



PCIe® 3.0 Electricals – Part III

Debendra Das Sharma
Intel Corporation



Agenda

- Problem Statement
- Existing Usage of K-Codes
- Metrics considered for evaluation
- Current Direction on Encoding
- Summary & Call to Action

Disclaimer: Information contained here is **Work In Progress**.
Specification is at Rev 0.3 level.

Problem Statement

- PCIe® 3.0 data rate decision: 8 GT/s
 - ✓ High Volume Manufacturing channel for client/ servers
 - Same channels and length for backwards compatibility assuming worst-case
 - ✓ Low power and ease of design
 - Avoid using complicated receiver equalization, etc.
- Requirement: **Double Bandwidth** from PCIe 2.0
 - ✓ PCIe 1.0a data rate: 2.5 GT/s
 - ✓ PCIe 2.0 data rate: 5 GT/s
 - Doubled the bandwidth from PCIe 1.x to PCIe 2.0 by doubling the data rate
 - ✓ Data rate gives us a 60% boost in bandwidth
 - ✓ Rest will come from **Encoding**
 - Replace 8b/10b encoding with a scrambling-only encoding scheme when operating at PCIe 3.0 data rate
- Double B/W: Encoding efficiency improvement of 1.25 X data rate improvement of 1.6 yields
- **Challenge:** 8b/10b encoded the 2^8 data patterns and 12 K-codes

Agenda

- Problem Statement
- Existing Usage of K-Codes
- Metrics used for evaluation
- Current Direction on Encoding
- Summary & Call to Action

Existing Usage of K-Codes

- Two flavors for K-code use
 - ✓ Packet Stream (independent of link width)
 - ✓ Lane Stream (per-lane)
- Packet Stream relates to Packet Framing (Link-Wide)
 - ✓ STP - Start of TLP
 - ✓ END - End (Good) of TLP
 - ✓ EDB - End Bad of TLP
 - ✓ SDP - Start of DLLP
- Lane Stream relates to Ordered Sets:
 - ✓ Training Set #1 & #2
 - Link training and negotiation
 - ✓ SKP Ordered Sets
 - Periodic link clock compensation
 - Recovery from bit slip/add
 - ✓ Electrical Idle Start/ Exit sequence
 - Power management
- New encoding scheme needs to accommodate these existing usages

Agenda

- Problem Statement
- Existing Usage of K-Codes
- **Metrics used for evaluation**
- Current Direction on Encoding
- Summary & Call to Action

Error Detection Ability

- Robustness against bit errors considered
 - ✓ Bit flip, bit slip/add
- **Basic Fault Model:**
 - ✓ Guaranteed error detection against random bit flips in any TLP or DLLP or IDL or Ordered Set
 - Must not alias to a TLP or a DLLP with up to three bit flips
 - Can cause data corruption or flow-control problems
- No guaranteed detection of error with bit slip/add
 - ✓ Same as 2.0 ability
- No self healing
 - ✓ Errors cause transition to Recovery
- Eventual guaranteed recovery in the presence of multiple errors above including bit slip/add
- Need to handle killer packets
 - ✓ Send a different bit stream on retry of a packet

Other Metrics

- Bandwidth Inefficiency must be low enough
 - ✓ 8b/10b had a 20% inefficiency
 - ✓ New scheme must be in the 1-2% range for inefficiency
 - Would result in close to 2X the bandwidth from PCIe 2.0
- Time Overhead through Recovery as well as L0s/L1 exit must be minimal
 - ✓ Enables better power management without performance penalty
- Bytes continue to be the unit of transmission
 - ✓ Enables single-wide/double-wide type of parallel implementation
 - E.g., no end TLP in bit 3 and a new TLP starts in bit 4 within a byte
 - ✓ Preservation of framing rules and length of TLP/DLLP
- Switch to new encoding after speed change from electrical idle in Recovery.Speed
- Minimal changes beyond PHY layer
 - ✓ Ease of implementation

Agenda

- Problem Statement
- Existing Usage of K-Codes
- Metrics used for evaluation
- **Current Direction on Encoding**
- Summary & Call to Action

Overall Scheme

- Current direction is to use two levels of encapsulation
 - ✓ A 128/130 code on individual lanes
 - ✓ Physical layer packetization to identify “packet” boundaries
- Lane Level 128/130 Code
 - ✓ 130 bit code called a Block
 - ✓ Used for block lock
 - Substitutes COM used for Symbol lock in 8b/10b
 - ✓ Differentiate certain packet types
- Physical Layer packetization identifies packet boundaries. Packet types:
 - ✓ Link Level (TLP or DLLP)
 - ✓ Lane Level (Ordered Sets or LIDL)
- Scrambling only (no 8b/10b) to provide edge density
 - ✓ Additive scrambling on a per-lane basis
 - ✓ Current direction degree 23 polynomial for LFSR with different taps for adjacent lanes (up to 8)
 - ✓ Electrical Idle Exit Ordered Set resets scrambler (Recovery/ Config)

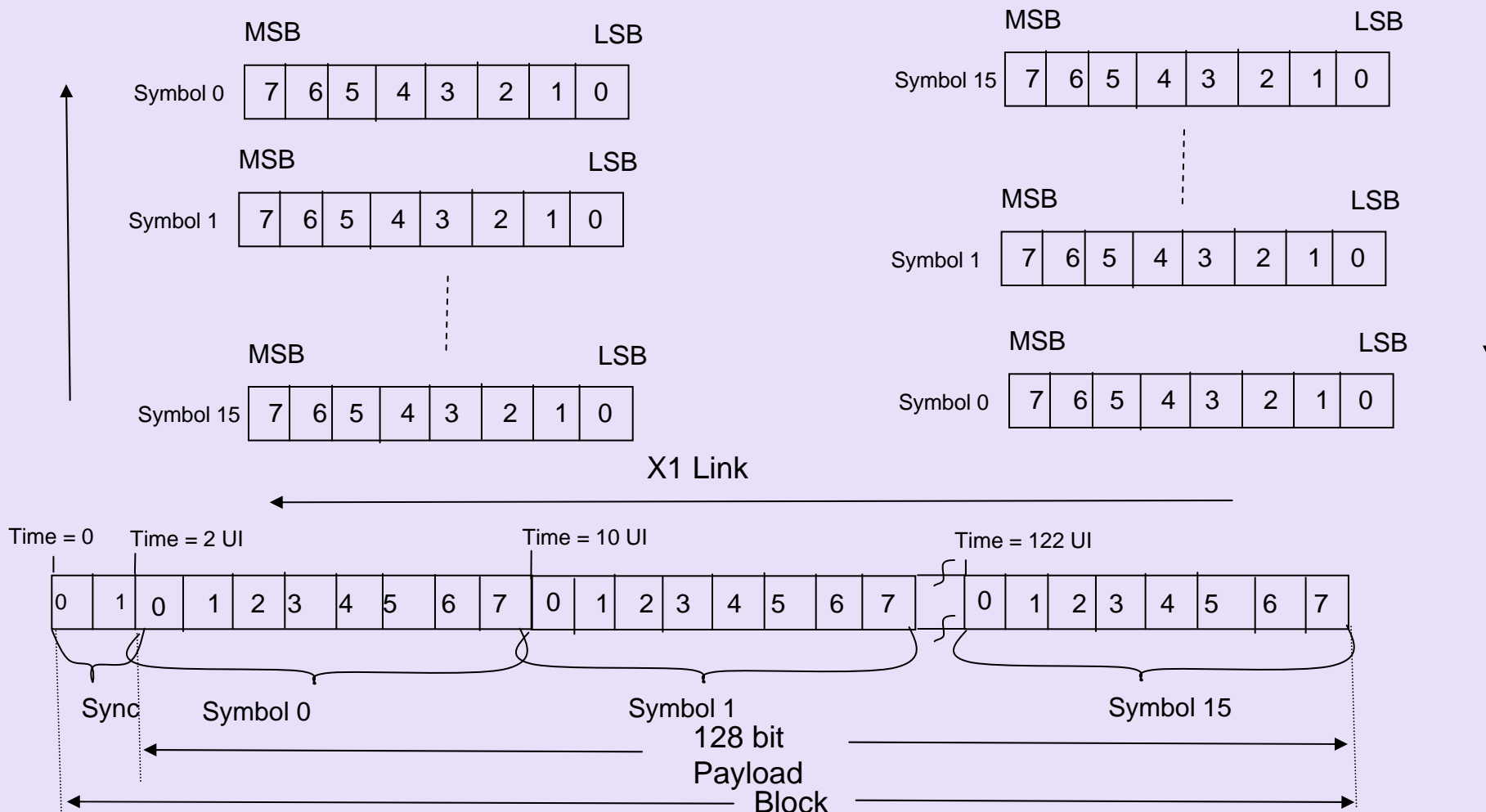
Lane Level Encoding

- Lane Level Encoding with a 128/130 code
 - ✓ 2-bit sync character followed by 128 bit payload
 - 128 bit payload may be any combination of
 - Part of one or more Link level packets (TLP, DLLP)
 - One or more lane level packets (IDL, Ordered Sets)
 - 2-bit sync character not scrambled
 - 128 bit payload may or may not be scrambled
 - ✓ 10_ sync character used for certain Ordered Sets
 - EIEOS not scrambled
 - 10_ [<00000000>_<11111111> ..8 times] represents EIEOS
 - Used for the dual purpose of low frequency patterns (in Recovery and Configuration) as well as establishing block lock
 - ✓ 01_ sync character is used when 128 bit information contains IDL, TLP, or DLLP
 - Each bit in the 128 bit payload is scrambled
 - ✓ Alternatives for other Ordered Set encodings considered

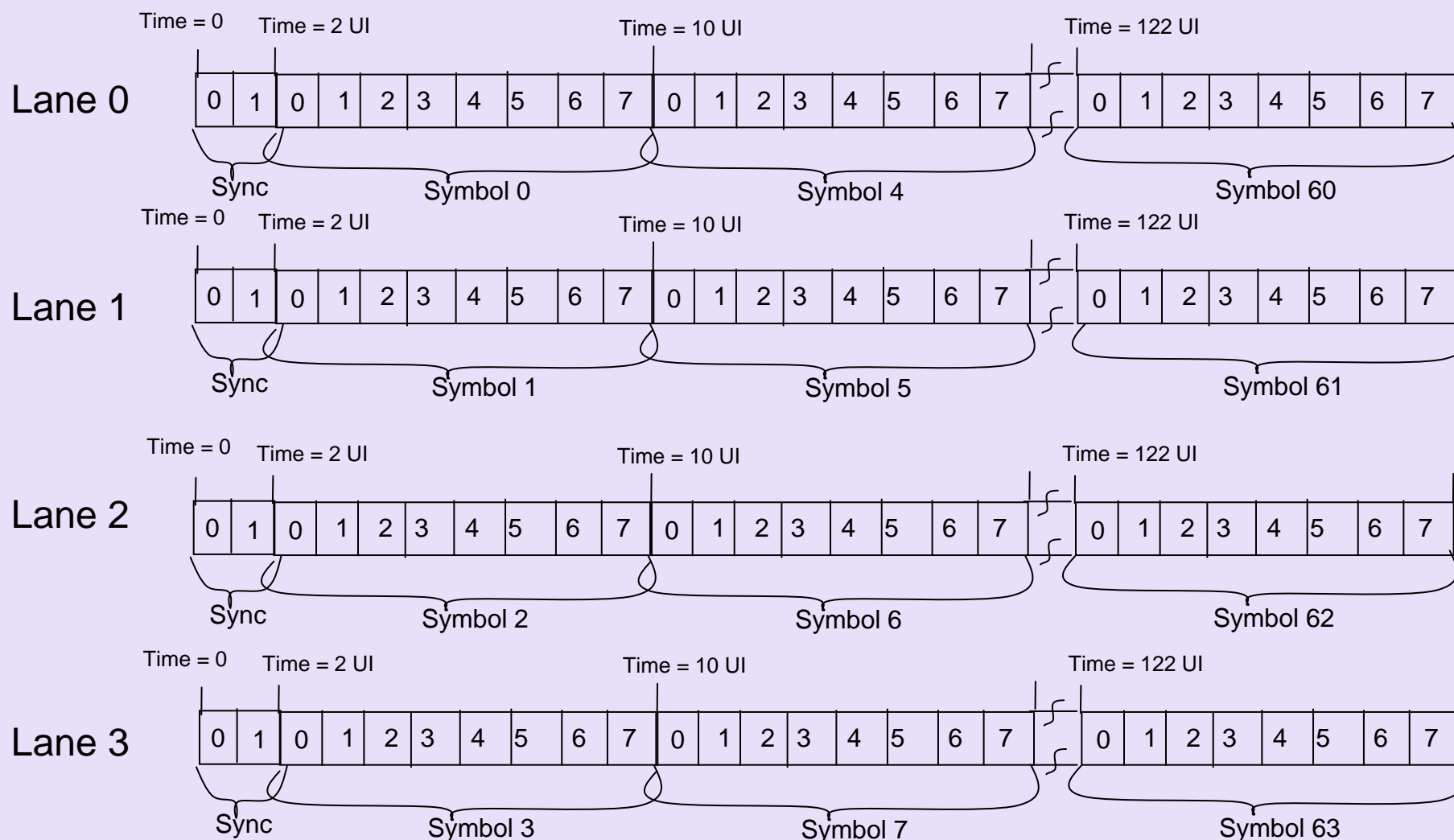
Mapping of bits on a x1 Link

Receive

Transmit



Mapping of bits on a x4 Link

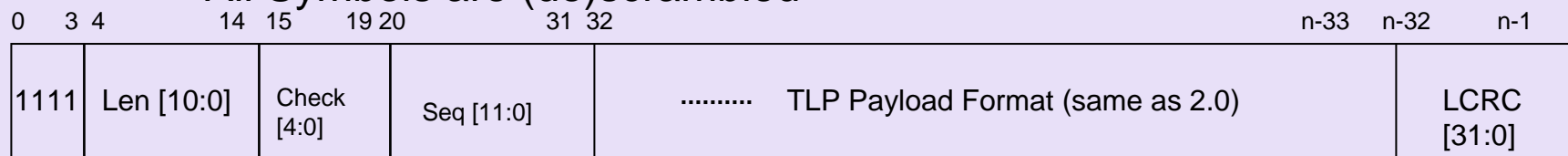


Physical Layer Encapsulation

- First Symbol (scrambled) indicates packet type:
 - ✓ 00000000 is Logical IDL
 - All subsequent lanes in same Symbol-time should be LIDL
 - Receivers check for all 0s (after descrambling) in LIDL
 - PAD functionality merged with LIDL
 - ✓ 1111xxxx is STP
 - Subsequent 11 bits (link wide) define the length
 - ✓ 00001111 is SDP
 - 2nd Symbol also gets a fixed encoding
 - ✓ 00000011 is EDB
 - EDB packet is 4 Symbols; each with the same value 00000011

P-Layer Encapsulation: TLP

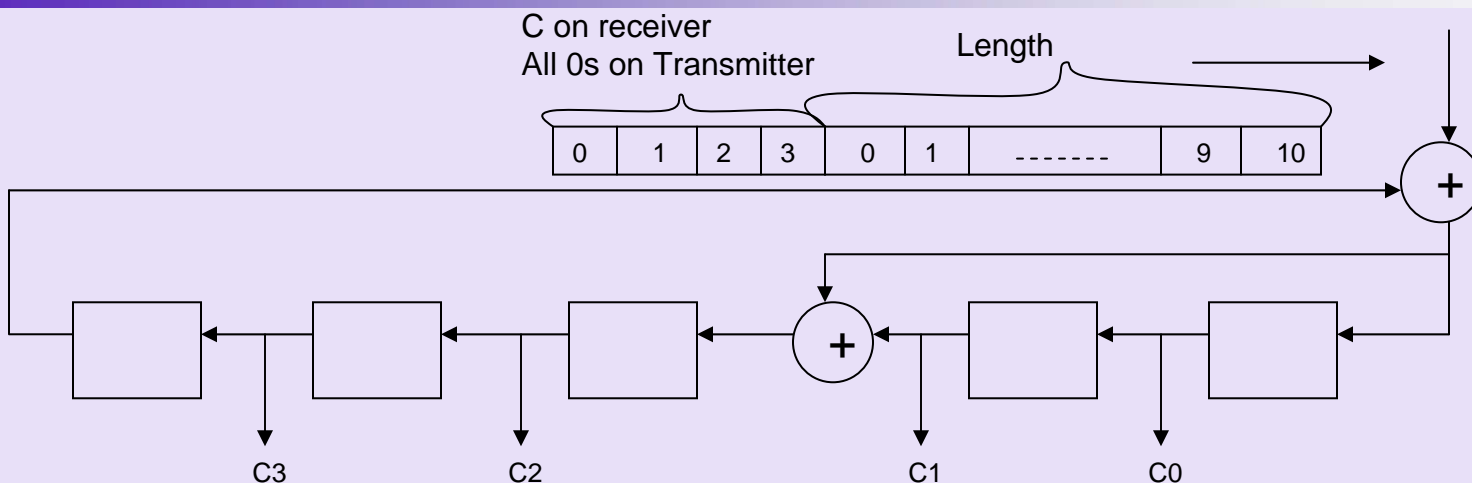
- Length known from the first 3 Symbols
 - ✓ First 4 bits are 1111 (bit[0:3] = 4'b1111)
 - ✓ Bits 4:14 has the length of the TLP (valid values: 5 to 1031)
 - ✓ Bits 15:19 is check bits to cover the TLP Length field
 - Primitive Polynomial ($X^4 + X + 1$) protects 15 bit field
 - Provides double bit flip detection guarantee (length 11 bits + CRC 4 bits)
 - Odd parity covers the 15 bits (length 11 bits + CRC 4 bits)
 - Guaranteed detection of triple bit errors (over 16 bits)
 - ✓ Sequence Number occupies bits 20:31
 - ✓ TLP payload is from the 4th Symbol position (same as 2.0)
 - ✓ No explicit END. Need to check first Symbol after TLP for implicit END vs an explicit EDB => Ensures triple bit flip detection
 - ✓ All Symbols are (de)scrambled



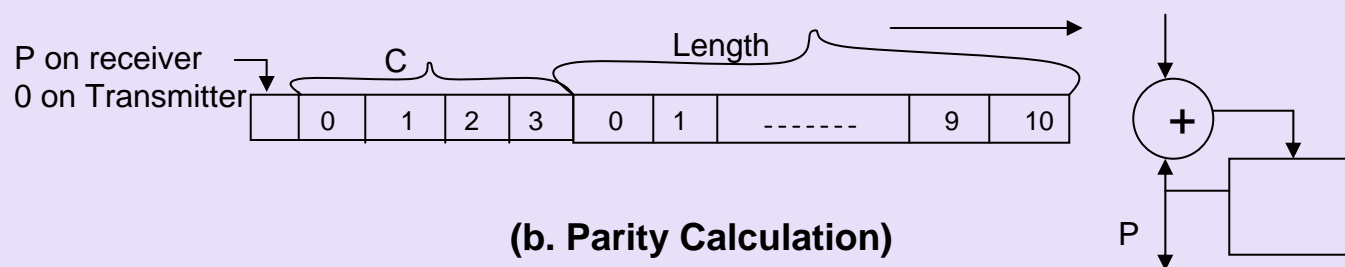
(TLP Layout)

[Len[10:0]: length of the TLP in DWs, Check[4:0]: Check Bits, 18:4, No END]

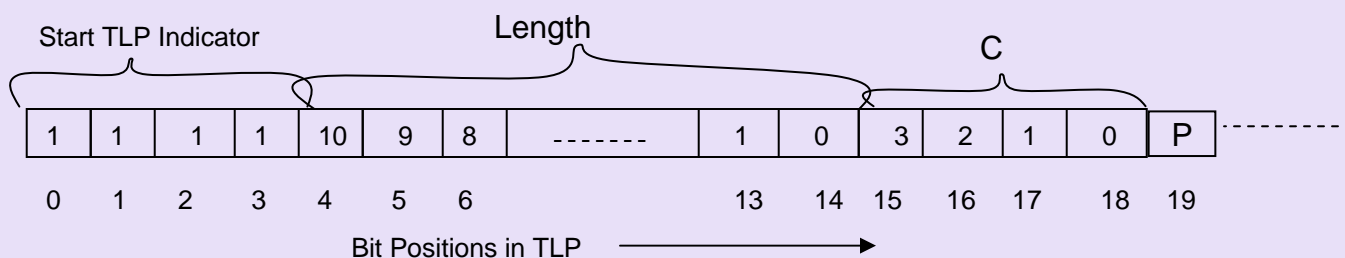
TLP Length Protection



(a. CRC Check Bit Calculation with x^4+x+1)



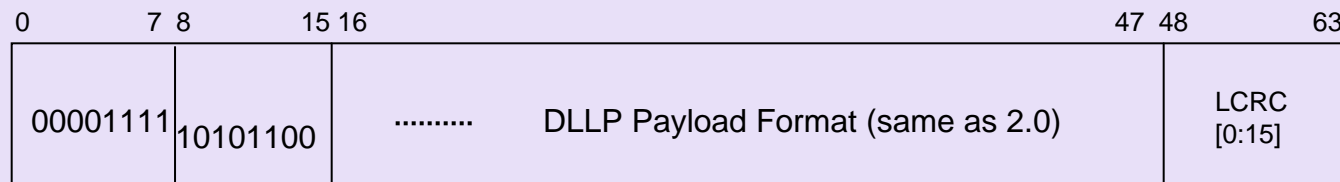
(b. Parity Calculation)



(c. Layout of bits in Packet)

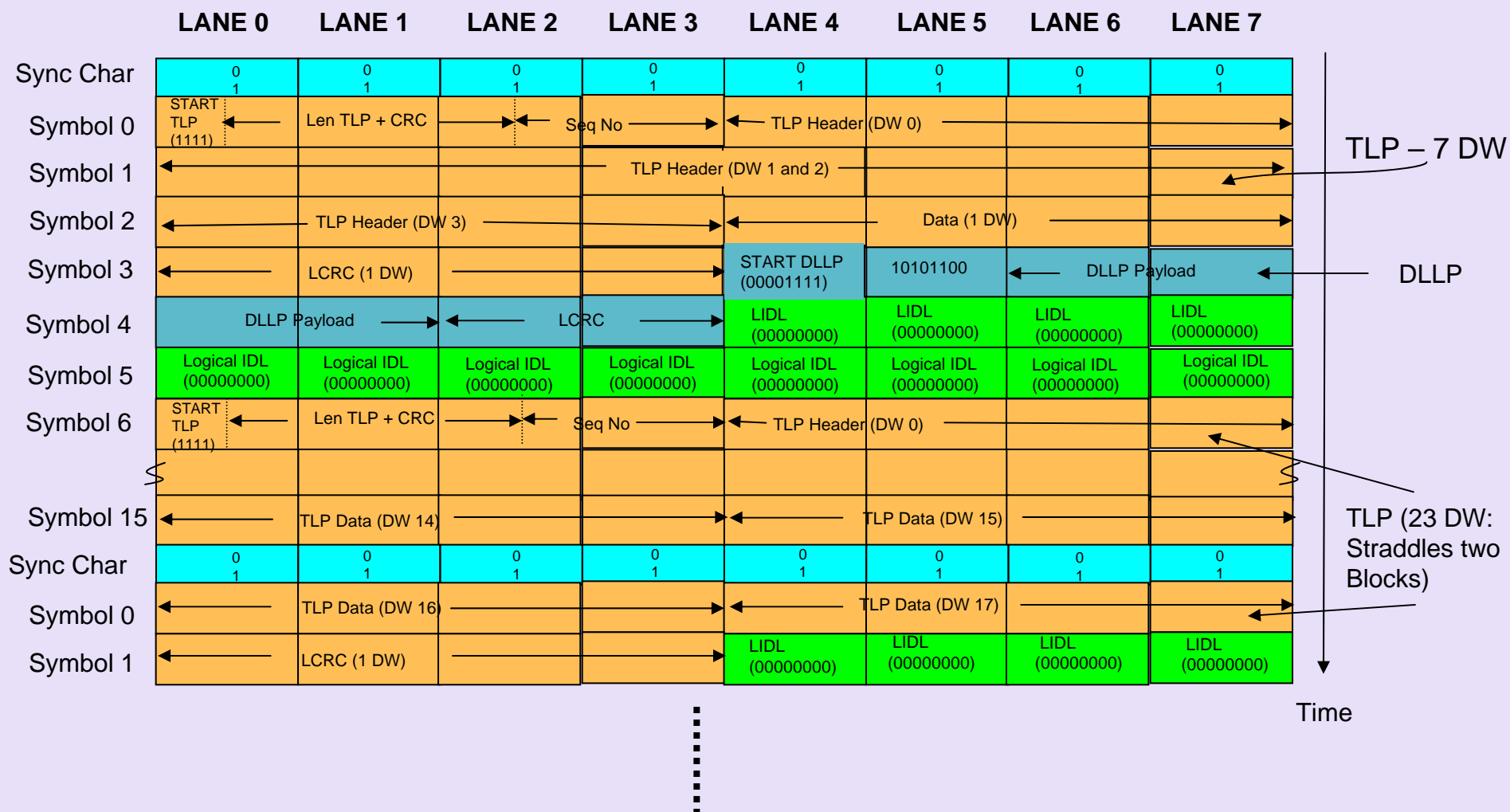
P-Layer Encapsulation: DLLP

- Preserve DLLP layout of 2.0 spec
- First Symbol is 0000_1111
- Second Symbol is ACh
 - ✓ Will allow to share encoding with some Ordered Sets if needed
- Next 4 Symbols (2 through 5) are the DLLP layout
- Next 2 Symbols (6 and 7): LCRC (identical to 2.0)
- No explicit END
- All Symbols are (de)scrambled

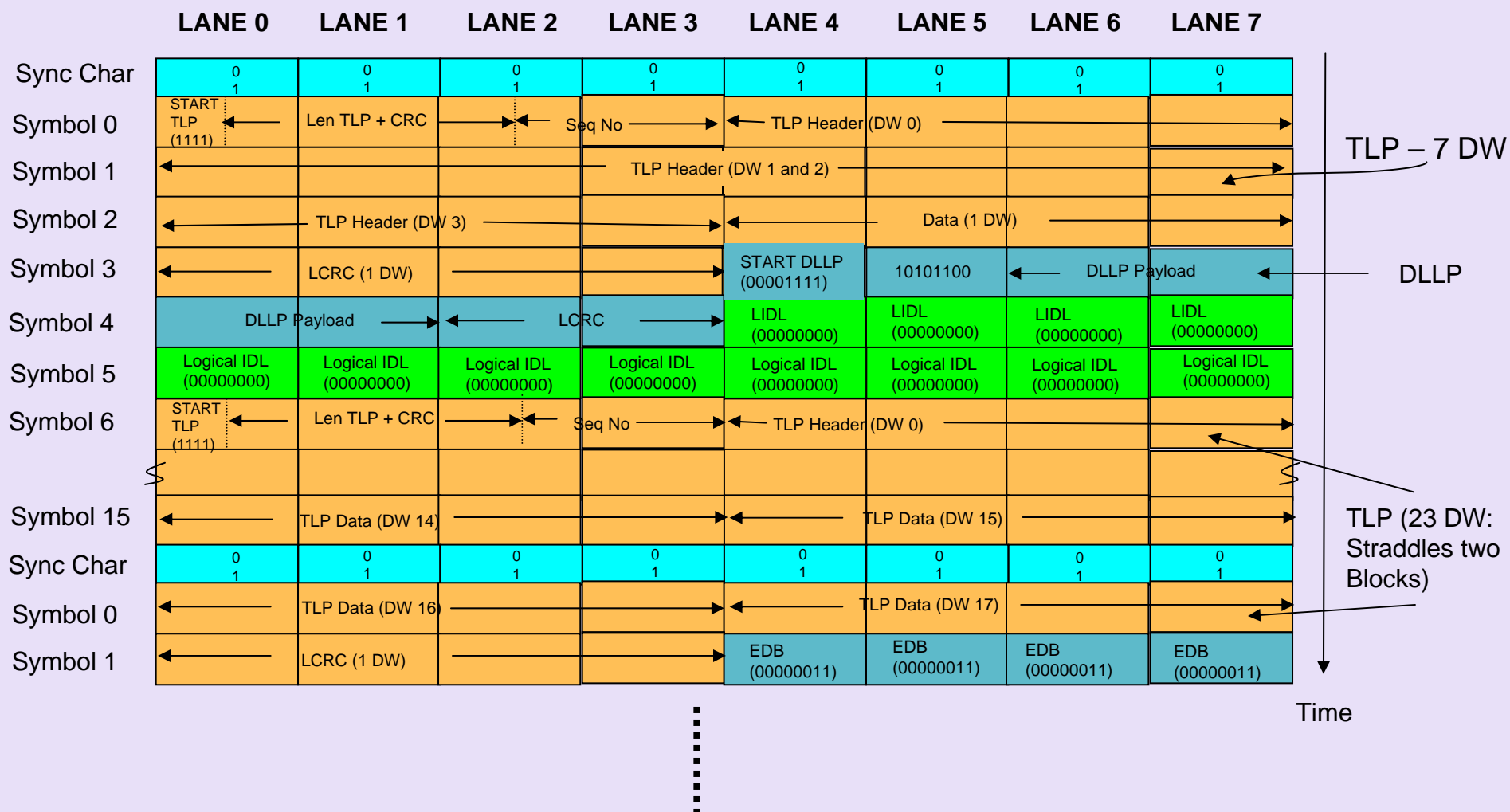


(DLLP Layout)

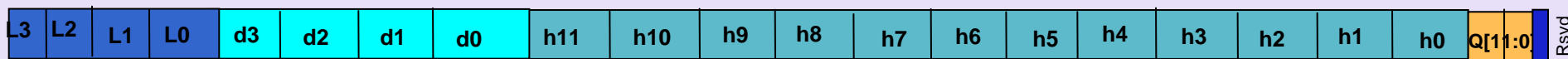
Ex: TLP/ DLLP/ IDLs in x8



TLP/ DLLP/ EDB/ IDLs in x8

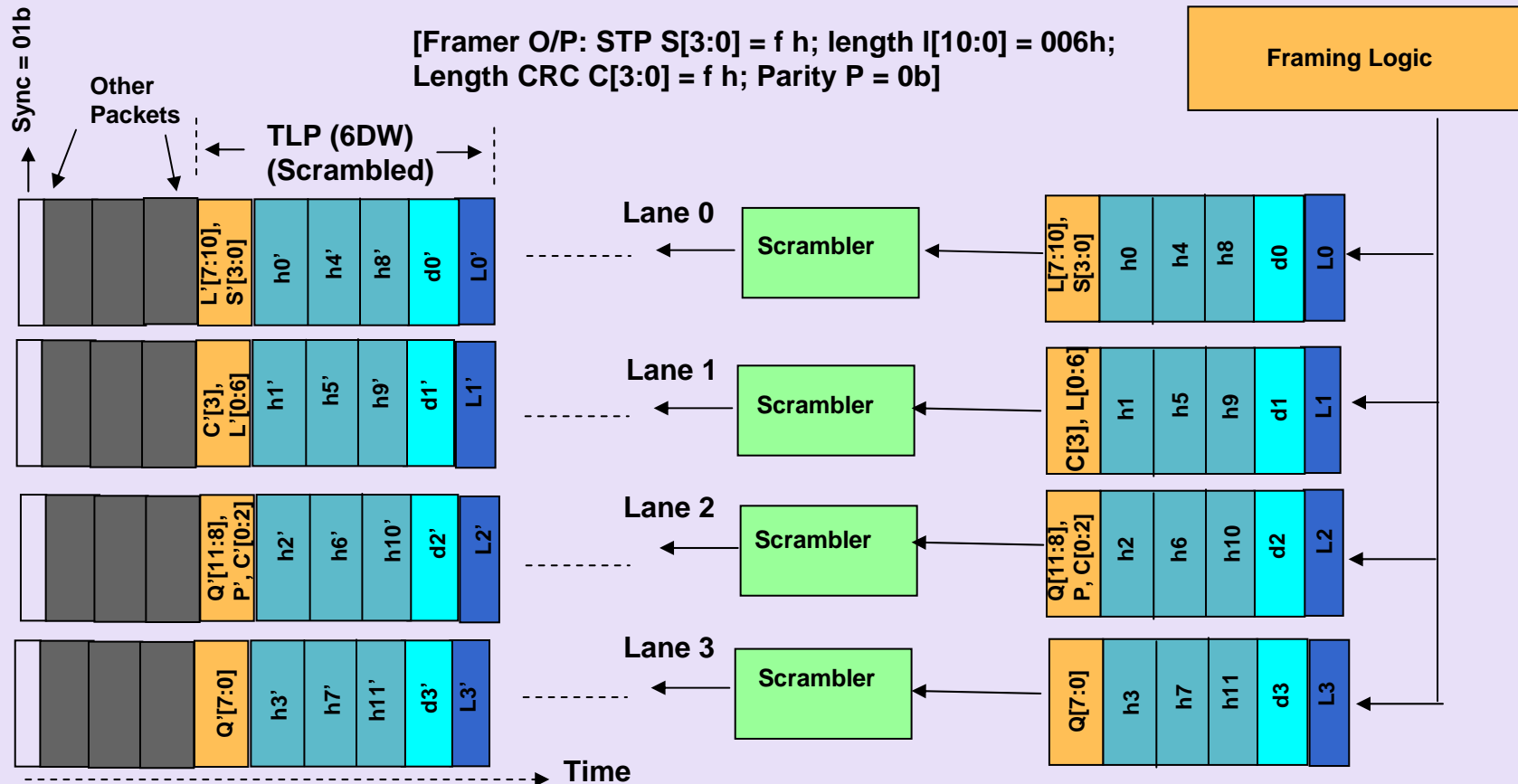


TLP Transmission in a X4 Link

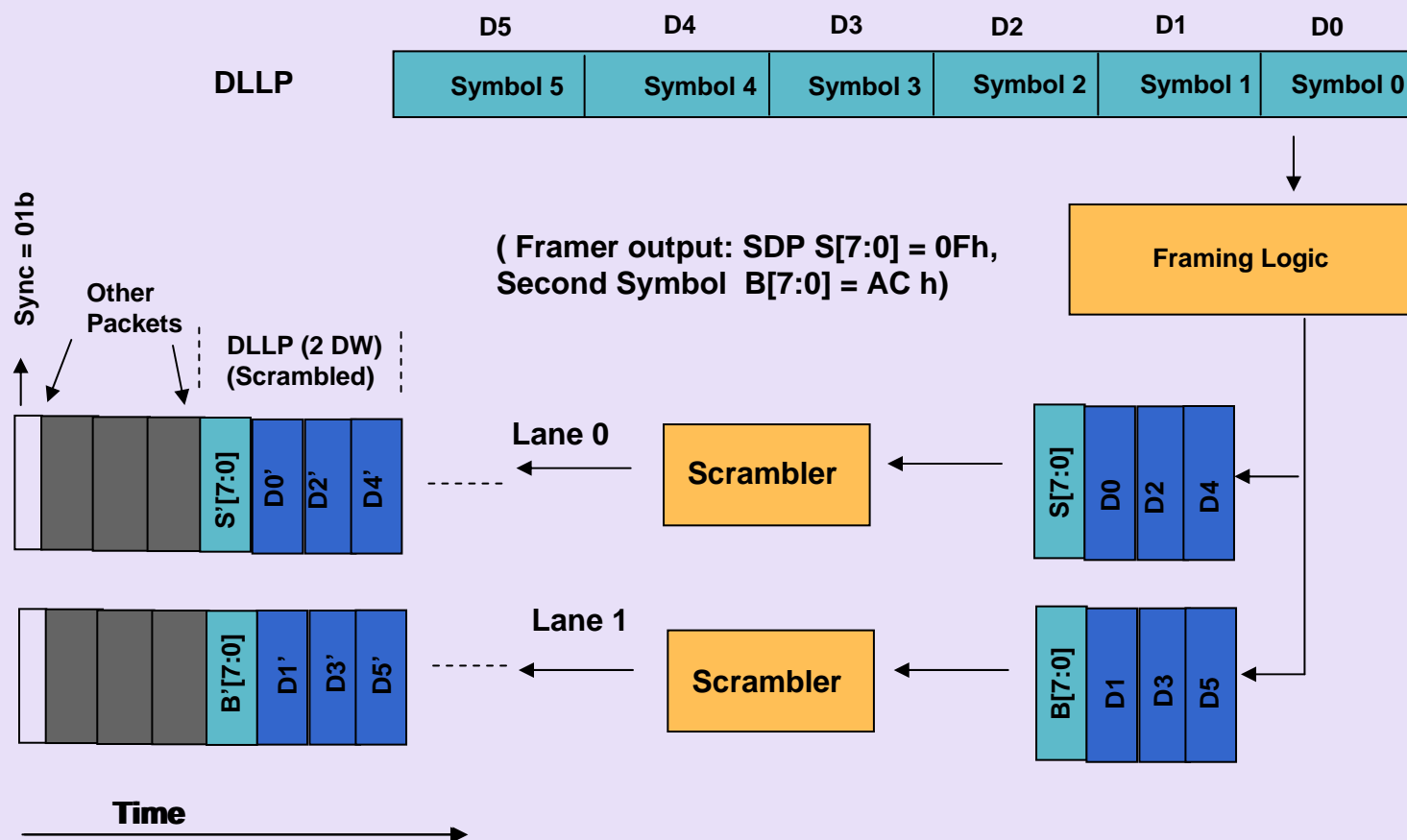


(TLP Transmitted: 3 DW Header (h0 .. h11) + 1 DW Data (d0 .. D3).
1 DW LCRC (L0 .. L3) and Q[11:0]: Sequence No from Link Layer)

[Framer O/P: STP S[3:0] = f h; length I[10:0] = 006h;
Length CRC C[3:0] = f h; Parity P = 0b]



DLLP Transmission in a x2 Link



Error Detection and Recovery

- Framing error
 - ✓ The first Symbol is not one of the three allowed sets
 - ✓ After block lock, if the sync character is not 01 or 10
- Any error requires directing LTSSM to Recovery
 - ✓ CRC error or framing error
 - ✓ Stop processing any received TLP/ DLLP after error until we get through Recovery
 - ✓ Block lock acquired with EIEOS
 - ✓ Scrambler reset with each EIEOS
- Error Detection Guarantees
 - ✓ Triple bit flip detection within each TLP/ DLLP/ IDL/ OS
- Killer Packets: In Recovery.Idle, mandate a variable number of IDL Symbols so that the same TLP retransmitted immediately after Recovery does not come out with the same bit pattern

SKP Ordered Sets

- SKP OS Requirement under discussion
- Multiple options being considered
- Requirements from retiming repeaters, switch, and Logic Analyzer trending towards having an explicit SKP Ordered Set
 - ✓ Variable Block Length
 - ✓ Not Scrambled
 - ✓ Details still being worked on and subject to change
- Loopback Solution would also use the SKP Ordered Set

Agenda

- Problem Statement
- Existing Usage of K-Codes
- Metrics used for evaluation
- Current Direction on Encoding
- **Summary & Call to Action**

Summary & Call to Action

- Encoding scheme decided and development in progress
- Offers advantage of 25% bandwidth for 8GT/s (and above) data rate over 8b/10b encoding
- Rev 0.3 Spec Completed
- Track the spec development and plan for products accordingly

Thank you for attending the
PCIe Developers Conference 2008

For more information please go to
www.pcisig.com