

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Proposed Text for Annex 90A

Richard Tse, Microchip Technology

Steve Gorshe, Microchip Technology

Marek Hajduczenia, Charter Communications

Presented at P802.3cx teleconference, May 25, 2021

27 **Supporters**

- 28 • Andras de Koos, Microchip Technology
- 29 • Clark Carty, Cisco Systems
- 30 • Mark Bordogna, Intel
- 31 • Sriram Natarajan, Cisco Systems

32 **Annex 90A (informative) Timestamping Accuracy Considerations**

33 **90A.1 High Accuracy Timestamping Introduction**

34 This annex provides information on supporting high accuracy timestamping for time synchronization
35 protocol (TimeSync) Client implementations compliant with Clause 90. This timestamping may be used
36 for time synchronization protocols including IEEE Std 1588 and IEEE Std 802.1AS.

37 **90A.2 High Accuracy Timestamping Background**

38 Ethernet support for time synchronization protocols (Clause 90) was not originally specified to support
39 high accuracy timestamping. Thus, implementation flexibility permitted by this standard prior to the
40 addition of certain registers (IEEE Std 802.3cx support, Timestamp reference, first symbol after SFD,
41 Multilane support, and TX/RX num_unit_change support registers shown in Table 45-235 and the fine
42 resolution path data delay registers located throughout subclauses 45.2.1 to 45.2.6) could lead to
43 timestamp accuracy impairments that might not satisfy high accuracy timing requirements.

44 Timestamping accuracy can be impaired when two TimeSync Clients do not account for a varying
45 physical layer device (PHY) path data delay in the same manner. Examples of PHY functions that cause
46 variation in the PHY path data delay include alignment marker (AM) or codeword marker (CWM)
47 insertion/removal, Idle insertion/removal, and multi-physical coding sublayer (PCS) lane
48 distribution/merging.

49 Timestamping accuracy can also be impaired when two TimeSync Clients do not use the same message
50 timestamp point. As specified in 90.7, this standard gives two options for the message timestamp point
51 (the beginning of the start of frame delimiter, the SFD, and the beginning of the first symbol after the
52 SFD) but recommends using the beginning of the first symbol after SFD, which is consistent with IEEE Std
53 1588 and IEEE Std 802.1AS.

54 Table 90A-1 shows the magnitude of potential timestamp accuracy impairments that could be
55 generated by the aforementioned causes.

Table 90A-1 – Magnitude of Potential Timestamp Accuracy Impairments

Ethernet Rate	Magnitude of Potential Timestamp Accuracy Impairments per Transmit or Receive Port (ns)			
	Mismatched Message Timestamp Point ¹	Idle Insertion/Removal ^{2,3}	AM/CWM Insertion/Removal ³	PCS Lane Distribution/Merging
10M	800	400	N/A	N/A
100M	80	40	N/A	N/A
1G	8	16 ⁴ , 8 ⁵	N/A	N/A ⁴ , 0 ⁵ , N/A ^{6,7}
2.5G	3.2	12.8	N/A	N/A ⁷
5G	1.6	6.4	N/A	N/A ⁷
10G	0.8	3.2	N/A	N/A ⁴ , 0 ⁵
25G	0.32	1.28	10.24	N/A
40G	0.2	1.6	6.4	4.8
100G	0.08	0.64	12.8	12.16
200G	0.04	0.32	2.56	N/A ⁷
400G	0.02	0.16	2.56	N/A ⁷

57 Notes:

- 58 1. The value shown only accounts for the time between the two message timestamp point options
59 when they are adjacent. See Annex 90A.3 for other factors that can affect some of these values.
60 2. The value shown corresponds to the effect of a single Idle insertion/removal.
61 3. The path data delay of a TimeSync message is only affected when the message coincides with an
62 AM, CWM, or Idle insertion/removal event.
63 4. For 1000Base-X or 10GBase-R
64 5. For 1000Base-T or 10GBase-X
65 6. For 10GBase-T
66 7. For these rates, the lane distribution/merging operation belongs only to the forward error
67 correction (FEC) function and not to the PCS function. The FEC lane distribution/merging
68 operation is not subject to potential timestamp accuracy impairments because its path data
69 delay determination was already clearly defined by the original specification, IEEE Std 802.3-
70 2018, and not subject to implementation flexibilities.

71 90A.3 Considerations for Use of Different Message Timestamp Points

72 If two TimeSync Clients use different message timestamp points, a timestamp accuracy impairment
73 equal to the time difference between the two message timestamp points will be incurred on the
74 TimeSync link delay measurement. The magnitude of this impairment, as shown in Table 90A-1, is the
75 time difference between the beginning of the SFD message timestamp point and the beginning of the
76 first symbol after the SFD message timestamp point when they are adjacent to each other, which they
77 normally are. For implementations that do not use the TX/RX num_unit_change support and Multilane
78 support registers (see Table 45-235), an additional impairment could result if these two message
79 timestamp point options are further separated at due to:

- 80 • Insertion of bytes between the two message timestamp points for AM or CWM functions

- 81 • Multi-PCS lane distribution delays
- 82 Implementations compliant to this version of the standard only suffer a timestamp accuracy impairment
83 of one byte time between the two message timestamp point options because:
- 84 • The effect of AM or CWM insertion is accounted for, using the Tx_num_unit_change and
85 Rx_num_unit_change primitives (see 90.4.3.3, 90.4.3.4, and Annex 90A.5)
- 86 • The multi-PCS lane path data delay is modelled as a constant value for all PCS lanes (see 90.7
87 and Annex 90A.4).

88 **90A.4 Considerations for Multi-PCS Lane Functions**

89 The general concept used to accommodate the delay variation of the multi-PCS lane
90 distribution/merging operation is explained in Annex 90A.7. This concept takes advantage of the fact
91 that the sum of the intrinsic delay variation of the transmit (Tx) multi-PCS lane distribution operation
92 and of the intrinsic delay variation of the receive (Rx) multi-PCS lane merging operation is a
93 predetermined constant for the given multi-PCS lane function.

94 This concept allows the intrinsic delay variations to be modelled as constant values, thus enabling the
95 static TimeSync PCS transmit path data delay register and TimeSync PCS receive path data delay register
96 to be used with high accuracy timestamping even when multi-PCS lane functions are present. As
97 explained in 90.7, the TimeSync PCS transmit path data delay register would use the greatest PCS lane
98 distribution delay as its constant value (which corresponds to the start of the Tx PCS lane distribution
99 function) and the TimeSync PCS receive path data delay register would use the smallest PCS lane
100 merging delay as its constant value (which corresponds to the start of the Rx PCS lane merging function).

101 Because the PCS transmit path data delay is modelled as a constant value, the minimum and maximum
102 TimeSync PCS transmit path data delay registers in an ideal implementation have the same value due to
103 the multi-PCS lane distribution operation. Likewise, because the PCS receive path data delay is modelled
104 as a constant value, the minimum and maximum TimeSync PCS receive path data delay registers in an
105 ideal implementation have the same value due to the multi-PCS lane merging operation. Having
106 identical minimum and maximum values in these registers indicates that there is no uncertainty in the
107 PCS path data delay.

108 The above consideration of the multi-PCS lane distribution/merging operation is consistent with that for
109 the multi-FEC lane distribution/merging operation.

110 **90A.5 Considerations for AM/CWM and Idle Functions**

111 Timestamp accuracy impairment can occur because AM, CWM, or Idle insertion and removal events
112 cause an instant change in the PCS path data delay. Unlike other PHY functions, these events do not
113 generate PCS path data delay variations that can be pre-determined and the Tx path data delay variation
114 is not mirrored by the Rx path data delay variation.

115 Each of these path data delay variations may be accounted for by using the Tx_num_unit_change and
116 Rx_num_unit_change primitives (see 90.4.3.3 and 90.4.3.4). These primitives allow the TimeSync Client
117 to compensate for the instant change in the path data delay. Because the primitives compensate for the
118 instant path data delay changes, the TimeSync PCS transmit path data delay register and TimeSync PCS

119 receive path data delay register can operate as static values, even when AM, CWM, or Idle
120 insertion/removal operations are present.

121 Examples that show how Tx_num_unit_change and Rx_num_unit_change may be used are given in
122 Annex 90A.5.1 and Annex 90A.5.2, respectively.

123 *[Note: do we need to add figures to Annex 90A.5.1 and Annex 90A.5.2 to help illustrate the examples?]*

124 **90A.5.1 Example use of Tx_num_unit_change**

125 1. Scenario without AM, CWM, or Idle insertion/removal event:

- 126 • Arrival time of a message timestamp point at the Tx xMII = T1
- 127 • Tx PCS path data delay = PDD1
 - 128 • The constant value, PDD1, is programmed into the TimeSync PCS transmit path data
 - 129 delay registers
- 130 • Calculated Tx departure timestamp = T1 + PDD1

131 2. Scenario with AM, CWM, or Idle insertion/removal in which Tx_num_unit_change is used to account
132 for the Tx PCS path data delay variation, allowing the Tx PCS path data delay to be modelled as a
133 constant:

- 134 • Adjusted arrival time of the message timestamp point at the Tx xMII = T1 +
135 Tx_num_unit_change*(nanoseconds/unit)
 - 136 • The arrival time at the Tx xMII is modified to reflect the AM, CWM, or Idle
 - 137 insertion/removal event (as if it happened before the Tx xMII, per 90.7)
 - 138 • The value of Tx_num_unit_change is positive when data is inserted ahead of the
 - 139 message timestamp point, increasing the Tx path data delay, and negative when data is
 - 140 removed ahead of the message timestamp point, decreasing the Tx path data delay
- 141 • Tx PCS path data delay = PDD1
 - 142 • The constant value, PDD1, programmed into the TimeSync PCS transmit path data delay
 - 143 registers does not change.
- 144 • Calculated Tx departure timestamp = T1 + (PDD1 + Tx_num_unit_change*(nanoseconds/unit))

145 **90A.5.2 Example use of Rx_num_unit_change**

146 1. Scenario without AM, CWM, or Idle insertion/removal event:

- 147 • Arrival time of a message timestamp point at the Rx xMII = T2
- 148 • Rx PCS path data delay = PDD2
 - 149 • The constant value, PDD2, is programmed into the TimeSync PCS receive path data
 - 150 delay registers
- 151 • Calculated Rx arrival timestamp = T2 – PDD2

152 2. Scenario with AM, CWM, or Idle insertion/removal in which Rx_num_unit_change is used to account
153 for the Rx PCS path data delay variation, allowing the Rx PCS path data delay to be modelled as a
154 constant:

- 155 • Adjusted arrival time of the message timestamp point at the Rx xMII = $T2 +$
156 $Rx_num_unit_change * (nanoseconds/unit)$
- 157 • The arrival time at the Rx xMII is modified to reflect the AM, CWM, or Idle
158 insertion/removal event (as if it happened after the Rx xMII, per 90.7)
- 159 • The value of Rx_num_unit_change is positive when data is inserted ahead of the
160 message timestamp point, increasing the Rx path data delay, and negative when data is
161 removed ahead of the message timestamp point, decreasing the Rx path data delay
- 162 • Rx PCS path data delay = PDD2
- 163 • The constant value, PDD2, programmed into the TimeSync PCS receive path data delay
164 registers does not change.
- 165 • Calculated Rx arrival timestamp = $T2 - (PDD2 + Rx_num_unit_change * (nanoseconds/unit))$

166 **90A.5.3 Considerations for Implementations without Tx_num_unit_change and Rx_num_unit_change**

167 For an implementation that does not compensate for the path data delay variation resulting from AM,
168 CWM, or Idle insertion/deletion removal events (e.g., without the Tx_num_unit_change and
169 Rx_num_unit_change primitives), the effect of the timestamp accuracy impairments that result from
170 these events can be evaluated to determine if they cause significant degradation in the TimeSync
171 system's performance. Some observations that might help this evaluation are given below:

- 172 • Typically, the probability that an AM, CWM, or Idle insertion/deletion removal event affects the
173 path data delay of a TimeSync message is small.
- 174 • A low-pass filter, which might be present in a TimeSync Client's time recovery algorithm, could
175 attenuate the effect of the resulting impairments.
- 176 • An implementation that does not transmit TimeSync messages in the region of AM/CWM
177 insertions or Idle insertions/removals might avoid the generation of the impairment at its Tx
178 port. However, this does not guarantee that the corresponding remote Rx port will not generate
179 impairments from its own Idle insertions/removals.

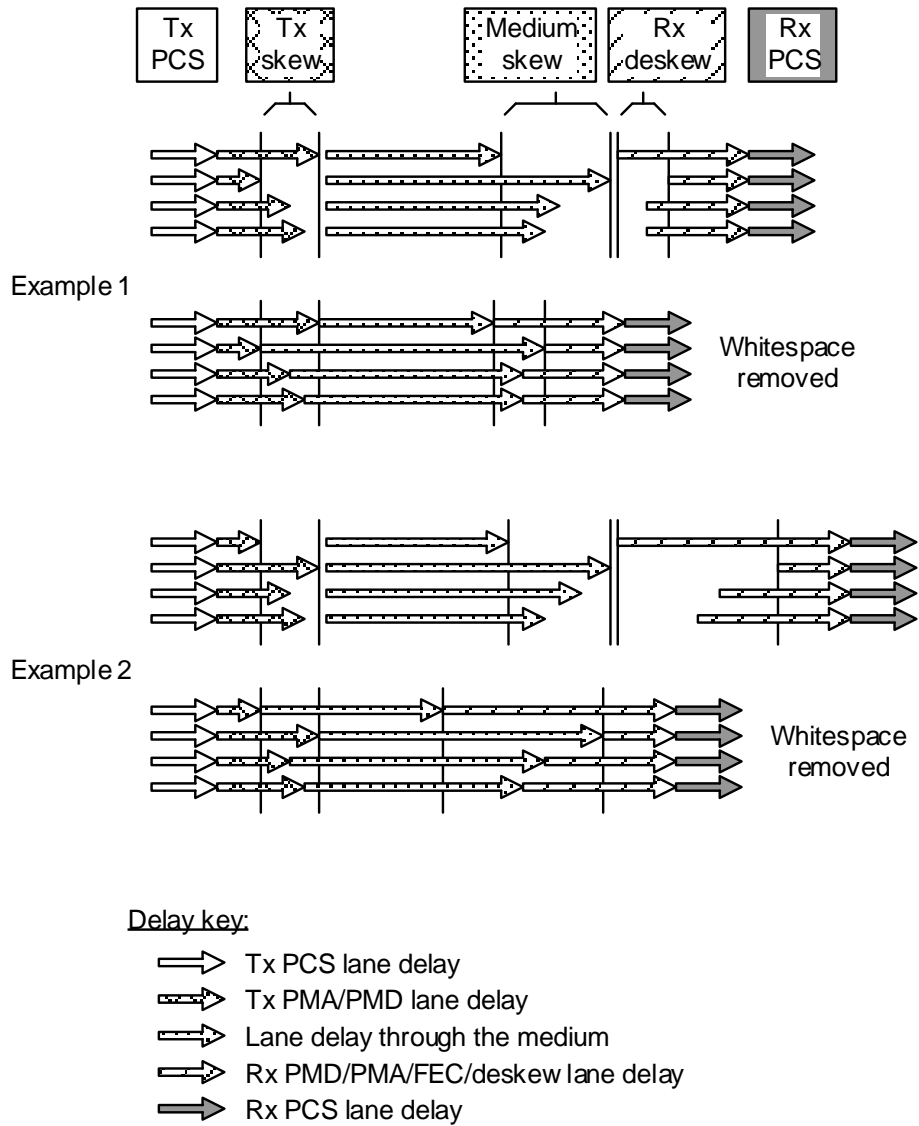
180 **90A.6 Considerations for Tx Skew**

181 For a multi-lane PHY, the receiver accounts for the skew of the medium by timestamping with respect to
182 the lane with the smallest deskew buffer delay (see 90.7). Thus, the medium delay used by the time
183 synchronization protocol is that of the lane with the greatest delay.

184 For these multi-lane PHYs, the presence of skew at the transmit Medium Dependent Interface (MDI) is
185 difficult to compensate for because this skew is entwined with but independent from the skew of the
186 medium. As shown in the examples of Figure 90A-1, the transmit skew in series with the medium skew
187 can either be additive or subtractive. Example 1 and Example 2 in Figure 90A-1 have the same transmit
188 skew and the same medium skew but, because these skews are associated differently, the total skew

189 seen at the receiver is different. By obscuring the latency of the medium, transmit skew can contribute
 190 to time synchronization error.

191 The per-lane transmit skew values might be compensated at the receiver, where the total skew of each
 192 lane can be observed at its Rx deskew buffers. By using the observed per-lane total skew values at the
 193 receiver and the per-lane transmit skew values, the actual skew of each lane of medium could be
 194 determined. To negate the need for this type of processing, it is recommended that multi-lane
 195 transmitter implementations try to minimize the lane skew at their MDI.



196

197 **Figure 90A-1 – Transmit PMA/PMD Skew in Series with Medium Skew**

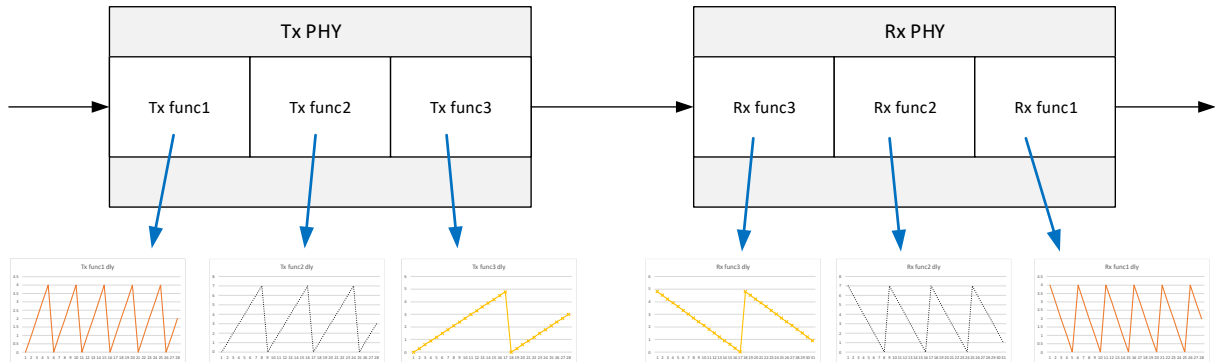
198 **90A.7 General Method for Dealing with Repeating Delay Variation Patterns**

199 Many PHY functions have varying intrinsic delays with the following characteristics:

- 200 • The Tx and the Rx intrinsic delay variations follow a known repeating pattern.

- 201 • The Tx intrinsic delay variation pattern is a mirror of the Rx intrinsic delay variation pattern and the
 202 sum of the two intrinsic delays is a known constant value. This is true because the data stream
 203 before the Tx function and the Rx stream after the Rx function are identical.

204 It is possible to take advantage of the above characteristics to simplify the path data delay modeling. For
 205 example, if a PHY has multiple functions with these delay characteristics, as shown in Figure 90A-2, its
 206 aggregated path data delay may be modelled as a constant value instead of the dynamically varying sum
 207 of multiple varying delays.

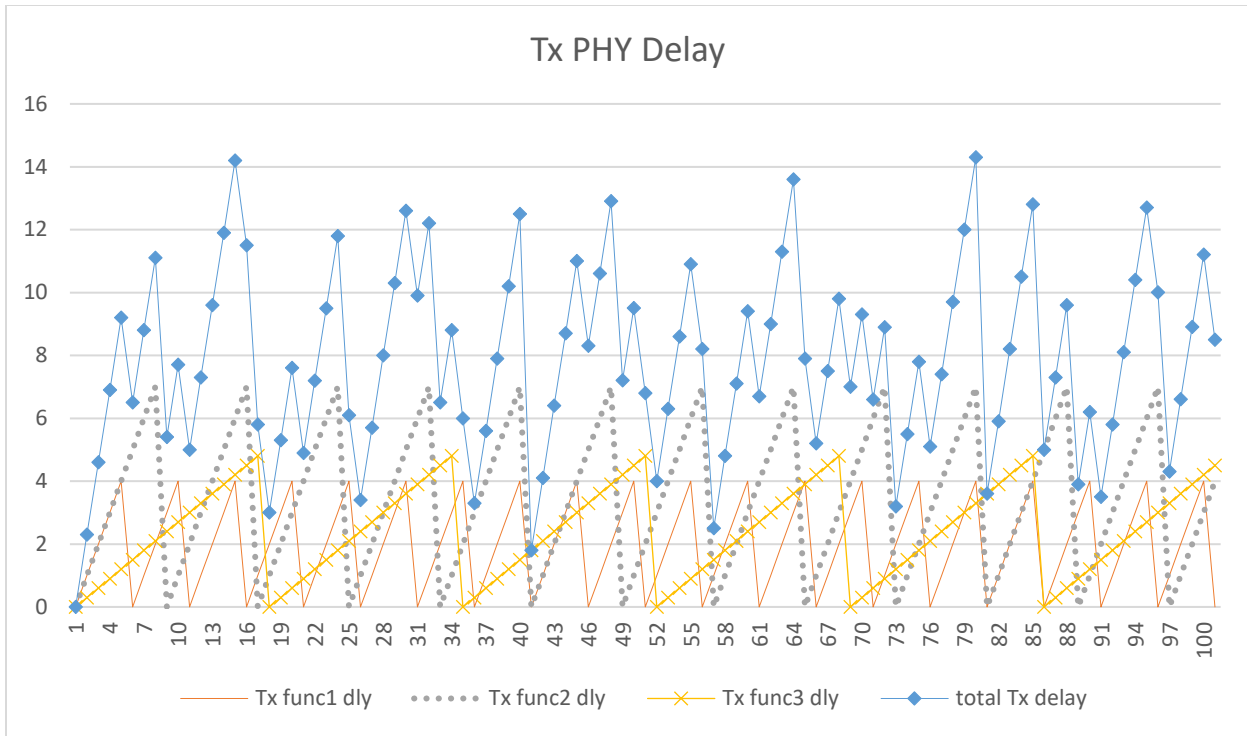


208

209

Figure 90A-2 – PHY with Cascaded Functions with Varying Delays

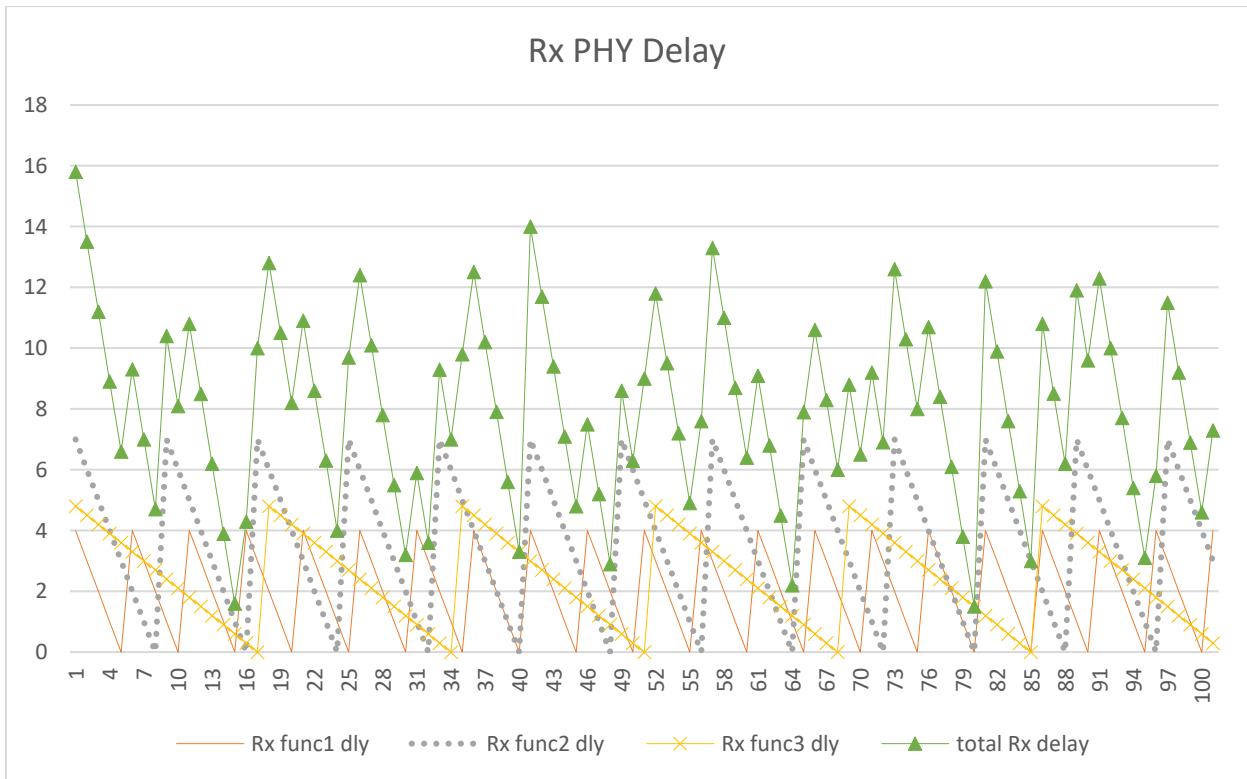
210 For the example shown in Figure 90A-2, the delay of each Tx function and the sum of them are shown in
 211 Figure 90A-3 and the delay of each Rx function and the sum of them are shown in Figure 90A-4. These
 212 sums have no easily discernable pattern and might require an implementation to determine the
 213 instantaneous path data delay for any chosen bit that corresponds to the message timestamp point of a
 214 TimeSync message in the Ethernet data stream.



215

216

Figure 90A-3 – Total Delay of Tx PHY with Cascaded Varying Delays



217

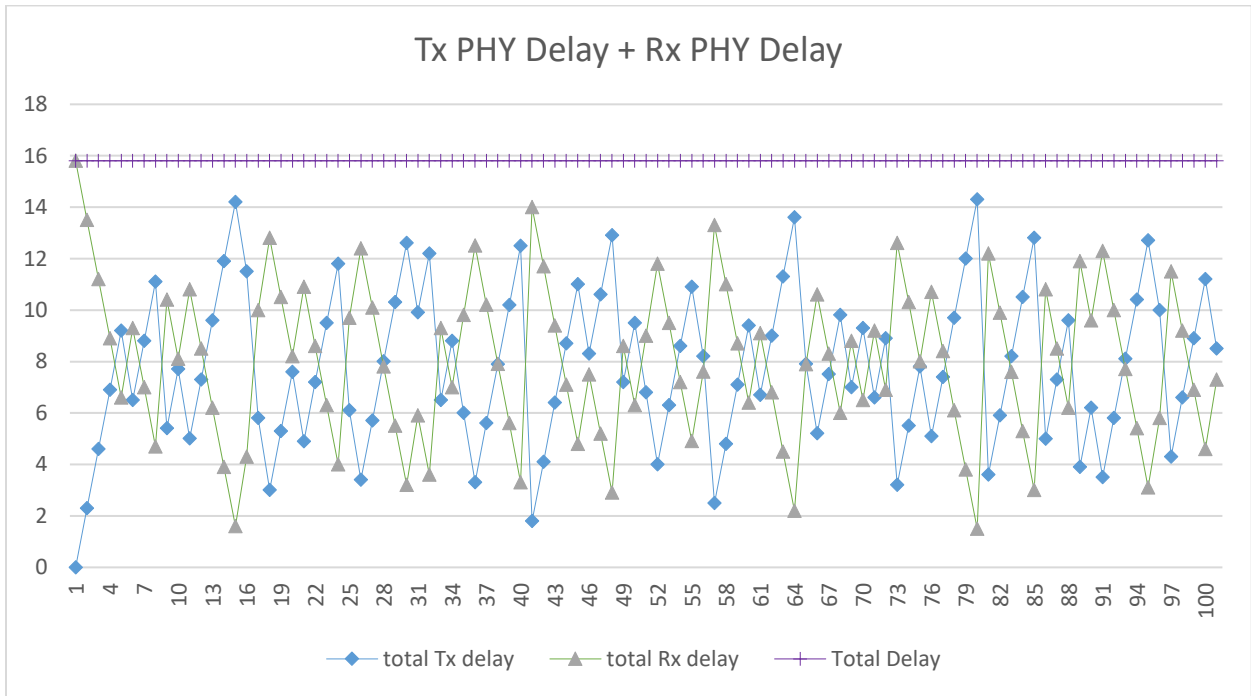
218

Figure 90A-4 – Total Delay of Rx PHY with Cascaded Varying Delays

219 Because the intrinsic varying delay in the Rx PHY is a mirror of the intrinsic varying delay in the Tx PHY,
220 the total intrinsic delay through both PHYs is a constant, as illustrated in Figure 90A-5. This eliminates
221 the need to track the varying delay of the message timestamp point of a TimeSync message through the
222 Tx PHY and the Rx PHY. Instead, it is possible to divide the aggregate constant delay into Tx and Rx
223 portions and allocate them to the individual PHY instances. The allocated portions of the constant total
224 intrinsic delay value are then added to the implementation-specific delays, which are also constant
225 values, of the corresponding Tx and Rx PHYs to compensate these intrinsic varying delays into
226 timestamps.

227 It is recommended to use this method to model all varying PHY delays of this nature.

228



229

230 **Figure 90A-5 – Tx PHY Delay, Rx PHY Delay, and Total Delay**

231

232 *[Note: Should we give examples of how this method can be used for existing basic functions such as*
233 *64B/66B encoding/decoding, 2x32B to 66B encoding/decoding, 256B/257B transcoding?]*

234