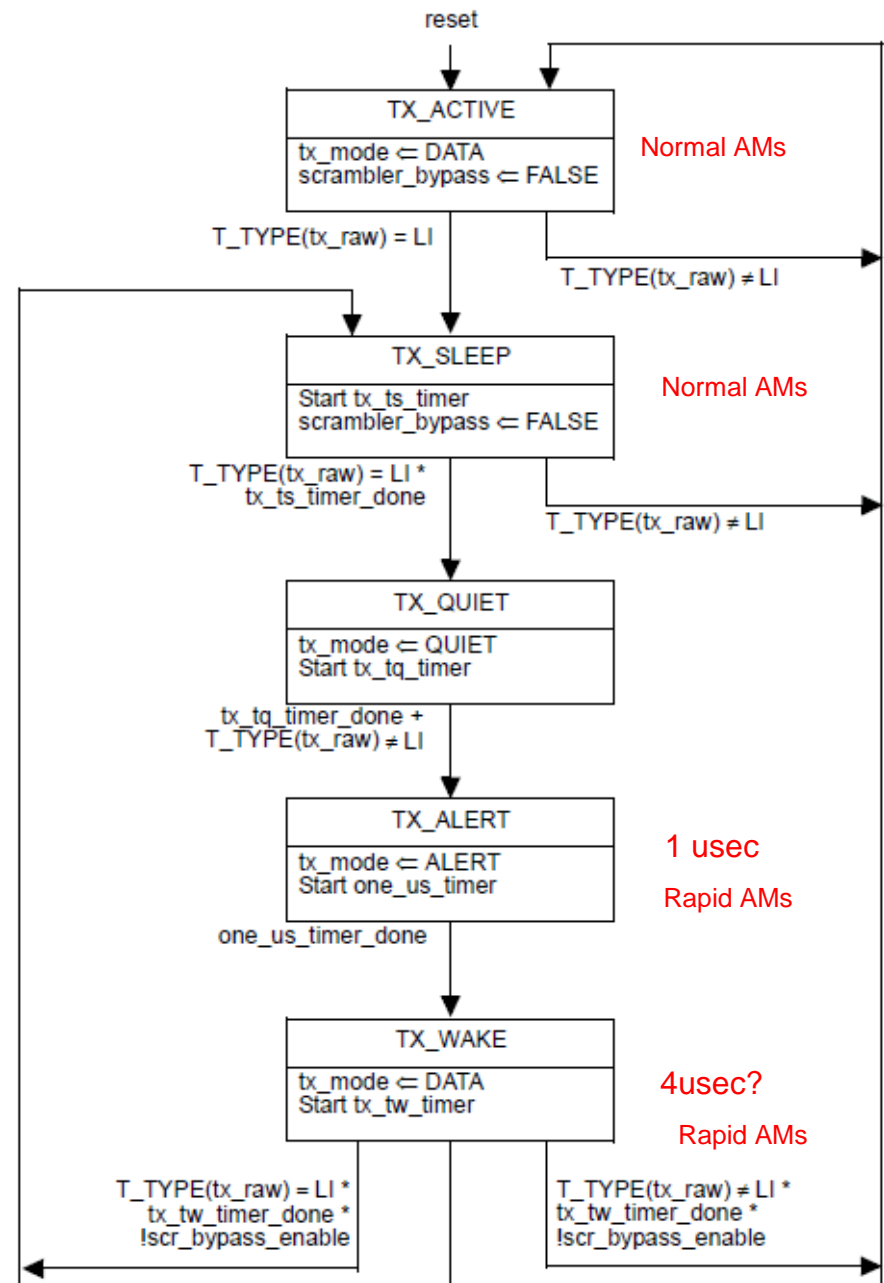
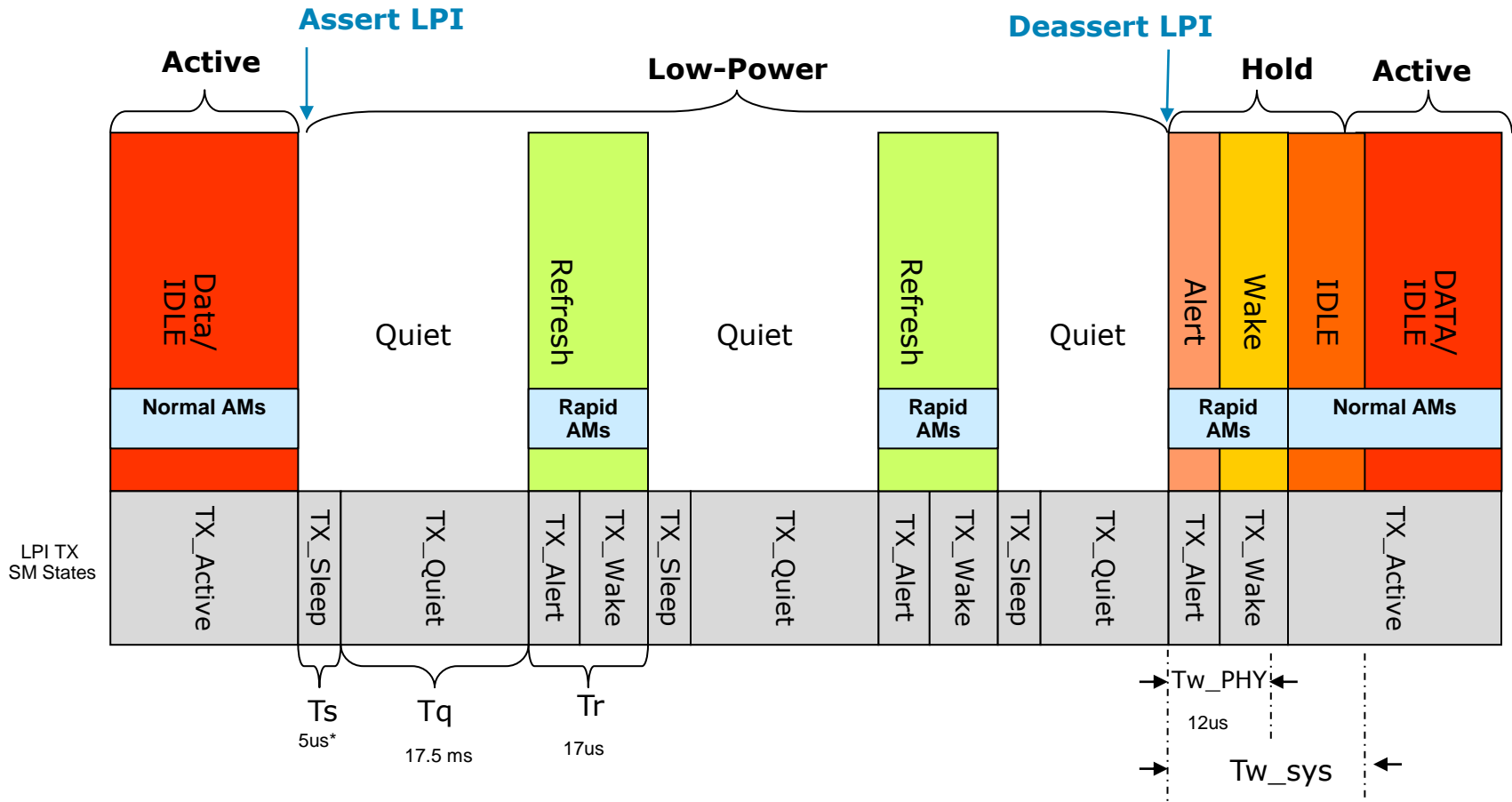


Figure 49-16– LPI Transmit state diagram

LPI with Rapid AMs



* Numbers are from 10GBSE-KR

Room for Rapid AMs?

- In 802.3ba the Alignment Markers are sent very infrequently, every 16k blocks on each PCS Lane
- This allows room for the AMs to be added into the data stream by deleting Idles periodically, just as is done for clock compensation
- If we send AMs rapidly, then we can still delete idles in order to send AMs?
 - Yes this works since rapid AMs are only sent on link re-start when transitioning out of LPI, so only LPI or Idle is being sent at that time
 - Proposal is that Either Idles or LPIs can be deleted to add in the AMs

Document Changes Required

If we add EEE to 100 Gb/s, what clauses have to be modified:

Clause 30: Management additions

Clause 45: MDIO register additions

Clause 74 (KR FEC): Minor changes if we include 40 Gb/s as a service to humanity

Clause 78: Add in overview of 100 Gb/s EEE, timing parameters etc

Clause 81 (RS/MII): Add in changes that are similar to those made in clause 46 for KR

Clause 82 (PCS): Add in changes that are similar to those made in clause 49 for KR, plus add in changes needed for the Rapid Alignment Marker support

Clause 83 (PMA): Add in signals that pass through the PMA (energy_detect for instance)

Clause 84 (40GBASE-KR4): EEE PMD changes if we include 40 Gb/s as a service to humanity

Clause 85 (40GBASE-CR4/100GBASE-CR10): EEE PMD changes

Any new clauses created for 802.3bj, PMD and others will require appropriate additions.

EEE for 100 Gb/s Summary

- The majority of the EEE protocol that was developed for 10 Gb/s can be applied directly to EEE at 100 Gb/s (Low Power Idle)
- The main protocol hurdle that 100 Gb/s has to overcome in order to support EEE is the Alignment Marker distance and startup time due to this distance
- Presented is a methodology to change to the Alignment Marker protocol which will greatly reduce the start up time and enable efficient EEE support at 100 Gb/s

Thanks!