



1000BASE-T PHY Control State Diagram Modifications

**Niall Fitzgerald, Adam Healey,
Brian Murray, Jacobo Riesco
LSI Corporation**

**Presented by Adam Healey
IEEE P802.3az Task Force Meeting
New Orleans, LA
January 2009**

Comments #98 and #101: Problem statement

- PHY Control state diagram allows a transition from UPDATE to WAKE to be forced at any time by the assertion of `loc_lpi_req = FALSE`
 - This results in continued transmission for `lpi_waketx_timer` followed by a period of silence (`tx_mode = SEND_Z`) no less than `lpi_wakemz_timer`
- This implies that the link partner's update of timing and adaptive filter coefficients could be interrupted at any time
- This permits pathological timing scenarios where `LP_IDLE` is asserted at the GMII such that the PHYs transitions to the UPDATE state and then `LP_IDLE` is de-asserted forcing the link partner to abort update of timing and adaptive filter coefficients
 - Repetitions of this timing cycle can starve the PHY of essential updates and degrade link performance
 - This issue could also be addressed by enforcing a minimum period that the “power management agent” must assert `LP_IDLE`

Approaches to comments #98 and #101 – 1

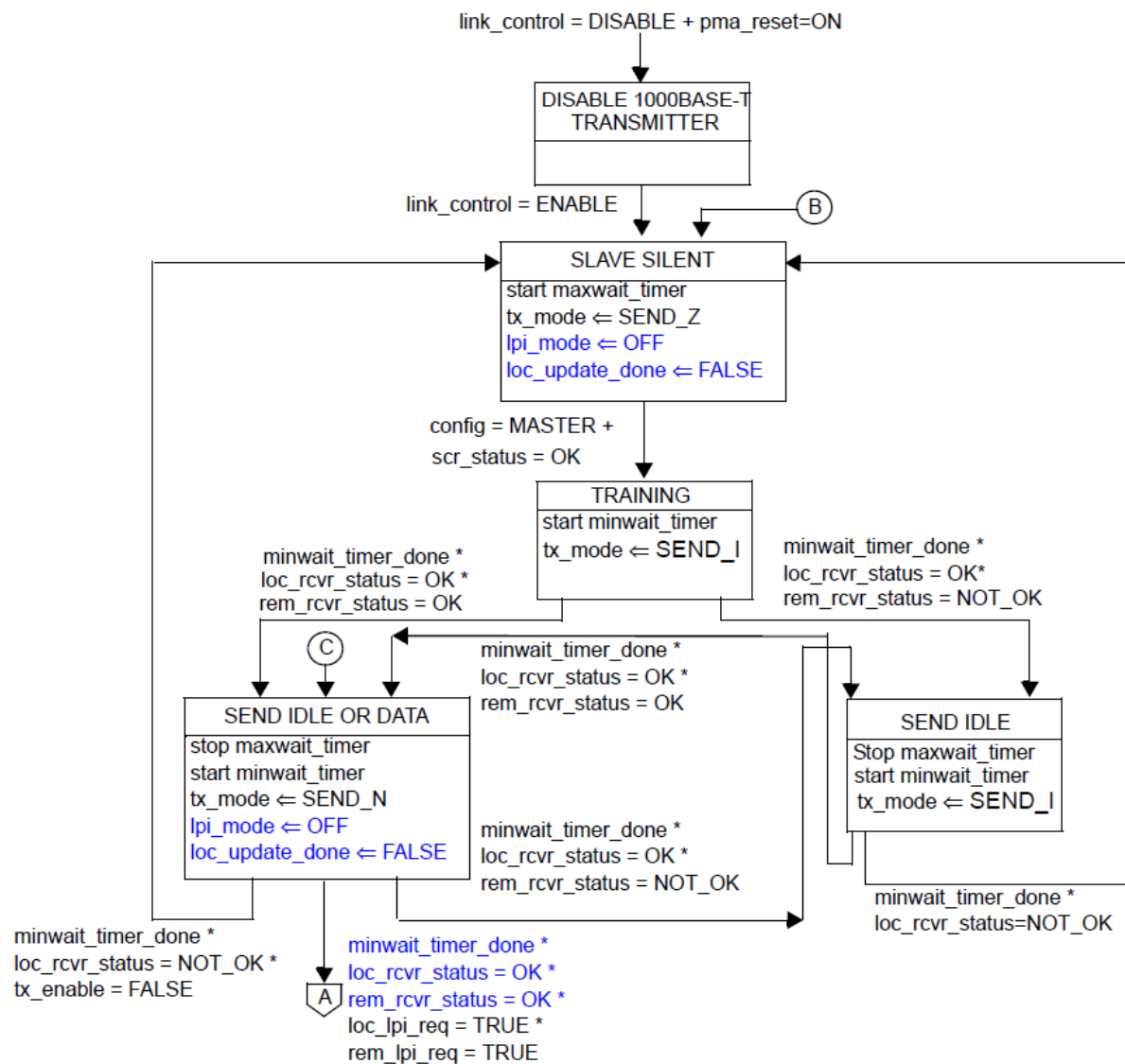
- Define the minimum time the “power management agent” must assert LP_IDLE
 - To ensure that both the local device and link partner both enjoy a period of uninterrupted transmission of a least lpi_update_timer (T_u)
 - No less than $2T_u(\text{min.}) + T_w(\text{min.})$, where T_w corresponds to lpi_wake_timer
 - This translates directly to the size of the buffer that must be maintained by the transmitter
- Define the minimum time the agent must wait between de-asserting LP_IDLE and asserting LP_IDLE again
 - Again, to ensure a period lpi_update_timer of uninterrupted transmission
 - No less than the minimum value of T_u

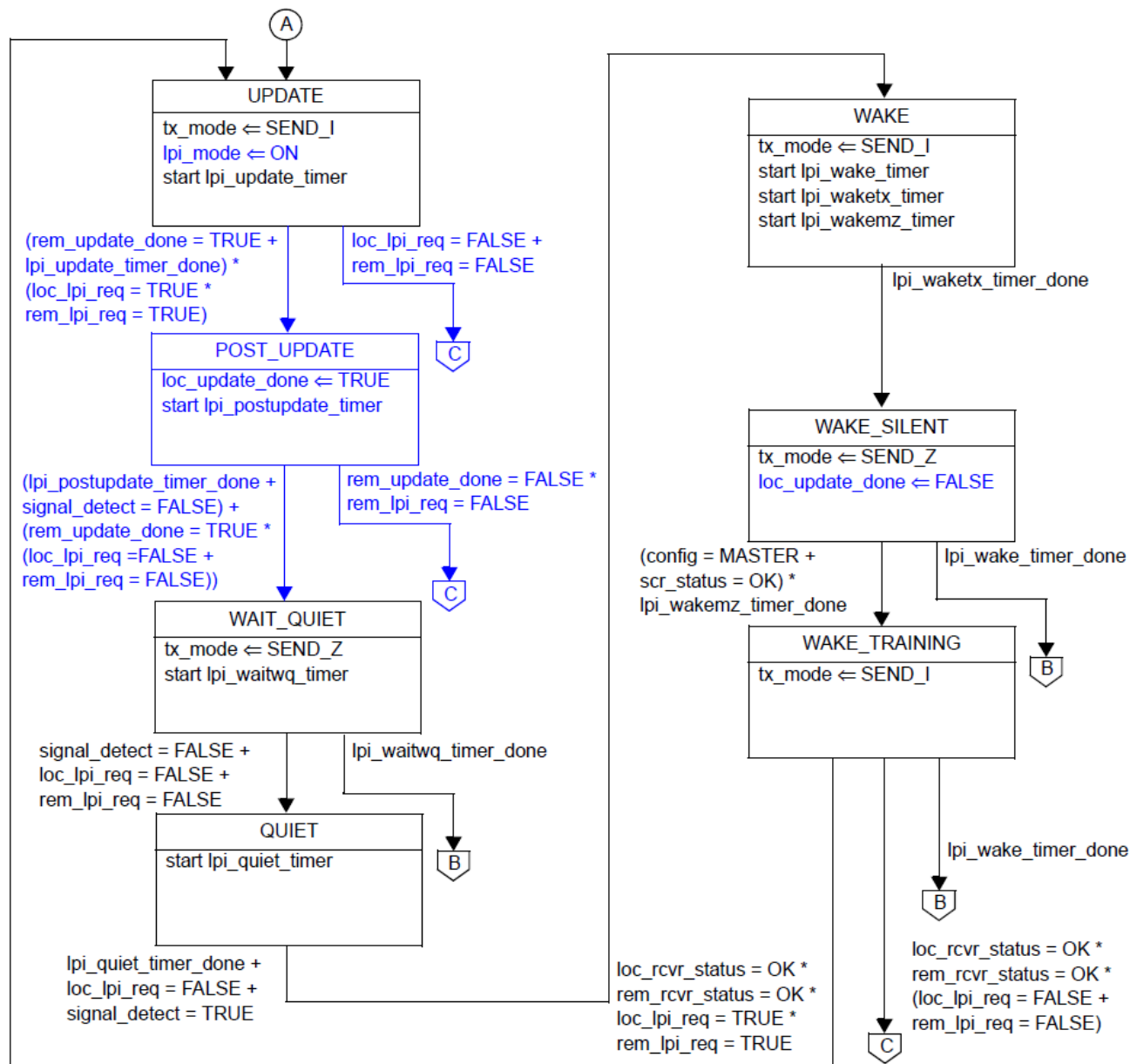
Approaches to comments #98 and #101 – 2

- However, these rules really address an issue with the 1000BASE-T PHY Control state diagram
 - Appropriate changes to the PHY Control state diagram would ensure proper operation of the PHY without any additional restrictions on the agent
 - Avoid unwanted dependencies between proper operation of the agent and proper operation of the PHY
 - Agent does not need to make special provisions for a 1000BASE-T PHY
 - Address the root cause of the issue rather than consider work-arounds

Summary of proposed changes

- Introduce new POST_UPDATE state, succeeding the UPDATE state, controlling transitions into WAIT_QUIET or SEND IDLE OR DATA
- Restore lpi_mode to its Draft 1.0 definition
- Introduce new variable loc_update_done
 - Indicates completion of timing and adaptive filter coefficient updates
 - Assigned a value of FALSE prior to entering the UPDATE state
 - Assigned a value of TRUE in the POST_UPDATE state
 - Communicated to the link partner and received as rem_update_done
 - Use the same encoding rules adopted for loc_lpi_mode (possibly modified by comment #9).
- Remove the transition from WAKE_TRAINING to WAKE_SILENT
 - It was added to combat a fall-through case in the Draft 1.0 state diagram which no longer exists in Draft 1.1
- Remove lpi_waitwt_timer
 - It was added to combat a fall-through case in the Draft 1.0 state diagram which no longer exists in Draft 1.1





Highlights

- A direct transition is provided from UPDATE (or POST_UPDATE) to SEND IDLE OR DATA if the link partner has not yet completed filter coefficient updates (e.g. `rem_update_done = FALSE`)
 - Update of adaptive filter coefficients may continue uninterrupted
- When the remote PHY has signaled completion of update then the transition through to the wake sequence is possible
- Duration of `lpi_postupdate_timer` is required to be greater than one round-trip delay
 - Propose a range of 2.0 and 2.2 microseconds
- If `loc_lpi_req = FALSE` during POST_UPDATE, then the local device must wait for `rem_update_done = TRUE` before proceeding to WAKE
 - This will not add time to the overall wake time budget

Comment #102: Problem statement

- Failure to achieve both `loc_rcvr_status = OK` and `rem_rcvr_status = OK` prior to `lpi_wake_timer_done` causes PHY Control to transition to the SLAVE SILENT state and initiate re-training
 - This will correspond to an interruption of service spanning hundreds of milliseconds
- What are the consequences of not waking within the allotted time?
 - Packet(s) transmitted immediately after `lpi_wake_timer_done` could be lost
 - For this reason, it is imperative to set PHY parameters so that the chances of failing to wake within the allotted time are very small
- During a refresh or when system wake time greatly exceeds the PHY wake time, the consequences are minor
 - No data loss, perhaps a very small compromise of power savings (e.g. refresh may be slightly longer on occasion)
- Consequences are considerably more severe in all cases when re-training is enforced

Summary of proposal

- Use lpi_wake_timer to monitor the health of the link
 - Define that lpi_wake_timer_done causes a new counter, “1000BASE-T wake error,” to be incremented
 - Counter is represented in the Clause 45 management register space and is cleared on read
 - System management reads the counter to understand if the link is failing to recover from low-power mode within the allotted time and takes corrective actions as necessary
- Define a new timer, lpi_link_fail_timer
 - Functionally replaces lpi_wake_timer in the PHY Control state diagram, e.g. expiration triggers re-training
 - Started in the WAKE state
 - Propose timer value to be 90 to 110 microseconds
- Add action “Stop lpi_wake_timer” to SEND IDLE OR DATA to prevent lpi_wake_timer_done from being satisfied after successful wake



Questions?