

A 103.125-Gb/s Reverse Gearbox IC in 40-nm CMOS for Supporting Legacy 10- and 40-GbE Links

Taejun Yoon, Joon-Yeong Lee, Jinhee Lee, Kwangseok Han, *Member, IEEE*, Jeong-Sup Lee, Sangeun Lee, Taeho Kim, Jinho Han, Hyosup Won, Jinho Park, and Hyeon-Min Bae, *Member, IEEE*

Abstract—This paper presents the first 103.125-Gb/s multilink gearbox (MLG) IC, which facilitates the transport of independent 10- and 40-GbE signals to 4×25.78 Gb/s physical layers, such as 100GBASE-xR4. The IC consumes only 1.37 W while implementing complicated reverse gearbox functionality. The measured TX jitter from 10- and 25-G lanes is 0.407 and 0.448 psrms, respectively. The measured input sensitivities for a BER of 10^{-12} of the 10- and 25-G RXs are 20 and 42 mVppd, respectively. The proposed gearbox IC, fabricated in a 40-nm CMOS process, occupies 3.7×3.4 mm². The power consumption of RX and TX in a 25-G interface is 50.9 and 52 mW, respectively, and those of a 10-G interface are 29 and 24.4 mW, respectively. MLG functionality is verified using embedded self-test logics.

Index Terms—10 GbE, 40 GbE, 100 GbE, CDR, delay- and phase-locked loop (D/PLL), low power, MLG 2.0, multilink gearbox (MLG), reference less, reverse gearbox IC, transceiver.

I. INTRODUCTION

RECENTLY, mobile, video streaming, and cloud services have triggered explosive growth in data traffic at data centers. The data traffic at data centers is expected to dominate all internet traffic. [1]. This trend strongly drives the demand for high-speed-and-power-efficient data center networks.

To satisfy the demand for higher bandwidth networks, the switch-to-switch connections in data centers are being upgraded from 10 to 100 GbE. The IEEE 802.3ba standardized the transportation of a 100-GbE frame over a 4×25 Gb/s physical layer [2]. In the initial adoption stage of the 100-GbE network, the electrical-side interface was based on 10-Gb/s signaling and the optical-side interface was based on 25-Gb/s signaling; 25 Gb/s is not an integer multiple of 10 Gb/s. Thus, a gearbox IC that provides 10:4 MUX and 4:10 DEMUX functionality was introduced to allow a communication between 10- and 25-Gb/s interfaces. However, it is expected that the

Manuscript received May 12, 2016; revised September 22, 2016 and November 19, 2016; accepted November 25, 2016. Date of publication January 9, 2017; date of current version March 3, 2017. This paper was approved by Associate Editor Anthony Chan Carusone. This work was supported by the National Research Foundation of Korea grant funded by the Korea Government Under Grant 2010-0028680.

T. Yoon, J.-Y. Lee, J. Han, H. Won, and H.-M. Bae are with the School of Electrical Engineering, Korea Advanced Institute of Science and Technology, Daejeon, 34141, South Korea.

J. Lee, K. Han, J.-S. Lee, S. Lee, T. Kim, and J. Park are with TeraSquare Inc., Seoul 06220, South Korea.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/JSSC.2016.2636858

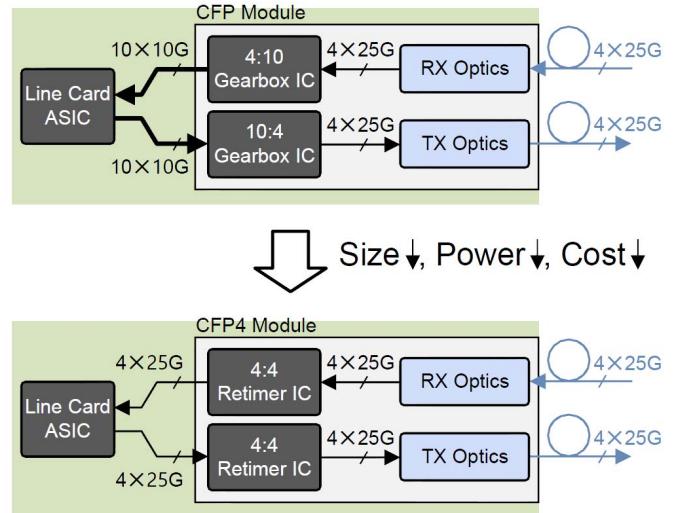


Fig. 1. Evolution of the network connection within a data center.

optical module will eventually be directly connected to the switch ASIC through retimed 25-Gb/s electrical interfaces to save power and area [3], [4]. Fig. 1 shows the progression of a 100-GbE optical module interface.

Ever since client side interfaces were upgraded to 25 G, the demand for new methods offering seamless connections of legacy 10- and 40-GbE links, originating from independent service providers has increased to achieve more efficient management of network resources and cost-effective upgrading in data centers. This demand has led to the new implementation agreement of a reverse gearbox IC, referred to as a multilink gearbox (MLG) [5], [6]. Fig. 2 shows the differences among three types of gearbox ICs. For the conventional gearbox IC, plesiochronous 10-GbE links (or a combination of 10- and 40-GbE links) cannot be aggregated into 4×25 Gb/s streams of one common clock domain, since it does not support the rate conversion process. Independent Ethernet traffics can have the frequency offsets by as much as ± 100 ppm, because the clock signals are generated locally. However, the proposed MLG 2.0 IC supports any combination of independent 10- and 40-GbE links thanks to the embedded rate conversion process. The MLG logic block enables the mapping of plesiochronous 10-GbE/40-GbE streams to an MLG stream of

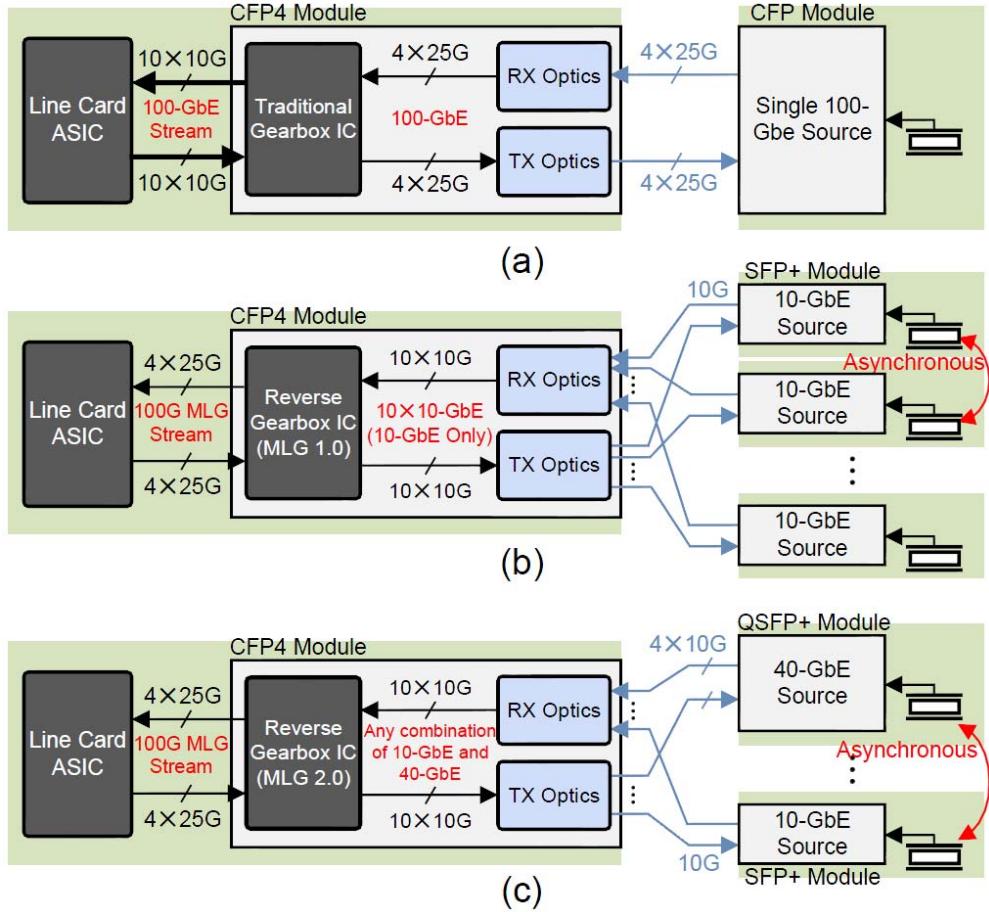


Fig. 2. Comparison between the conventional and the reverse gearbox IC. (a) Conventional gearbox. (b) MLG 1.0. (c) MLG 2.0.

a common clock domain while preserving the lane and bit ordering of individual 10-G lanes, and vice versa [5]. Note that MLG 1.0 supports only up to ten independent 10-GbE links.

Fig. 3 shows two application examples of the MLG 2.0 IC: a port expander and a virtual link [6]. The MLG 2.0 IC provides a seamless interface between 10-Gb/s links and a switch ASIC with 25-Gb/s I/Os. As a result, the number of 10-Gb/s ports connected to a single 25-Gb/s I/O pin increases by 2.5 times as compared with that connected to a single 10-Gb/s I/O in a switch ASIC. In the virtual link application, the reverse gearbox IC enables fiber-efficient data center interconnection by aggregating multiple 10- and 40-GbE links and transmitting over 25-Gb/s links. This functionality reduces the capital/operation expenses in a data center by reducing the number of fiber/optical components needed to satisfy the bandwidth requirements.

The primary technical challenge for implementing the reverse gearbox IC is that multiple 10- and 40-GbE links from independent networks operate independently, with slightly different data rates up to ± 100 ppm. To support this plesiochronous operation, the reverse gearbox IC should integrate complex digital logic functions, such as physical coding sublayer (PCS), typical in network ASICs, together with high-speed mixed mode SerDes. The integration of large digital logic blocks with high-speed multiple transceivers

causes a considerable inefficiency in terms of power and area due to complicated high-speed clock routing across the synthesized digital logic block. Yet another physical layer design challenge is that entire parallel 10- and 25-G CDRs should operate independently without a reference clock signal. It is because 10-G links can be originated from independent service providers and thus the number of activated links can vary in time. VCO-based parallel CDRs are inefficient in terms of area, and phase interpolator (PI)-based CDRs do not operate independently in the absence of a reference clock signal. The aforementioned design challenge escalates due to low-power requirement, since such complex IC should be integrated in a minuscule module form factor, such as CFP4 or QSFP28.

In this paper, we propose a low-power, reference-less 103.125-Gb/s reverse gearbox IC satisfying Optical Inter-networking Forum (OIF) MLG 2.0 standards. The reverse gearbox IC includes: 1) 10x10 G and 4x25 G transceivers operating independently without an external reference clock signal and 2) a synthesized logic block enabling transport of multiple 10- and 40-GbE data streams across 4x25 G lanes. This IC is fabricated in a 40-nm CMOS process and consumes 1.37 W when ten 10-G and four 25-G transceivers operate at 10.31 and 25.78 Gb/s, respectively, in normal MLG mode operation.

Section II describes the architecture of the proposed reverse gearbox IC and provides a detailed description of

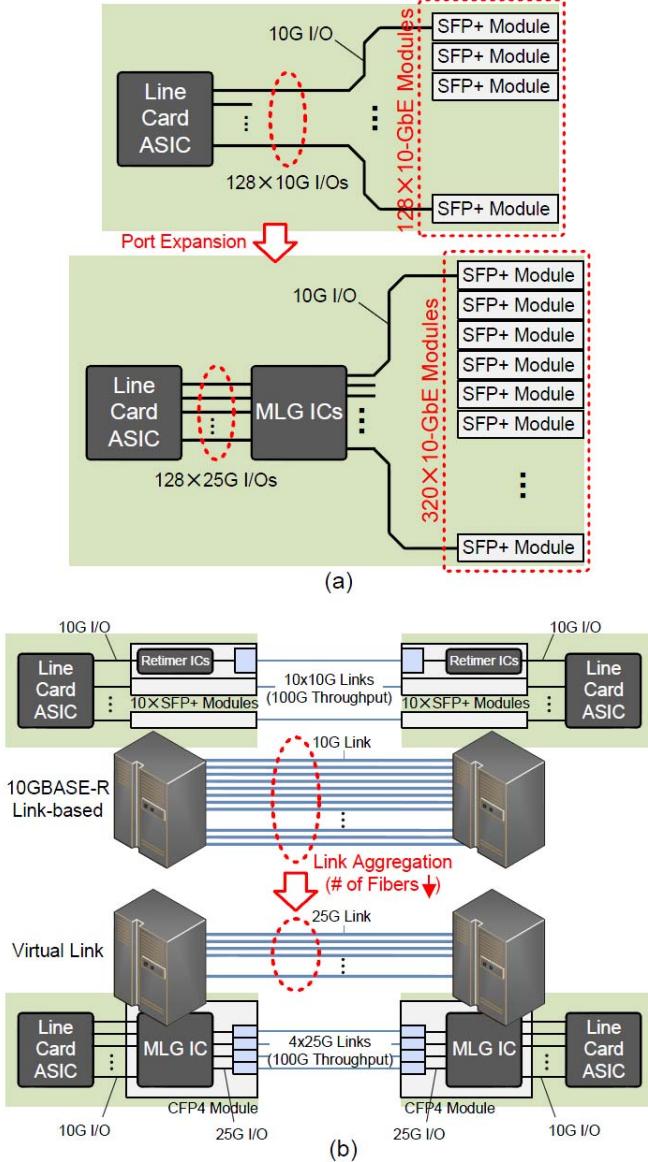


Fig. 3. Application examples of the reverse gearbox IC. (a) Port expander [75] {R3-2}. (b) Virtual link.

the key components of the 10- and 25-G transceivers. Section III describes the logic operations of the synthesized logic blocks. Section IV provides measurement results and comparisons with the recent studies. Finally, Section V concludes this paper.

II. 103.125-Gb/s REVERSE GEARBOX IC ARCHITECTURE

Fig. 4 shows the overall architecture of the proposed reverse gearbox IC and indicates the direction of signal flows in the MLG operation mode. The IC consists of a 10-G interface, including ten parallel 10-G transceivers and a 25-G interface, including four parallel 25-G transceivers, two clock generators, and an MLG 2.0 logic core. The MLG logic core is synthesized and implements an MLG MUX and DEMUX. The MLG MUX multiplexes ten plesiochronous 10.31-Gb/s streams into four synchronous 25.78-Gb/s streams; the MLG DEMUX reverses the operation of the MLG MUX. The control signals for each

10- and 25-G domain TX PI are provided by the corresponding RX (see Fig. 4). Note that the 10- and 25-G clock domains are not synchronized in general and thus the MLG 2.0 logic core should manage the dynamic frequency offset for error-free operation. This issue is discussed in more detail in Section IV. The 10-GbE frame generators and checkers are implemented in the MLG 2.0 logic core to verify the MLG logic functionality using an external loopback. The proposed reverse gearbox IC can also be configured as a conventional gearbox.

A PI-based reference-less CDR in each lane operates independently while sharing a differential clock signal from a VCO [7]. The PI-based CDR architecture is chosen, because it is better suited than its VCO-based counterparts for parallel transceiver design, because of its power-and-area efficiency and robustness to interference. To achieve reference-less frequency acquisition while minimizing power overhead, each RX lane employs a stochastic reference clock generator (SRCG) with a jitter suppression loop (JSL) [9]–[11] (see Fig. 5). Quasi-periodic reference clock signals generated by the SRCGs are fed to frequency detectors (FDs) in the common clock generator. The up/down signals from FDs are aggregated in the digital domain to control the VCO, such that the VCO is frequency locked to the average frequency of the incoming signals. Thus, this frequency acquisition scheme achieves lane-independent reference-less clock acquisition in the presence of incoming data in any of the RX lanes [7]. The 25-G clock generator incorporates two VCOs to support not only 4×25.78 Gb/s transmission but also 4×20.63 Gb/s transmission that maps two 40GBASE-R signals over four 20.63-Gb/s physical lanes [6]; 12.89-GHz (or 10.31 GHz) half-rate differential global clock signals are distributed from the 25-G clock generator to each 25-G transceiver via an on-chip transmission line (T-line) with a multidrop scheme; quarter-rate clock signals are locally generated. As compared with a four-phase clock distribution scheme, the proposed two-phase distribution scheme, i.e., differential clock distribution, does not suffer from I/Q phase mismatch and saves a significant amount of power, because only a single differential clock generator and a differential distribution buffer are used for all RX and TX lanes in the 25-G interface. For the 10-G interface, each transceiver receives a full-rate differential 10.31-GHz clock signal from the 10-G clock generator and locally generates half-rate clock signals, because the half-rate scheme suffices for the adoption of CMOS logic gates.

A. All-Digital Open-Loop Controlled D/PLL-Based Transceiver Architecture

The 10- and 25-G transceivers employ an all-digital delay and phase-locked loop (D/PLL) to provide decoupled RX and TX bandwidths: The RX has a wide jitter tracking bandwidth and the TX has a narrow jitter transfer bandwidth. The proposed D/PLL scheme employs an open-loop controlled PI (see Fig. 6) for efficient implementation of the reverse gearbox IC that integrates multiple parallel transceivers and a complex digital logic core. Fig. 7 shows the D/PLL-based transceiver architecture proposed in [7]. The D/PLL in [7] filters out the jitter from the input data stream without substantial

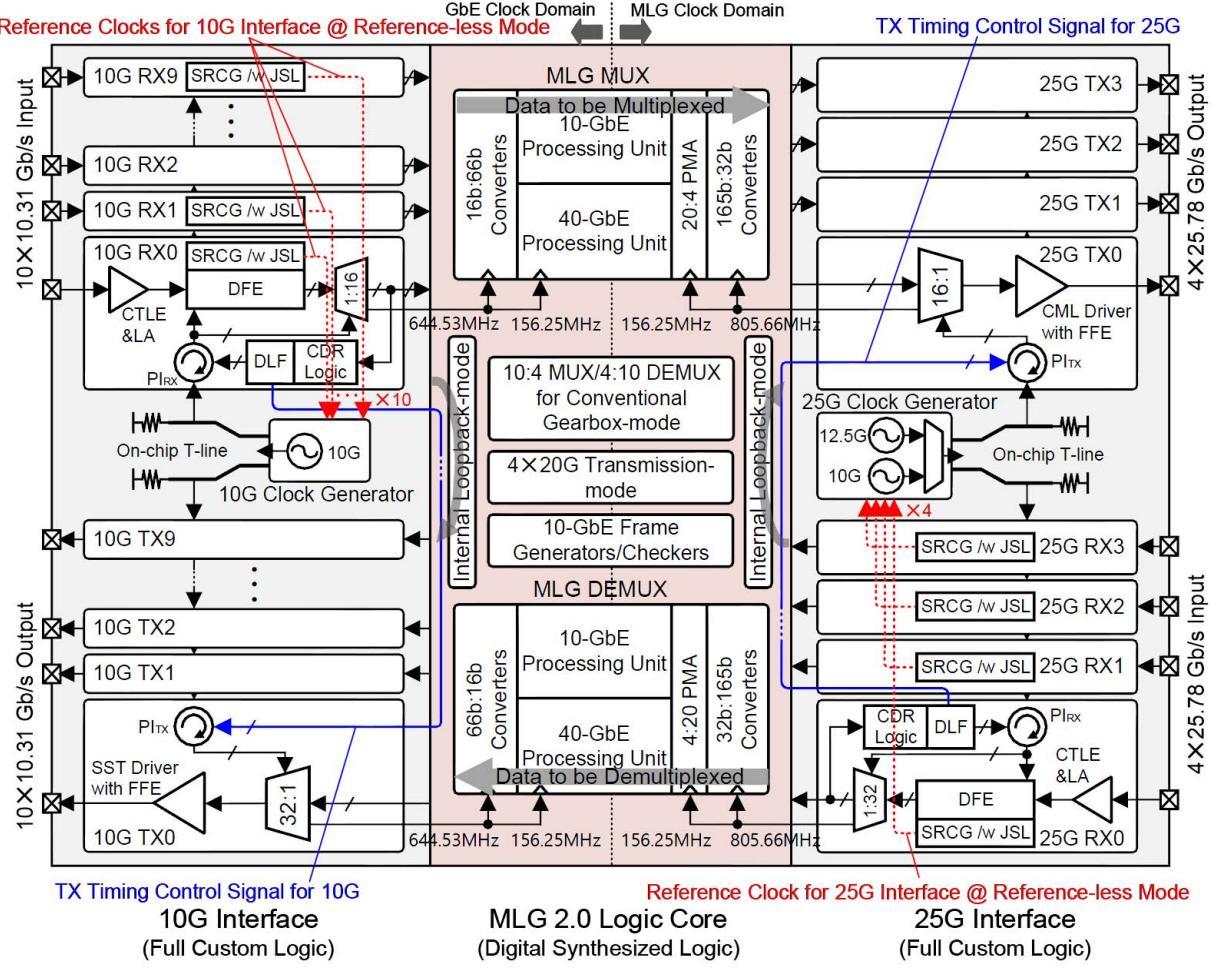


Fig. 4. Overall architecture of the reverse gearbox IC and signal flow in the MLG operation mode.

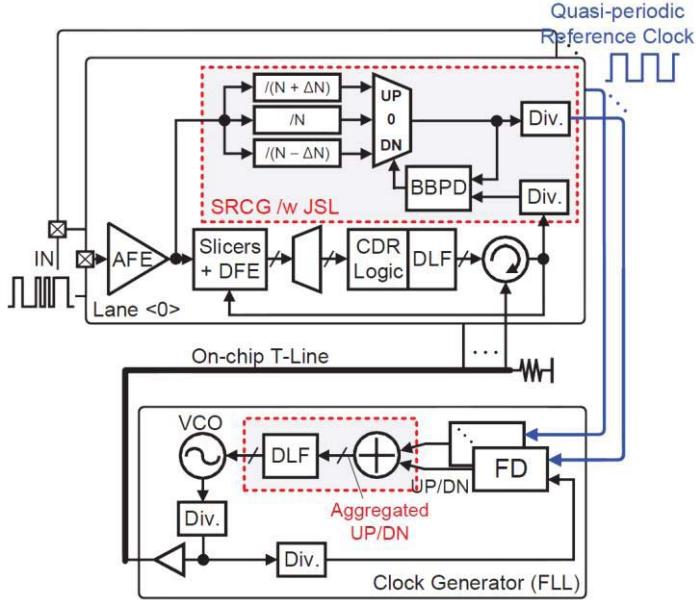


Fig. 5. Block diagram of the lane-independent reference-less clock acquisition scheme.

power penalty unlike the conventional D/PLL-based transceiver [12], because a power-hungry VCO and line-rate delay cells are substituted with two cascaded closed-loop controlled

low-power PIs. However, because both RX and TX PIs are placed at the RX side, 7-GHz four-phase TX clock signals should be distributed to the TX side across the digital

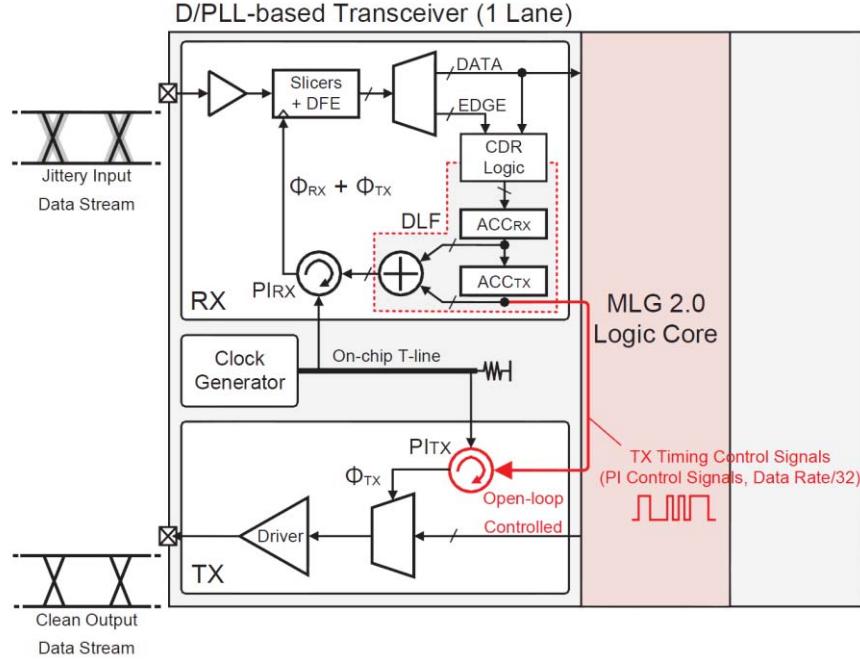


Fig. 6. Proposed all-digital D/PLL-based transceiver with an open-loop controlled PI.

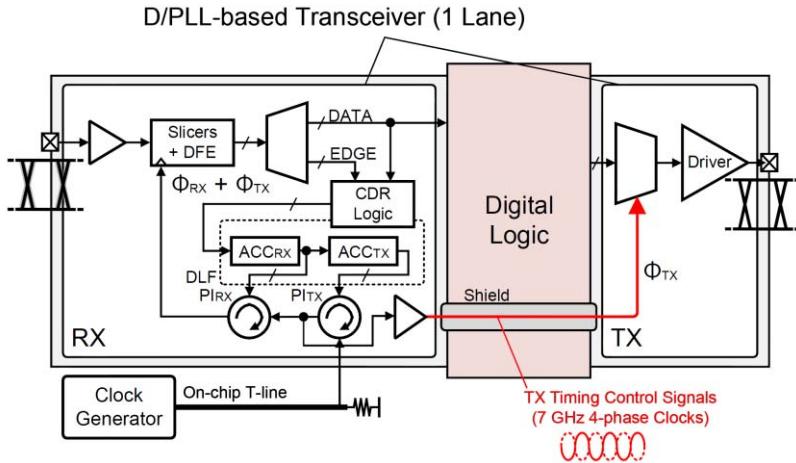


Fig. 7. All-digital D/PLL-based transceiver proposed in [7] and [8].

logic area. This kind of clock signal routing complicates automated physical design, i.e., floorplanning, placement, and routing, of the digital logic core, because a grounded shield layer under the entire clock path is necessary for signal integrity. In addition, the physical distance between the RX and the TX causes considerable power consumption in the clock drivers. Previous works [13]–[15] employing D/PLL-based transceiver architectures suffer from the same power issues, because high-speed clock signals should be routed from the VCO to the phase shifter.

Fig. 8 shows the layout floorplan of the proposed reverse gearbox IC, where the clock distribution path and the digital domain timing control signal paths between RXs and their corresponding TXs are highlighted. The floorplan in Fig. 8 clearly shows that the length of the TX timing control signal

is minimal. The minimum routing distance between the RX and the TX of the reverse gearbox IC is roughly four times longer than that shown in the previous design [7], because the large MLG logic block is surrounded by 10- and 25-G transceivers; this arrangement makes the closed-loop D/PLL scheme impractical because of its excessive power consumption.

To overcome the aforementioned design issues, the proposed D/PLL-based transceiver employs an open-loop controlled PI scheme for power-efficient remote control of the TX clock signal, as shown in Fig. 6. The TX PI is physically located at the TX side and directly receives the VCO clock signal using a multidrop scheme as the RX PI. Then, the TX PI is controlled using a low-frequency TX timing control signal from the digital loop filter. This scheme can be implemented

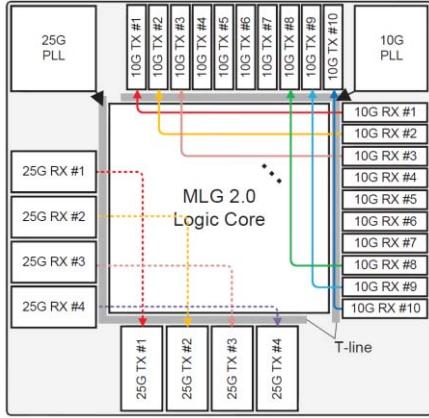


Fig. 8. Layout floor plan of the proposed reverse gearbox IC and TX timing control signal paths.

without distributing the high-frequency clock signals across the digital logic area.

Fig. 9 shows the linearized s-domain models of the D/PLL in the previous [7] and the proposed designs. A bang–bang phase detector is represented by a linearized gain, K_{BBPD} , which is given by [8]

$$K_{BBPD} = \frac{1}{\sqrt{2\pi}\sigma_J} \left[1 + e^{-\frac{1}{2}\left(\frac{\beta\theta}{\sigma_J}\right)^2} \right] \quad (1)$$

where σ_J denotes the standard deviation of the input Gaussian jitter. β and γ represent the integral gain of the RX accumulator, ACC_{RX} , and the TX accumulator, ACC_{TX} , respectively. f_s is the operating frequency of the accumulators. θ is the gain of the PI, which is given by

$$\theta = \frac{1UI}{2R_{interpolator}} \quad (2)$$

where $R_{interpolator}$ is the resolution of the PI. The proposed D/PLL scheme replaces the phase domain addition [7] with a digital domain addition; thus the TX PI is open-loop controlled. As a result, the overall phase-domain transfer function [7] remains unchanged. The closed-loop transfer function $\emptyset_{TX}/\emptyset_{IN}$ is given by

$$\frac{\emptyset_{TX}(s)}{\emptyset_{IN}(s)} = \frac{K_{BBPD}f_s^2\gamma\theta}{s^2 + K_{BBPD}f_s\beta\theta s + K_{BBPD}f_s^2\gamma\theta} \quad (3)$$

which clearly shows that the single closed-loop zero is located at the origin. The relationship between $\emptyset_{TX} + \emptyset_{RX}$ and \emptyset_{IN} is

$$\frac{\emptyset_{TX}(s) + \emptyset_{RX}(s)}{\emptyset_{IN}(s)} = \frac{K_{BBPD}f_s\beta\theta s + K_{BBPD}f_s^2\gamma\theta}{s^2 + K_{BBPD}f_s\beta\theta s + K_{BBPD}f_s^2\gamma\theta} \quad (4)$$

and the relation between $\emptyset_{TX} + \emptyset_{RX}$ and \emptyset_{TX} becomes

$$\frac{\emptyset_{TX}(s)}{\emptyset_{TX}(s) + \emptyset_{RX}(s)} = \frac{1}{1 + \frac{\beta}{f_s\gamma}s}. \quad (5)$$

From (5), it is clear that the TX bandwidth can be made lower than the RX bandwidth by adjusting design parameters β and γ [7]. The proposed open-loop D/PLL design features the key properties of the conventional D/PLL architecture while reducing power consumption and achieving design simplicity.

B. Receiver and Transmitter Implementation

The 25-G interface adopts a quarter-rate clocking scheme. CMOS logic gates are used extensively to save power and area. CML gates are used only for the analog front end (AFE) and 50- Ω drivers. Fig. 10 shows the block diagram of the 25-G RX. The RX includes an analog equalizer, a limiting amplifier (LA), a one-tap loop-unrolled DFE, an SRCG with a JSL, a DEMUX, and a CDR logic block implemented in the digital domain. The analog equalizer provides up to 23 dB of boosting gain to cover not only short-reach application but also medium-reach applications, such as CEI-28G-MR and CAUI-4 C2C (\sim 20-dB loss at 14 GHz). To compensate for the loss up to 20 dB, the analog equalizer and the one-tap loop-unrolled DFE are employed. The LA provides a dc gain of 22.3 dB. The RX uses a quarter-rate edge sampling scheme [20] to reduce power consumption. Internal eye opening monitoring (EOM) is implemented at the output of the LA using an auxiliary sampler with adjustable reference voltage and sampling phase. The optimum sampling phase of each quarter-rate clock signal is determined using EOM.

The 25-G TX in Fig. 11 also adopts a quarter-rate clocking scheme using a single-stage 4:1 MUX for power saving [7]. A phase trimming scheme is employed to eliminate duty cycle distortion and phase misalignment among quarter-rate clock signals. A three-tap FIR filter is incorporated for pre-emphasis up to 7 dB. An inductive peaking scheme is employed at the output stages to enhance the bandwidth. The 10-G RX shares most of the 25-G RX architecture except that it uses a half-rate clocking scheme. The 10-G TX in Fig. 12 employs a voltage-mode output driver scheme to reduce power consumption. Theoretically, the voltage-mode output driver consumes four times less power than its CML-based counterparts at the same output swing condition. The output swing is made variable from 0.25 to 0.9 V_{ppd} by introducing shunting slice units between the two output nodes [16]. The output impedance is made tunable to ensure 50- Ω matching under PVT variations. The 10-G TX supports two-tap pre-emphasis up to 10 dB.

III. MLG 2.0 LOGIC CORE

The MLG 2.0 logic core is the core block of the reverse gearbox IC. The MLG 2.0 logic core enables any combinations of multiple asynchronous 10- and 40-GbE links to be transported to 4×25 G lanes while preserving their frame structure and ordering, without the help of a GbE framer [6]. Thus, the reverse gearbox IC can be used as a port expander or a virtual link, unlike the traditional gearbox IC which provides only MUX and DEMUX functionality. The MLG functionality is implemented using a 100GBASE-R PCS and a physical medium attachment (PMA), which enables the construction of a 100-GbE frame under a variety of physical lane widths [6]. Fig. 13 shows an exemplary block diagram implemented in a line card ASIC for the conversion of a 100-GbE stream into 10×10.31 Gb/s streams [17]. The 64 b/66 b encoder in the PCS encodes the Ethernet MAC data into a continuous stream of 66-b blocks. The encoded stream is divided into twenty PCS lanes of 5.16 Gb/s using a round-robin algorithm.

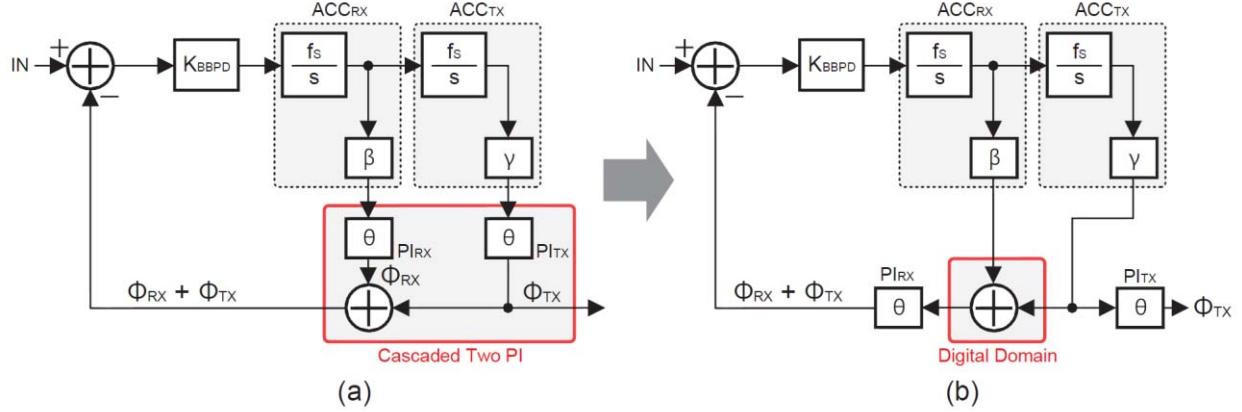


Fig. 9. Linearized s-domain model of (a) D/PLL in [7] and [8] and (b) proposed D/PLL.

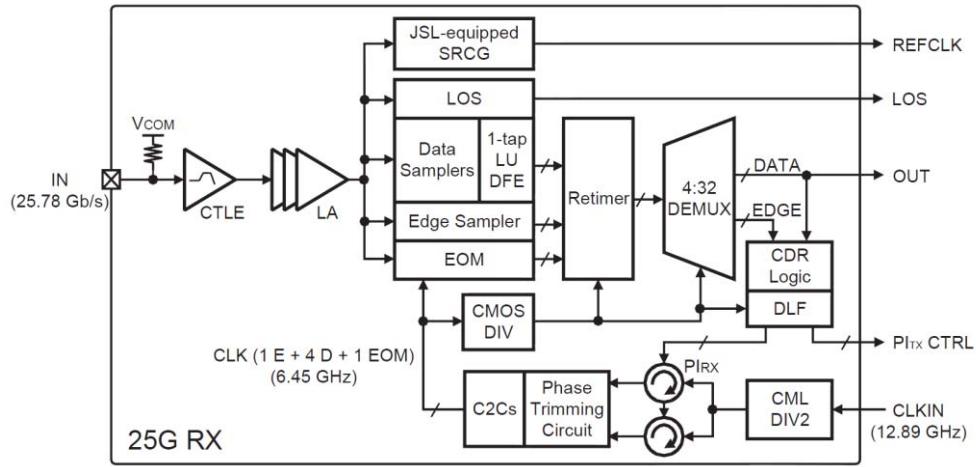


Fig. 10. Block diagram of the 25-G RX.

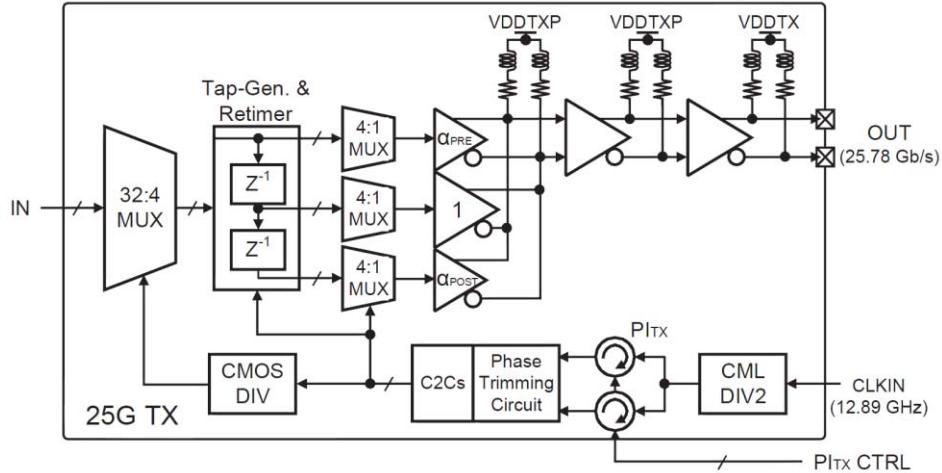


Fig. 11. Block diagram of the 25-G TX.

Then, a special 66-b block referred to as an alignment marker is added once every 16,384 blocks in each lane. The alignment markers are used for identification, deskewing, and reordering of the PCS lanes at the receiver side. Note

that twenty PCS lanes can be mapped into or unmapped from diverse formats of physical layers depending on the subsequent PMA. For 100GBASE-xR4 in Fig. 13, a pair of PCS lanes are multiplexed by the 20:10 PMA into a

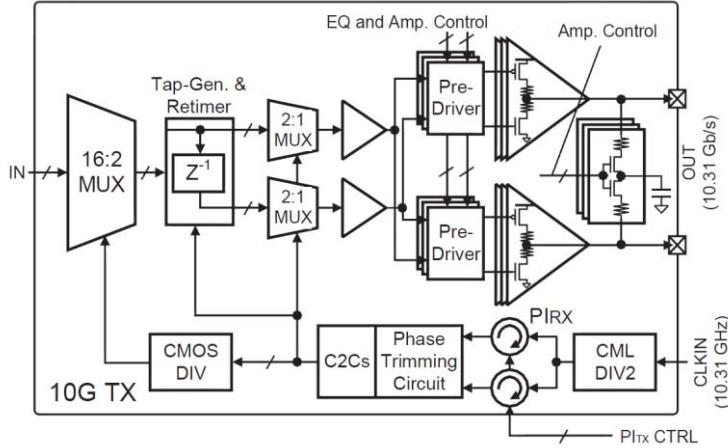


Fig. 12. Block diagram of the 10-G TX with an SST output driver.

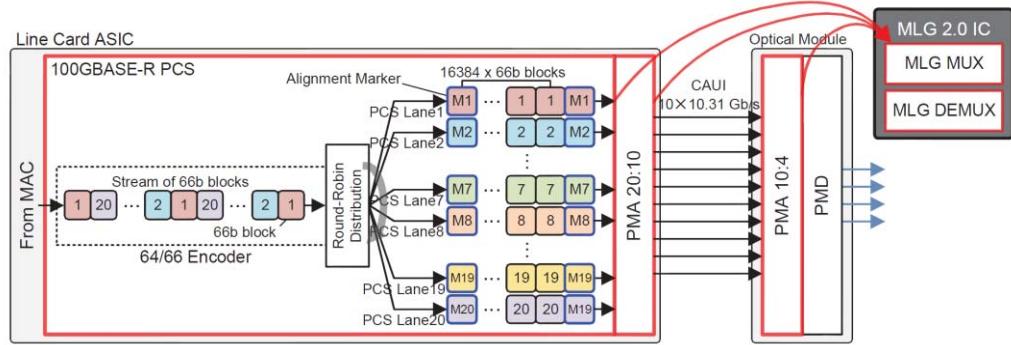


Fig. 13. Exemplary block diagram implemented in a line card ASIC for the conversion of a 100-GbE stream into 10×10.31 Gb/s streams.

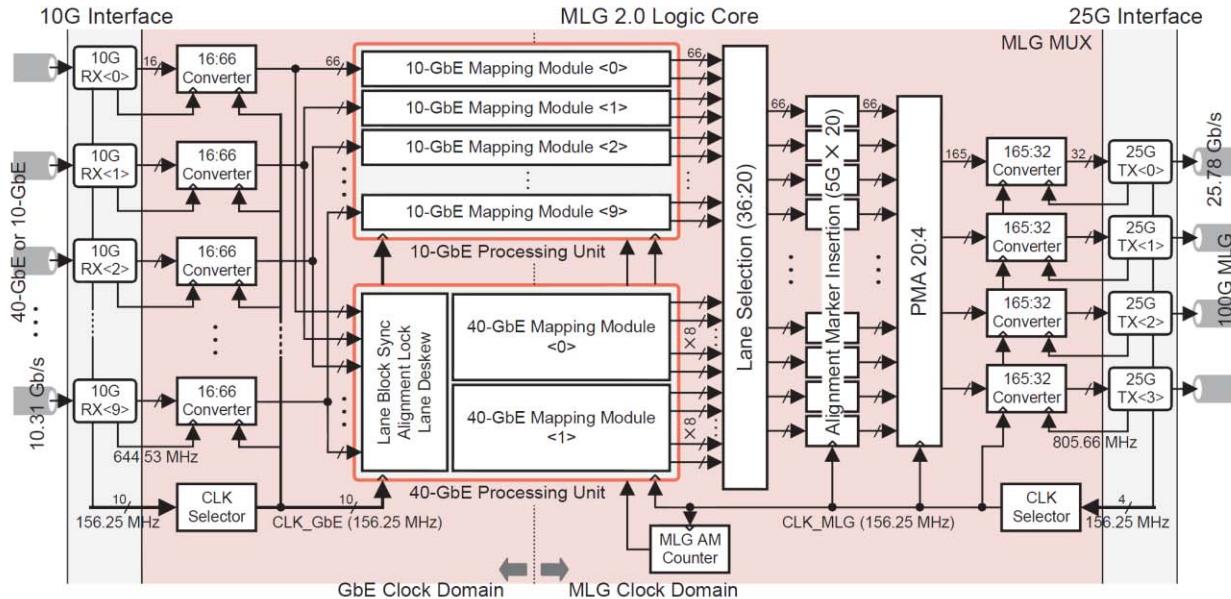


Fig. 14. Block diagram of the implemented MLG MUX.

single lane of 10.31 Gb/s and transmitted to the optical module via a CAUI electrical interface. Then, the 10:4 PMA within the optical module serializes the 10×10.31 Gb/s

electrical lanes into 4×25.78 Gb/s optical lanes. The aforementioned PCS and PMA functions are implemented in the proposed reverse gearbox IC. Fig. 14 presents a block diagram

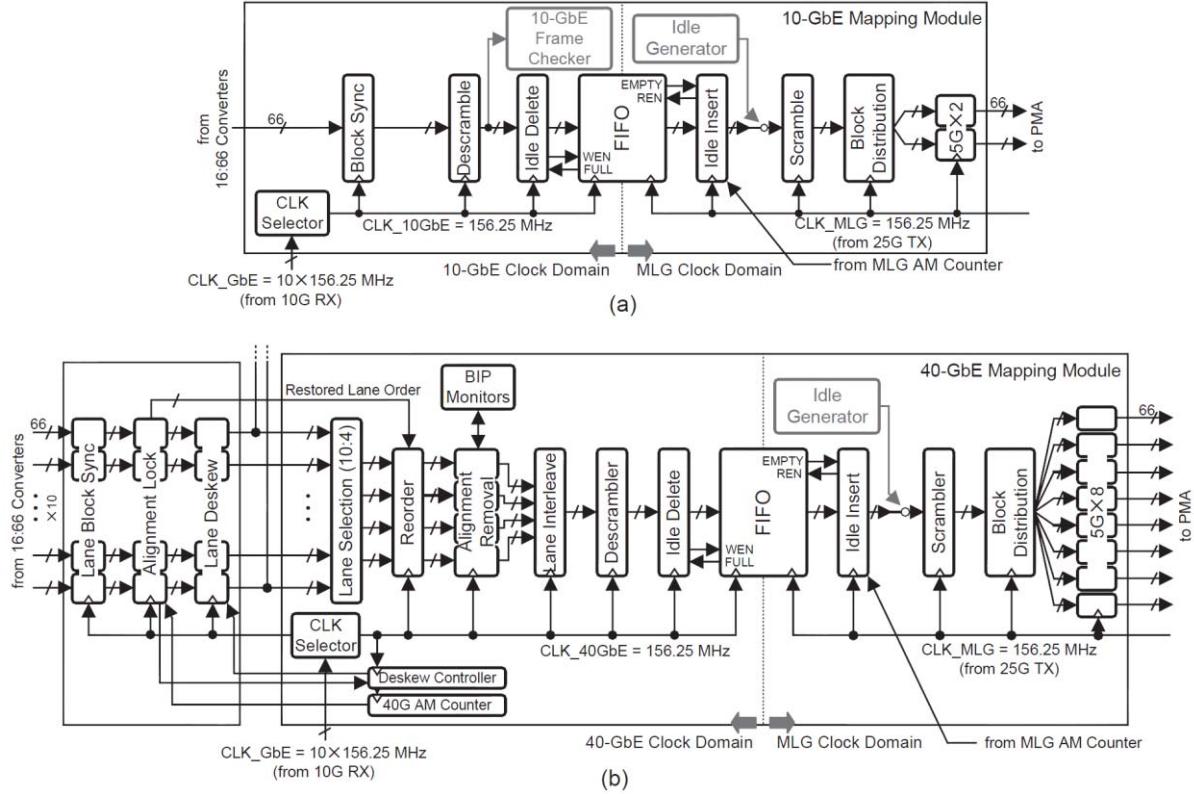


Fig. 15. Ethernet frame mapping modules in the MLG MUX. (a) 10-GbE mapping module. (b) 40-GbE mapping module.

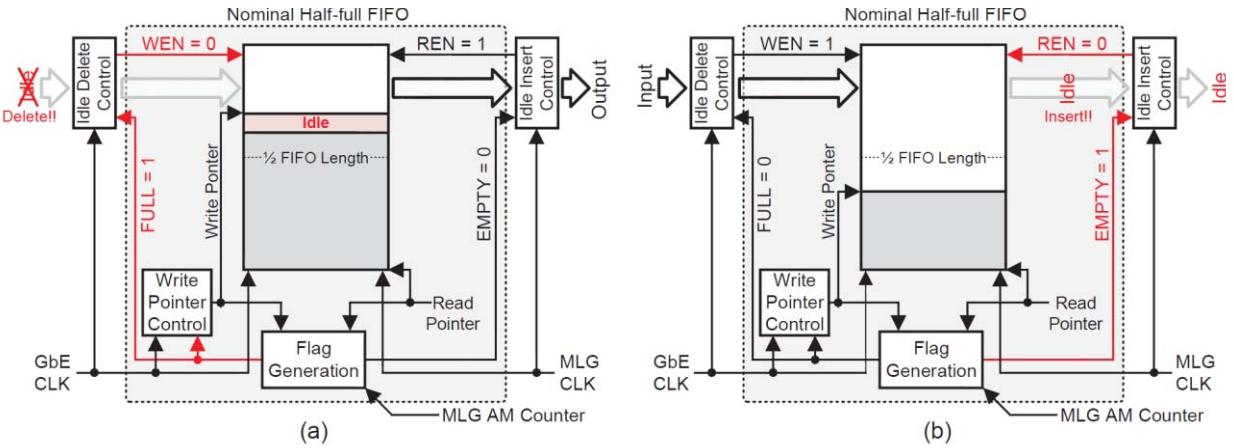


Fig. 16. Operation of the half-full FIFO when (a) GbE clock frequency > MLG clock frequency and (b) GbE clock frequency < MLG clock frequency.

and the signal flow of the MLG MUX. The MLG MUX aggregates recovered 10- and 40-GbE data from 10-G CDRs and converts them into 103.13-Gb/s data composed of 20×5.16 Gb/s MLG lanes. In each 10-G RX, the recovered data are deserialized by a factor of 16 for the digital domain interface. A 644.53-MHz 16-b bus is reordered by a 16-to-66-b converter into a 156.25-MHz 66-b bus, because the Ethernet signal is composed of the sequences of encoded 66-b blocks. To support up to ten 10GBASE-R and up to two 40GBASE-R data streams, the MLG MUX includes two Ethernet frame processing units, i.e., a 10-GbE unit, including ten mapping modules, and a 40-GbE unit, including two

mapping modules. The 10- and 40-GbE mapping modules convert the 10- and 40-GbE streams into two and eight MLG streams, respectively. The mapping modules operate with two different clock signals, one from the corresponding 10-G RX and the other from the 25-G TX. The two clock domains can have a rate difference of up to ± 100 ppm because of the plesiochronous nature of the Ethernet signal and the periodic insertion of an MLG lane marker. These rate differences are compensated using an FIFO within the mapping module. Unused mapping modules are disabled based on the physical lane configuration so as to save power consumption. The lane selection logic chooses the enabled MLG lanes from

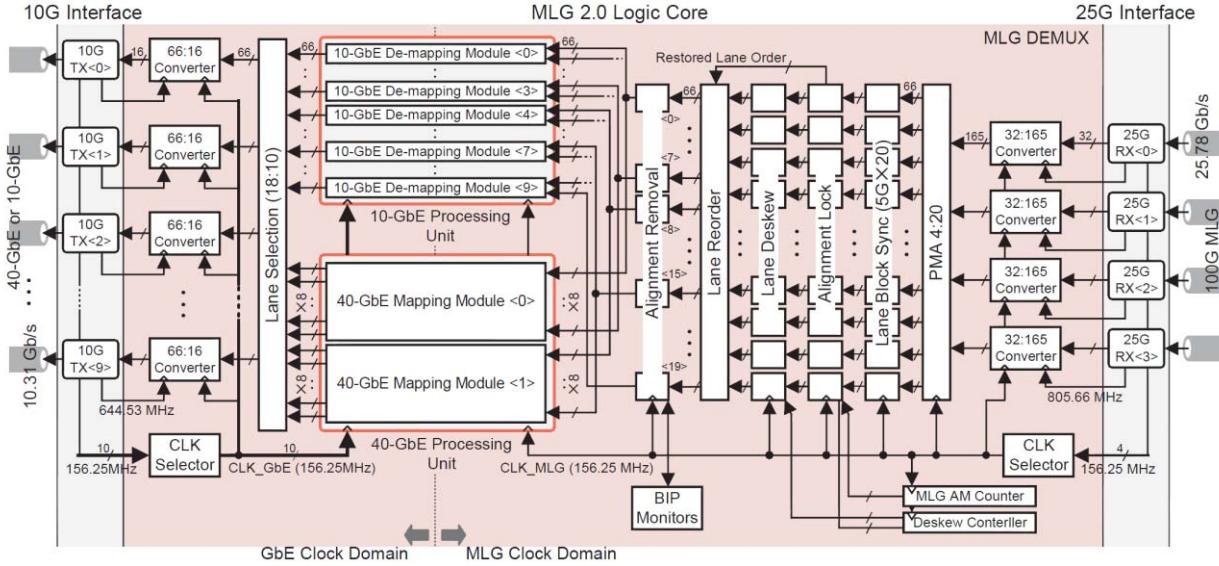


Fig. 17. Block diagram of the implemented MLG DEMUX.

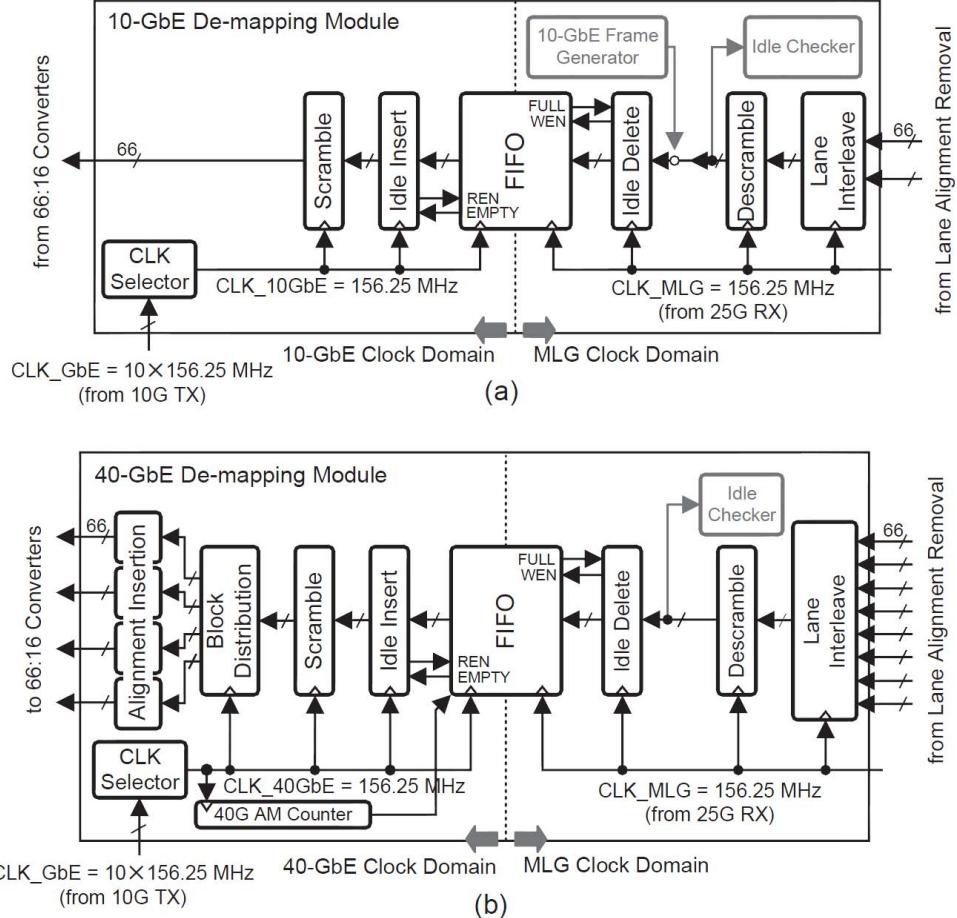


Fig. 18. Block diagram of Ethernet frame demapping module in the MLG DEMUX. (a) 10-GbE demapping module. (b) 40-GbE demapping module.

the mapping modules. Then, alignment markers are inserted into each MLG lane for proper deskewing and reordering at the MLG DEMUX [18]. The resulting MLG lanes are bit multiplexed by a 20:4 PMA into four lanes, each with a 165-b bus. A 165 to 32-b converter maps the 165-b parallel data onto

a 32-b parallel bus for serialization. The combined 100-Gb/s stream is structured in a similar way to a 100-GbE signal comprising twenty PCS lanes. However, the main difference is that the MLG 2.0 stream contains MLG 2.0 lane markers replacing idle bits [17].

Fig. 15 shows a detailed block diagram of the 10- and 40-GbE mapping modules in the 10- and 40-GbE processing units. For a 10-GbE stream [Fig. 15(a)], 66-b block synchronization is performed by monitoring the 2-b sync header. Then, the 10-GbE signal is descrambled except for the sync header and mapped into a common MLG clock domain, which is generally asynchronous with the 10-GbE domain. A conventional half-full FIFO [19] and its associated control logic blocks perform seamless data transfer between two different clock domains without any corruption, as shown in Fig. 16(a). Upon reset, the FIFO is initialized to be half full. The read and write pointers are located at the initial and the middle of the FIFO, respectively. When the 10-GbE clock is faster than the MLG clock, the FIFO will be filled with incoming data, because write events occur more frequently than read events. Then, the distance between the two pointers will be greater than half the length of the FIFO. In this case, the FIFO generates a FULL status flag to avoid FIFO overflow. Under the FULL status flag, eight bytes of idle characters after another eight bytes of idle characters are deleted at the write side until the FIFO returns to its (initial) half-full state. However, if the 10-GbE clock is slower than the MLG clock, the distance between the two pointers will become smaller than half the size of the FIFO. In this case, a group of eight bytes of idle characters is inserted after another eight bytes of idle characters at the read-side to prevent underflow [see Fig. 16(b)]. Finally, the 10-GbE stream is scrambled and distributed over two MLG lanes. The polynomial of the scrambler is given by [17]

$$G(x) = 1 + x^{39} + x^{58}. \quad (6)$$

For verification of the logical functionalities of the test chip, an internal 10-GbE frame generator is implemented.

The 40-GbE signal is mapped into eight MLG lanes via a similar process performed in the 10-GbE case. The differences are alignment locking, deskewing, reordering, and interleaving of the received 40-GbE frame prior to the descrambling process. Once 66-b block synchronization is completed, the mapping module achieves alignment lock by detecting 40-GbE PCS alignment markers on all lanes and eliminates skewing between lanes. The bit interleaved parity (BIP) value within the received alignment marker is compared with the calculated ideal BIP to detect infrequent error events [6]. The deskewed lanes are reordered based on the information gathered during the alignment lock process. Once the PCS lane alignment markers are removed from the 40-GbE signal, the PCS lanes are interleaved and descrambled prior to clock matching between 40 GbE and MLG domains. Then, the rate-matched data are scrambled and distributed to eight MLG lanes.

In contrast, the MLG DEMUX receives 4×25.78 Gb/s data consisting of twenty MLG lanes from a client-side interface and reconstructs the original 10- and 40-GbE streams, as shown in Fig. 17. In each 25-G RX, the recovered MLG data are deserialized to 32 parallel 805.66-Mb/s streams and is then mapped to a 156.25-MHz 165-b parallel bus using a 32-to-165-b converter in the digital domain. Four lanes of 165 parallel bits are demultiplexed by a 4:20 PMA into

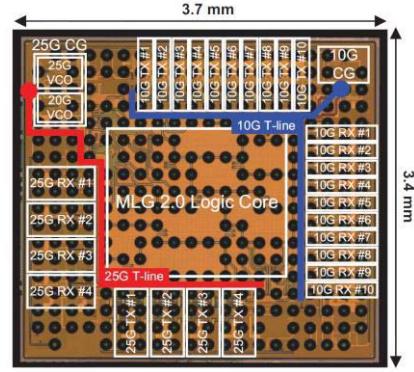


Fig. 19. Die photograph of the reverse gearbox IC.

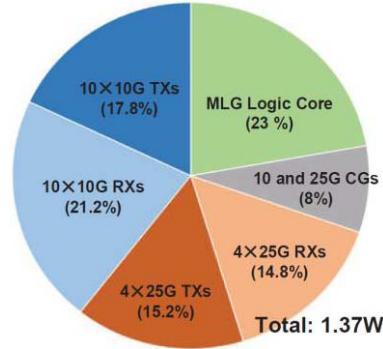


Fig. 20. Power breakdown of the reverse gearbox IC.

twenty 5.16-Gb/s MLG lanes with a 66-b interface, because each 25-G lane carries five MLG lanes, and 66-b block synchronization is performed for each of the MLG lanes. Then, the MLG DEMUX finds the MLG lane alignment markers from among all lanes for deskewing and reordering. The BIP value stored within the alignment marker is monitored in the same manner as in the 40-GbE MLG MUX case, so as to detect infrequent error events. Once the MLG lane alignment markers are removed from the MLG streams, two Ethernet frame processing units, i.e., a 10-GbE unit, including ten demapping modules, and a 40-GbE processing unit, including two demapping modules rebuild a 10-GbE stream from the two MLG lanes and a 40-GbE stream from the eight MLG lanes. For reconstructing Ethernet streams from the MLG streams, demapping modules employing FIFOs compensate for the instantaneous rate difference between two different clock domains. Ten sets of 66-b bus from the enabled demapping modules are distributed to ten parallel 66-to-16-b converters using the lane selection logic. Finally, 16 parallel bits in each lane are serialized by the 10-G TX to form an Ethernet stream.

Fig. 18 shows a detailed block diagram of the 10- and 40-GbE demapping modules. A 10-GbE stream is constructed by combining two MLG lanes [Fig. 18(a)] in a demapping module. Two MLG lanes without the alignment marker are interleaved/descrambled and mapped to the 10-GbE clock domain. Finally, the reconstructed 10-GbE signal is scrambled.

A 40-GbE signal is reconstructed from eight MLG lanes using similar process as described in the 10-GbE case

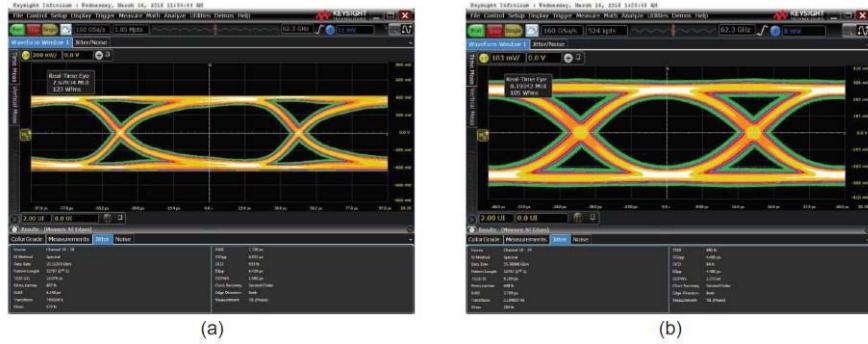


Fig. 21. Measured eye diagram and jitter decomposition. (a) 10-G TX. (b) 25-G TX.

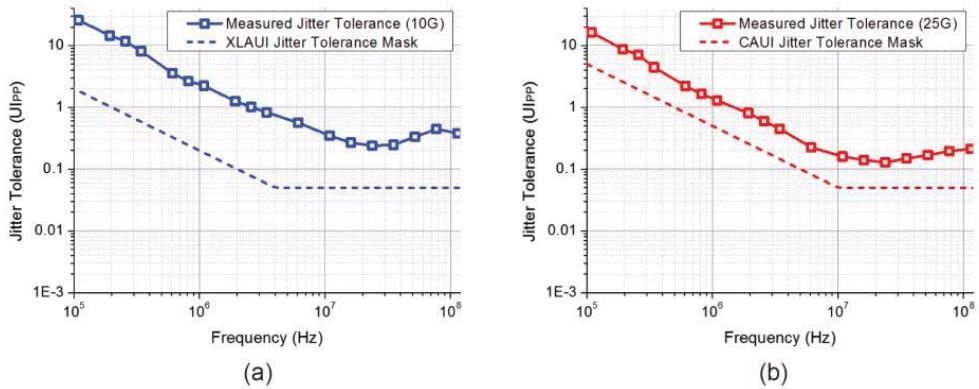


Fig. 22. Measured sinusoidal jitter tolerance of (a) 10- and (b) 25-G interface.

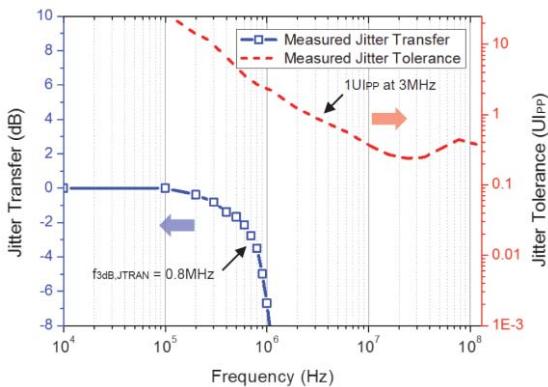


Fig. 23. Measured jitter transfer function of the 10-G transceiver.

[Fig. 18(b)]. After the idles are inserted or deleted for synchronization with the 40-GbE clock domain, each lane is scrambled and block distributed to four 10-G lanes. Then, 40GBASE-R PCS alignment lane markers are inserted into all lanes.

IV. EXPERIMENTAL RESULTS

Fig. 19 shows a microphotograph of the proposed reverse gearbox IC. The reverse gearbox IC was fabricated in a 40-nm CMOS process and was flip-chip packaged; 10-lane 10-G transceivers, 4-lane 25-G transceivers, two clock generators, and the MLG 2.0 logic core occupy an area of $3.7 \times 3.4 \text{ mm}^2$. The clock signals for 10- and 25-G interfaces are distributed from common 10- and 25-G clock generators

to each transceiver block via on-chip transmission lines. The reverse gearbox IC consumes 1.37 W of power in total; its power breakdown is shown in Fig. 20. The 10-G interface, 25-G interface, clock generators, and MLG 2.0 logic core dissipate 534.3 mW (10-G TXs: 17.8% and 10-G RXs: 21.2%), 411.1 mW (25-G TXs: 15.2% and 25-G RXs: 14.8%), 109.6 mW (8%), and 315.1 mW (23%), respectively. Compared with the conventional gearbox IC in [23], the proposed reverse gearbox IC dissipates 25% less power while supporting both conventional gearbox and MLG 2.0 functionality.

Fig. 21 shows the measured TX eye diagrams of the 10- and 25-G interfaces for 10.31- and 25.78-Gb/s PRBS-15 patterns, and 4-dB pre-emphasis at 12.5 GHz is applied to compensate for the package, evaluation board, and cable losses in the measurements. For the 25-G TX, the clock phases are manually calibrated using phase trimming circuitry to minimize Deterministic jitter (DJ) caused by phase mismatch among multiphase clocks. The measured eye opening amplitudes are 750 and 483 mV_{ppd} for 10- and 25-G interfaces, respectively. The 10-G TX has a 407-fs_{rms} Random jitter (RJ) and a 6.35-ps DJ, and the 25-G TX has a 448-fs_{rms} RJ and a 3.79-ps DJ. The measured input sensitivities of 10- and 25-G RXs are 20 and 42 mV_{ppd}, respectively for a BER of 10^{-12} .

Fig. 22 shows the sinusoidal jitter tolerance performances of the implemented CDRs at a BER of 10^{-12} with a PRBS-31 data pattern. The jitter tolerance is measured in a 5-dB loss channel consisting of PCB traces, connectors, and a cable.

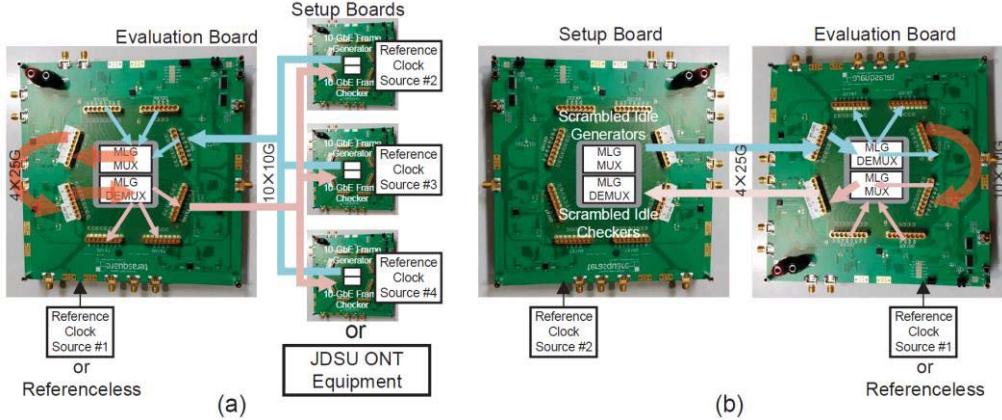


Fig. 24. MLG 2.0 function test configuration.

TABLE I
PERFORMANCE SUMMARY AND COMPARISON

Design	This Work	G. Ono ISSCC 2011 [22]	M. Harwood ISSCC 2012 [23]	J.-Y. Jiang ISSCC 2013 [24]	U. Singh ISSCC 2014 [25]
Process	40nm CMOS	65nm CMOS	40nm CMOS	65nm CMOS	40nm CMOS
Data-rate (Gb/s)	10×10 → 4×25 4×25 → 10×10	10×10 → 4×25 4×25 → 10×10	4×28 → 4×28	10×10 → 4×25 4×25 → 10×10	4×28 → 4×28
Supported Logic Function	MLG 2.0 /100-GbE Gearbox /4×20G Transmission	100-GbE Gearbox	N/A	100-GbE Gearbox	N/A
TX Output Jitter (psrms)	RJ: 0.448 psrms @ 25.78 Gb/s	0.43 psrms (10k-to-100MHz)	0.35 psrms (100k-to-1GHz)	0.187 psrms (100-to-1GHz)	0.165 psrms (10k-to-100MHz)
Input Sensitivity (mVppd)	42 @ 25.78 Gb/s	34.4 @ 25.78 Gb/s	-	-	27 @ 28 Gb/s
Power Consumption (W)	1.37 (4×25G TRX, 10×10G TRX, Logic Area)	1.99 (4×25G TRX, 10×10G TRX, Logic Area)	0.9 (4×28G TRX)	1.84 (4×25G TRX)	0.78 (4×28G TRX)
Chip Area (mm ²)	3.7×3.7	6.3×3.7	2.4×1.5	TX: 1.2×1.1×2 RX: 1.9×1.3×2	-
Supply Voltage (V)	0.9/ 2.5 (BGR, Bias Blocks)	1.0/1.8	-	1.2	0.9

Fig. 23 shows the measured jitter transfer function of the 10-G transceiver using the input jitter amplitude of 0.5 UIpp. The jitter transfer bandwidth, $f_{3dB,JTRAN}$, is approximately 800 kHz, whereas the jitter tolerance corner frequency is approximately 10 MHz. This result clearly demonstrates the jitter filtering capability of the proposed D/PLL-based transceiver without sacrificing jitter tracking performance.

Fig. 24 shows two test setups that are used to verify the MLG functionalities for 10- and 40-GbE frame transmissions. The transmission of 10-GbE frames is verified using the integrated 10-GbE frame generators and checkers, because conventional PRBS patterns cannot be directly applied for the verification of the logic functionality associated with the MLG. However, the PRBS patterns are still useful for the characterization of the PMA, including the SerDes, AFEs of the transceiver, and the CDR. The 10 × 10 GbE test patterns

are generated by multiple setup boards. The setup boards are synchronized to different reference clock sources with less than ±100 ppm differences to emulate a plesiochronous environment. The generated test patterns are transmitted to the DUT in the evaluation board. The test patterns pass through the external loopback path between the MLG MUX and MLG DEMUX, and return to the checker in the setup board. In this measurement, error-free transmission is obtained, which implies that the MLG MUX/DEMUX decodes the lane order and routes the 10-GbE signals properly. Logical functionalities of the implemented PCS are verified using JDSU OTN test equipment.

To verify the logical functionality of the data path handling the 40-GbE frame within the MLG MUX and MLG DEMUX, a scrambled idle test pattern generator and a checker are used in the remote loop back operation [6]. The MLG functionality

is evaluated by monitoring the states of the block lock, the alignment marker, and deskewing, in addition to the actual descrambled idle patterns. The scrambled idle test pattern is generated by injecting idle patterns into the scrambler in the MLG MUX. The test pattern is descrambled in the MLG DEMUX for evaluation. In this case, the descrambled outputs, including the sync header should all be idle bits and any mismatches, are counted as errors [6]. Table I compares the proposed reverse gearbox IC with previously published works that support data rates of over 25 Gb/s. Only the proposed IC supports MLG 2.0 functionality; additionally, it consumes the least power among the gearbox ICs.

V. CONCLUSION

This paper presents the industry's first 103.125-Gb/s reverse gearbox IC satisfying the OIF MLG 2.0 standard in 40-nm CMOS. The proposed IC includes ten parallel 10-G transceivers and four parallel 25-G transceivers, and enables the transmission of multiple asynchronous 10- and 40-GbE data streams across 4×25 G physical lanes. Each transceiver adopts an all digital open-loop controlled PI5-based D/PLL. The proposed D/PLL architecture enables power-and-area efficient implementation while achieving acceptable jitter filtering. All PI-based transceivers operate independently without a reference clock signal. To improve the power efficiency, the 10- and 25-G transceivers employ quarter-rate and half-rate clocking schemes, respectively, to extensively utilize the CMOS logic gates while minimizing the amount of parallelization. The MLG 2.0 logic core supports transition of up to ten 10-GbE stream and up to two 40-GbE stream to 4×25.78 Gb/s MLG streams by applying 100-GbE PCS and PMA functions. The proposed reverse gearbox IC consumes only 1.37 W while implementing complex MLG 2.0 functionalities. The power consumption is roughly 25% less than that of existing regular gearbox ICs supporting only simple MUX and DEMUX functionality.

REFERENCES

- [1] Cisco Global Cloud Index: Forecast and Methodology 2013–2018. [Online]. Available: http://www.cisco.com/c/en/us/solutions/collateral/service-provider/global-cloud-index-gci/Cloud_Index_White_Paper.pdf
- [2] IEEE Standard for Information Technology—Local and Metropolitan Area Networks—Specific Requirements—Part 3: CSMA/CD Access Method and Physical Layer Specifications Amendment 4: Media Access Control Parameters, Physical Layers, and Management Parameters for 40 Gb/s and 100 Gb/s Operation, IEEE Standard 802.3ba-2010 (Amendment to IEEE Standard 802.3-2008), Jun. 2010.
- [3] D. Law, D. Dove, J. D'Ambrosia, M. Hajduczenia, M. Laubach, and S. Carlson, “Evolution of Ethernet standards in the IEEE 802.3 working group,” *IEEE Commun. Mag.*, vol. 51, no. 8, pp. 88–96, Aug. 2013.
- [4] Data Center Design Considerations With 40GbE and 100GbE. [Online]. Available: http://en.community.dell.com/techcenter/extras/m/white_papers/20434277/download
- [5] Optical Internetworking Forum. (May 2012). Multi-link Gearbox Implementation Agreement-IA# OIF-MLG-01.0. [Online]. Available: http://www.oiforum.com/public/documents/OIF_MLG-01.0.pdf
- [6] Optical Internetworking Forum. (Apr. 2013). Multi-Link Gearbox Implementation Agreement-IA# OIF-MLG-02.0. [Online]. Available: <http://www.oiforum.com/public/documents/OIF-MLG-02.0.pdf>
- [7] H. Won *et al.*, “A 0.87 W transceiver IC for 100 Gigabit Ethernet in 40 nm CMOS,” *IEEE J. Solid-State Circuits*, vol. 50, no. 2, pp. 399–413, Feb. 2015.
- [8] N. Da Dalt, “Markov chains-based derivation of the phase detector gain in bang-bang PLLs,” *IEEE Trans. Circuits Syst. II, Express Briefs*, vol. 53, no. 11, pp. 1195–1199, Nov. 2006.
- [9] R. Inti, W. Yin, A. Elshazly, N. Sasidhar, and P. K. Hanumolu, “A 0.5-to-2.5 Gb/s reference-less half-rate digital CDR with unlimited frequency acquisition range and improved input duty-cycle error tolerance,” *IEEE J. Solid-State Circuits*, vol. 46, no. 12, pp. 3150–3162, Dec. 2011.
- [10] J. Han, J. Yang, and H.-M. Bae, “Analysis of a frequency acquisition technique with a stochastic reference clock generator,” *IEEE Trans. Circuits Syst. II, Express Briefs*, vol. 59, no. 6, pp. 336–340, Jun. 2012.
- [11] J. Yang *et al.*, “Clock recovery, receiver, and communication system for multiple channels,” U.S. Patent 13,829,566, Mar. 15, 2013.
- [12] T. H. Lee and J. F. Bulzacchelli, “A 155-MHz clock recovery delay-and phase-locked loop,” *IEEE J. Solid-State Circuits*, vol. 27, no. 12, pp. 1736–1746, Dec. 1992.
- [13] J. Kenney *et al.*, “A 9.95–11.3-Gb/s XFP Transceiver in 0.13- μ m CMOS,” *IEEE J. Solid-State Circuits*, vol. 41, no. 12, pp. 2901–2910, Dec. 2006.
- [14] W. S. Titus and J. G. Kenney, “A 5.6 GHz to 11.5 GHz DCO for digital dual loop CDRs,” *IEEE J. Solid-State Circuits*, vol. 47, no. 5, pp. 1123–1130, May 2012.
- [15] G. Shu *et al.*, “A reference-less clock and data recovery circuit using phase-rotating phase-locked loop,” *IEEE J. Solid-State Circuits*, vol. 49, no. 4, pp. 1036–1047, Apr. 2014.
- [16] W. Dettloff *et al.*, “A 32mW 7.4Gb/s protocol-agile source-series-terminated transmitter in 45nm CMOS SOI,” in *Proc. IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2010, pp. 370–371.
- [17] G. Nicholl, M. Gustlin, and O. Trainin, “A physical coding sublayer for 100GbE [applications & practice],” *IEEE Commun. Mag.*, vol. 45, no. 12, pp. 4–10, Dec. 2007.
- [18] IEEE Standard for Ethernet, IEEE Standard 802.3-2015 (Revision of IEEE Standard 802.3-2012), Mar. 2016.
- [19] J. Winkles. (2003). *Elastic Buffer Implementations in PCI Express Devices*. [Online]. Available: <http://www.docin.com/p-118906415.html>
- [20] M.-S. Chen, Y.-Y. Shih, C.-L. Lin, H.-W. Hung, and J. Lee, “A fully-integrated 40-Gb/s transceiver in 65-nm CMOS technology,” *IEEE J. Solid-State Circuits*, vol. 47, no. 3, pp. 627–640, Mar. 2012.
- [21] G. Ono *et al.*, “A 10:4 MUX and 4:10 DEMUX gearbox LSI for 100-gigabit Ethernet link,” *IEEE J. Solid-State Circuits*, vol. 46, no. 12, pp. 3101–3112, Dec. 2011.
- [22] M. Harwood *et al.*, “A 225mW 28Gb/s SerDes in 40nm CMOS with 13dB of analog equalization for 100GBASE-LR4 and optical transport lane 4.4 applications,” in *Proc. IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2012, pp. 326–327.
- [23] J.-Y. Jiang, P.-C. Chiang, H.-W. Hung, C.-L. Lin, T. Yoon, and J. Lee, “100Gb/s Ethernet chipsets in 65nm CMOS technology,” in *Proc. IEEE Int. Solid-State Circuits Conf. (ISSCC) Dig. Tech. Papers*, Feb. 2013, pp. 120–121.
- [24] U. Singh *et al.*, “A 780 mW 4 × 28 Gb/s transceiver for 100 GbE gearbox PHY in 40 nm CMOS,” *IEEE J. Solid-State Circuits*, vol. 49, no. 12, pp. 3116–3129, Dec. 2014.



Taehun Yoon received the B.S. degree in electronic and electrical engineering from Hanyang University, Seoul, South Korea, in 2005, and the M.S. degree in electronic and electrical engineering from the Pohang University of Science and Technology, Pohang, South Korea, in 2007. He is currently pursuing the Ph.D. degree with the Korea Advanced Institute of Science and Technology, Daejeon, South Korea.

From 2007 to 2011, he was a Research Engineer with Hynix Semiconductor, Icheon, South Korea.

Mr. Yoon was a co-recipient of the Presidential Prize at the 14th Korea Semiconductor Design Contest Award in 2013.



Joon-Yeong Lee received the B.S. and M.S. degrees in electrical engineering from the Korea Advanced Institute of Science and Technology, Daejeon, South Korea, in 2011 and 2013, respectively, where he is currently pursuing the Ph.D. degree.

His current research interests include clock and data recovery circuits, phase-locked loops, high-speed serial links, frequency synthesizers, spread spectrum clock generators, and low-power mixed signal circuits.



Jinhee Lee was born in Seoul, South Korea. He received the B.S., M.S., and Ph.D. degrees in electrical engineering from Seoul National University, Seoul, in 1999, 2001, and 2008, respectively.

He was a Senior Engineer with GCT Semiconductor, San Jose, CA, USA, from 2008 to 2012. He was a Principal Engineer and a Lead Engineer of multilink gearbox with TeraSquare Inc., Seoul. Since 2014, he has been the Senior Director of Technology with HivICs, Seongnam, South Korea. His current

research interests include Ethernet PHY layers, low-power mixed-signal circuits, and digital phase-locked loops.



Sangeun Lee was born in Seoul, South Korea. She received the B.S. degree in electrical engineering from Dankook University, Seoul, in 2003, and the M.S. degree in electrical engineering from The State University of New York at Buffalo, Buffalo, NY, USA, in 2006.

Since 2006, she has been with TeraSquare, Inc., Seoul, as an Analog Mixed-Signal Design Engineer, where she is involved in analog front ends, clock and data recovery circuits, and low-power mixed-signal circuits for high-speed serial transceiver.



Taeho Kim received the B.S., M.S., and Ph.D. degrees in electronics engineering from Inha University, Incheon, South Korea, in 2007, 2009, and 2012, respectively.

From 2013 to 2015, he was with TeraSquare Inc., Seoul, South Korea, where he was involved in the design of 100-G Ethernet transceiver. His current research interests include mixed-mode high-speed serial links interface and signal integrity.



Kwangseok Han (S'99–M'05) received the B.S., M.S., and Ph.D. degrees in electrical engineering from the Korea Advanced Institute of Science and Technology, Daejeon, South Korea, in 1998, 2000, and 2004, respectively.

From 2004 to 2007, he was with Samsung, Yongin, South Korea, as a Senior RF Circuit Designer for GSM Transceiver Development. From 2007 to 2013, he was with Cambridge Silicon Radio, Cambridge, U.K., as a Principal RF/Mixed Signal Designer, where he developed transceivers for the

IEEE 802.11 wireless local area network, Bluetooth, NFC, and Bluetooth Low Energy for mass production. In 2013, he joined TeraSquare Inc., Seoul, South Korea, as the Senior Director of Technology, where he has been leading the development of a low-power 100-Gb/s Ethernet IC. His current research interests include thermal noise modeling for CMOS, RF/analog circuit design, such as low-noise amplifiers, mixers, ADCs, voltage-controlled oscillators, phase-locked loops, clock data recovery, high-speed IOs, and system-level modeling.



Jinho Han was born in Incheon, South Korea. He received the B.S., M.S., and Ph.D. degrees in electrical engineering from the Korea Advanced Institute of Science and Technology, Daejeon, South Korea, in 2006, 2009, and 2014, respectively.

From 2013 to 2014, he was a Senior Analog Circuit Design Engineer with TeraSquare Inc., Seoul, South Korea, where he was involved in the development of low-power 100-Gb/s Ethernet IC and multilink gearbox. Since 2014, he has been a Senior Engineer with the Digital IP Development Team, Samsung Electronics, Hwasung, South Korea. His current research interests include high-speed serial links, frequency synthesizers, digital phase-locked loops, multiplying delay-locked loops, clock and data recovery circuits, and low-power mixed-signal circuits.

Dr. Han was a recipient of the Honor Prize at the 18th Humantech Paper Award in 2011 and a co-recipient of the Presidential Prize at the 14th Korea Semiconductor Design Contest Award in 2013. He currently serves as a Reviewer of the IEEE TRANSACTIONS ON VERY LARGE SCALE INTEGRATION SYSTEMS.



Jeong-Sup Lee received the B.S. and M.S. degrees in electrical engineering from the Korea Advanced Institute of Science and Technology, Daejeon, South Korea, in 2007 and 2009, respectively. He is currently pursuing the Ph.D. degree at the University of Michigan, Ann Arbor, MI, USA.

From 2009 to 2012, he was with TLi Inc., Korea, as a Technical Research Personnel (alternative military service), where he was involved in the design of high-speed video interface. From 2012 to 2015, he was with TeraSquare Inc., Seoul, South Korea, where he was involved in the development of low-power 100 Gb/s solutions.



Hyosup Won received the B.S. degree in microelectronics from Tsinghua University, Beijing, China, in 2009, and the Ph.D. degree in electrical engineering from the Korea Advanced Institute of Science and Technology, Daejeon, South Korea, in 2016.

He has been involved in the design of high-speed serial links, clock and data recovery circuits, and low-power mixed-signal circuits.

Dr. Won was a co-recipient of the Presidential Prize at the 14th Korea Semiconductor Design Contest Award in 2013.



Jinho Park received the B.S. degree in electrical engineering from Seoul National University, Seoul, South Korea, in 1996, and the Ph.D. degree in electrical engineering from the University of Washington, Seattle, WA, USA, in 2003.

He was with Marvell Semiconductor, Santa Clara, CA, USA, from 2003 to 2012, and leading the analog and RF design aspects of the world's first 802.11ac mobile MIMO IC publicly announced in 2012. He is currently with TeraSquare Co., Ltd., Seoul, fabless 100-Gb/s Ethernet IC company, as a CEO, and an Adjunct Professor of Electrical Engineering with the Daegu Institute of Science and Technology, Daegu, South Korea. He has co-authored the *Parasitic-Aware Optimization of CMOS RF Circuits* (Kluwer Academic, 2003) and has over 60 IEEE publications and U.S. patents.



Hyeon-Min Bae (M'09) received the B.S. degree in electrical engineering from Seoul National University, Seoul, South Korea, in 1998, and the M.S. and Ph.D. degrees in electrical and computer engineering from the University of Illinois at Urbana-Champaign, Champaign, IL, USA, in 2001 and 2004, respectively.

From 1995 to 1996, he served his military duty in Dokdo, East Sea. From 2001 to 2007, he led the analog and mixed-signal design aspects of OC-192 MLSE-based EDC ICs with Intersymbol Communications, Inc., Champaign. From 2007 to 2009, he was with Finisar Corporation (NASDAQ: FNSR), Sunnyvale, CA, USA, after its acquisition of Intersymbol Communications Inc. Since 2009, he has been with the Faculty of Electrical Engineering, Korea Advanced Institute of Science and Technology, Daejeon, South Korea, where he is currently an Associate Professor. In 2010, he founded Terasquare, Inc., Seoul, a venture-funded fabless semiconductor start-up, which provided low power all digital 100-Gb/s IC solutions. Terasquare, Inc. was acquired by Gigpeak (NYSE:GIG) in 2015. In 2013, He also founded OBElab, Inc., Seoul, a bio start-up that manufactures portable functional brain imaging systems. His current research interests include a wide range of topics in wire line communication and medical imaging systems.

Dr. Bae received the Excellence Award from the National Academy of Engineering of Korea in 2013 and the 2006 IEEE JOURNAL OF SOLID-STATE CIRCUITS Best Paper Award.