

Burst-Mode Clock and Data Recovery for Optically Interconnected Data Centers

Bhavin J. Shastri, and David V. Plant

Photonic Systems Group, Dept. of Electrical and Computer Eng., McGill University, Montreal, QC H3A 2A7, Canada
shastri@ieee.org

Abstract—We propose a novel burst-mode clock/data recovery (BM-CDR) architecture for optical data center applications. Our design is based on a hybrid topology of a CDR (feedback) and clock phase aligner (feed-forward) utilizing multi-phase clocks.

I. INTRODUCTION

Data centers or large clusters of servers are currently being aggressively deployed in a number of institutions to harness petaflops of computational power and petabytes of storage in a cost-efficient manner [1]. Consequently, there exists a worldwide research interest in designing such large data centers for optimally supporting various applications including scientific computing, financial analysis, data analysis and warehousing, and large-scale network services.

Data centers in general follow a tiered architecture in which network devices (switches or routers) are organized into two or three layers. The highest layer—core tier—is at the root of the tree, whereas the lowest layer—edge tier—is at the leaves of the tree. Between these layers, an aggregation tier may exist when the number of devices is large. The need for highly-specialized ASICs is undeniable [2], with clock and data recovery (CDR) being a critical function in backplane routing and chip-to-chip interconnects. The data received on the aggregation and edge node links is inherently bursty [3] with asynchronous phase steps $j - j - 2$ rad, that exist between the consecutive k^{th} and $(k + 1)^{\text{th}}$ packet. This inevitably causes conventional CDR circuits to lose pattern synchronization leading to packet loss. Preamble bits can be inserted at the beginning of each packet to allow the CDR feedback loop enough time to settle down and thus acquire lock. However, the use of a preamble introduces overhead, reducing the effective throughput and increasing delay. Consequently, to deal with bursty data, these nodes require a burst-mode CDR (BM-CDR). The most important characteristic of the BM-CDR is its phase acquisition time which must be as short as possible. In this paper, we present a novel BM-CDR architecture based on a hybrid topology; that is, a combination of feedback and feed-forward.

II. NOVEL BM-CDR ARCHITECTURE

A block diagram of the proposed BM-CDR is shown in Fig. 1. The BM-CDR is composed of a phase-tracking CDR and a clock phase aligner (CPA). The CDR senses data D_{in} , and generates a synchronized clock CK , with a voltage-controlled oscillator (VCO) in a phase-locked (feedback) loop (PLL). The phase and frequency of CK is compared to D_{in} in

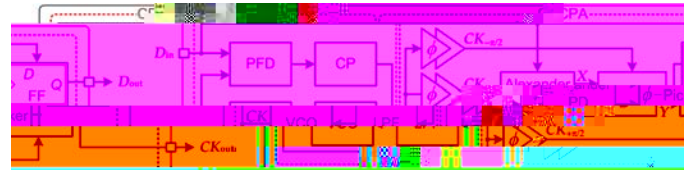


Fig. 1. at the frequency of interest.

Burst-mode functionality is obtained with the CPA which utilizes multi-phase clocks and a phase picking algorithm based on an “early-late” detection principle. This CPA is based on a feed-forward topology, and comprises of phase ($-$) shifters, an Alexander PD, a $-$ -picker, and a D flip (D-FF). The $-$ -shifters utilize the clock recovered from the CDR CK , to provide multiple clocks: CK_0 , CK_+ and $CK_{+ =2}$, with low skew and different phases: $- =2$ rad, and $+ =2$ rad, respectively, with respect to CK_0 . Next, an Alexander PD [4] which inherently exhibits *bang* (binary) characteristics is used to strobe the data form D_{in} , with consecutive clock CK_0 edges, at multiple points in the vicinity of expected transitions [see Fig. 2(a)], resulting in three data samples: previous bit A , current bit B , and a sample of the current bit at the zero crossing T . Depending on the phase difference between the consecutive clock edges, the PD aided by these samples, $X = T - A$ and $Y = A - T$, can determine the location of the clock edge with respect to the data edge as follows: (a) if $A \neq T = B$ ($X \neq 0, Y = 0$) CK_0 lags D_{in} —is late—when $0 < X < + =2$ rad [see Fig. 2(b)]; (b) if $A = T \neq B$ ($X = 0, Y \neq 0$) CK_0 leads D_{in} —is early—when $0 < Y < + =2$ rad [see Fig. 2(c)]; (c) if $A = T = B$ ($X = 0, Y = 0$) CK_0 is in phase with D_{in} .

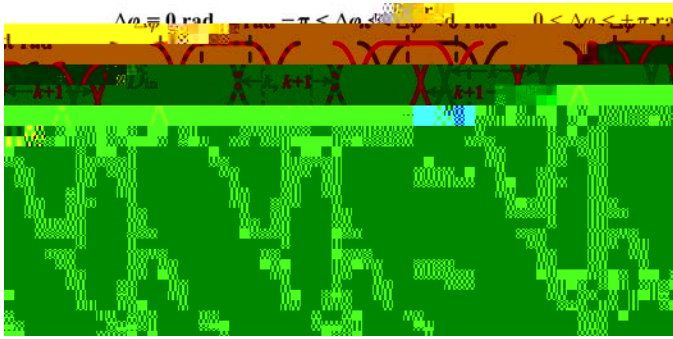


Fig. 3. CPA phase picking algorithm.

CK_o early-late information (X and Y) together with the two multi-phase clocks, CK_{-2} and CK_{+2} , is provided to the ϕ -picker. The idea then behind the phase picking algorithm is depicted with the aid of eye diagrams in Fig. 3. When there is no phase difference between the consecutive packets, $\Delta\phi = 0$ rad, either of the clocks, CK_{-2} and CK_{+2} , will correctly sample the data bits of the phase shifted $(k+1)^{\text{th}}$ packet [see Fig. 3(a)]. This is also true for an antiphase step $\Delta\phi = \pi$ rad—not shown as this is a modulo- 2π process. For a phase step $-\pi < \Delta\phi < 0$ rad, clock CK_{+2} will sample the bits on or close to the transitions of the data eye, whereas clock CK_{-2} will correctly sample the data [see Fig. 3(b)]. Similarly for a phase step $0 < \Delta\phi < +\pi$ rad, clock CK_{-2} will sample the bits on or close to the transitions, whereas clock CK_{+2} will correctly sample the data [see Fig. 3(c)]. That is, regardless of any phase step, there will be at least one clock, either CK_{-2} or CK_{+2} , that will yield an accurate sample. The ϕ -picker then selects the most accurate clock CK_{out} , from these two possibilities for driving the D-FF to retime the data; that is, sample the noisy data, yielding an output D_{out} with less jitter. The foregoing concepts on the Alexander PD and the ϕ -picker are summarized in Table I, leading to the circuit topology in Fig. 4.

III. HARDWARE IMPLEMENTATION

The BM-CDR is being implemented for operation at 10 Gb/s. The main building blocks include a CDR from Centellax (Part #TR1C1-A) and a CPA built by integrating individual chips from Hittite Microwave on a custom designed printed circuit board (PCB). More specifically, the PCB is populated with three 4-bit digital ϕ -shifters (Part #HMC543), an Alexander PD comprised of four D-FFs (Part #HMC673LC3C) and two XOR gates (Part #HMC671LC3C), and a ϕ -picker comprised of an AND gate (Part #HMC672LC3C) and a 2:1