

# Ground-Referenced Signaling for Intra-Chip and Short-Reach Chip-to-Chip Interconnects

Walker J. Turner<sup>1</sup>, John W. Poulton<sup>1</sup>, John M. Wilson<sup>1</sup>, Xi Chen<sup>2</sup>, Stephen G. Tell<sup>1</sup>, Matthew Fojtik<sup>1</sup>, Thomas H. Greer III<sup>1</sup>, Brian Zimmer<sup>2</sup>, Sanquan Song<sup>2</sup>, Nikola Nedovic<sup>2</sup>, Sudhir S. Kudva<sup>2</sup>, Sunil R. Sudhakaran<sup>2</sup>, Rizwan Bashirullah<sup>3</sup>, Wenxu Zhao<sup>4</sup>, William J. Dally<sup>2</sup>, C. Thomas Gray<sup>1</sup>

<sup>1</sup>NVIDIA, Durham, NC, <sup>2</sup>NVIDIA, Santa Clara, CA, <sup>3</sup>University of Florida, Gainesville, FL, <sup>4</sup>Now with Broadcom, Irvine, CA

**Abstract**—While high-speed single-ended signaling maximizes pin and wire utilization within on- and off-chip serial links, problems associated with conventional signaling methods result in energy inefficiencies. Ground-referenced signaling (GRS) solves many of the problems of single-ended signaling systems and can be adapted for signaling across RC-dominated channels and LC transmission lines. The combination of GRS and clock forwarding enables simple but efficient signaling across on-chip communication fabrics, off-chip organic packages, and off-package printed circuit boards. Various methodologies compatible with GRS are presented in this paper, including design considerations and various circuit architectures. Experimental results for multiple generations of GRS-based serial links are presented, which includes a 16Gb/s 170fJ/b/mm on-chip link, a 20Gb/s 0.58pJ/b link across an organic package, and a 25Gb/s 1.17pJ/b link signaling over a printed-circuit board.

**Keywords**—Ground Referenced Signaling, High-Speed Serial Link, Single Ended Signaling, Transceiver

## I. INTRODUCTION

As the semiconductor industry approaches the end of the ‘Moore’s-Law era’, the ability to integrate increased functionality onto a single chip is becoming limited. In the past, multi-chip modules (MCM’s) were used in advance of semiconductor scaling to implement systems too large for a single die. MCM’s may become more important to the industry because of ultimate limitations in transistor scaling. The ability to split complex functionality across multiple chips requires energy-efficient chip-to-chip signaling that can support I/O bandwidth approaching on-chip bisection bandwidth. Power budgeting requires both on- and off-chip communication circuitry to consume a small fraction of total die power, while packaging technology severely limits the number of I/O pins available for chip-to-chip communication.

These problems have led to increasing interest in low-power intra-chip links [1-3] and short-reach off-chip links [4-8], where some of these experimental links give up the numerous advantages of differential signaling to conserve on-chip routing resources and signaling I/O pins. Although single-ended (SE) systems can provide a 2x advantage in pin utilization, the signal-integrity and circuit-design challenges make it difficult to achieve the data rates of differential signaling systems. SE signaling is also fundamentally disadvantaged in energy efficiency; for the same investment in energy, SE systems develop only half the receiver voltage as differential systems.

This research was, in part, funded by the U.S. Government under the DARPA CRAFT program. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government.

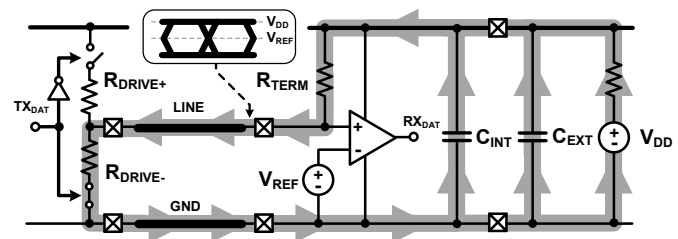


Fig. 1. Typical single-ended signaling interconnect.

Problems encountered in conventional SE systems are covered in Section II. Section III describes ground-referenced signaling (GRS), a novel signaling method that avoids many of these problems. Section IV introduces the channel properties of on- and off-chip short-reach links. Section V describes short-reach link organization and clocking architectures. Sections VI and VII detail the circuit design of GRS transmitters and receivers, respectively. Section VIII describes three GRS experimental links: an on-chip communication fabric, an on-package chip-to-chip link for organic packages, and a package-to-package link that can operate over a short printed-circuit board channel. Section IX concludes the paper.

## II. SINGLE-ENDED SIGNALING DESIGN CONSIDERATIONS

Fig. 1 shows a typical SE signaling system. Data on a link is represented by a series of voltages, each transmitted for a unit interval (UI) of time and referred to a common terminal, usually ground. For non-return-to-zero (NRZ) signaling, two voltages represent logic 1 and 0. Typically, current is drawn continuously from one power supply for one logic value (logic 0 in Fig. 1), while transient currents flow in the supply networks at both ends of the link during switching events. At the receive end, signal current passes through a terminating resistor  $R_{TERM}$  to develop a voltage that can be detected by the receiver. Current returns to the transmitter over the common reference conductor(s) to complete the circuit. The single-ended nature of this scheme has three major problems that make it difficult to build a robust high-speed SE link without sacrificing energy-efficiency.

### A. Simultaneous Switching Noise

Simultaneous switching noise (SSN) is the best-known problem of SE signaling and results from data-dependent current consumption. While transient current can flow within the power delivery networks (PDN) during data transitions, many SE systems terminate the line to the positive supply rail and only

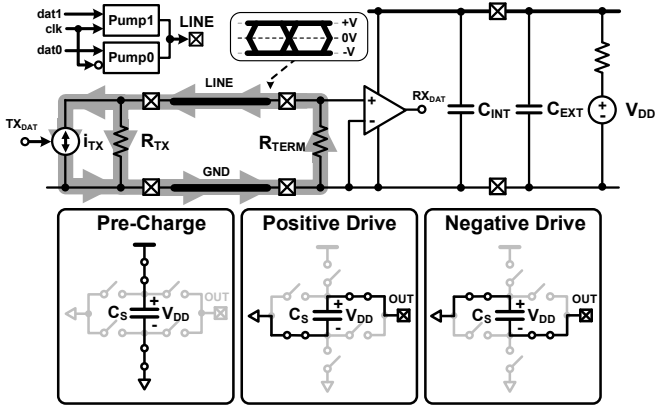


Fig. 2. Ground-referenced signaling operation.

draw supply current when transmitting a logic 0. Variations in supply current generate a fluctuating voltage drop across the PDN impedances, which appears as supply noise to adjacent lanes. This noise is additive when lanes switch simultaneously and is proportional to signal amplitude, thus it cannot be overcome with increased signal swing.

SSN occurs at both Tx and Rx ends of the link, and is often the dominant noise source in SE systems. In addition to affecting link performance, SSN adds power supply design challenges, since the supply must handle large instantaneous changes in load current. While data encoding methods, such as data-bus inversion, can reduce SSN up to 50% through balancing either DC or AC currents, one or more additional lanes are needed to transmit the coded data, thus reducing link pin-efficiency. The remaining recourse for reducing SSN is to provide more supply pins, and thus increase packaging cost.

### B. Signal Reference Offset

Since data is represented by a voltage on a single wire, the receiver requires a reference voltage for comparison with the signal to recover the data. The reference needs to be accurately maintained midway between the received voltage amplitudes. It is difficult in practice to perfectly match this reference voltage between two ends of a link. Additionally, variations in operating temperature, supply voltage, and transistor aging can cause the reference to drift over time. Reference voltage errors can only be overcome by increasing signal amplitude, often far more than needed to overcome random noise. Therefore, reference error is a major source of energy inefficiency in SE systems.

### C. Signal Return Path

In theory, signal current flows down the signal line, through the Rx terminator, and back to the transmitter through a single reference conductor. However, the signal return paths for most SE systems are more complicated. Typical termination schemes cause return current to split between supply and ground portions of the PDN, which includes internal and external bypass capacitors. Because the resulting signal return path is frequency dependent, it is difficult to achieve a high-quality impedance match at the Rx terminator and thus signal reflections are increased. Additionally, return current flowing through the PDN increases cross-talk due to shared supplies—another aspect of

SSN. This problem also exists at the transmitter, since signal current mainly flows through the supply when sending a logic 1 and through the ground network for a logic 0.

## III. GROUND REFERENCED SIGNALING

Ground referenced signaling (GRS) overcomes the main difficulties of SE signaling outlined in the previous section by adopting the ground network as the signal reference voltage. Since the ground network is usually the most robust and lowest impedance supply network in a system, mismatch in the reference voltage (*i.e.* ground) can be made very small between the transmitter and receiver. In GRS systems, ground is the one and only signal return network, which ensures a high-quality termination at both ends of the link by removing the PDN from the return path.

The operating principle of GRS is shown in Fig. 2. Signal voltages are generated symmetrically above and below the ground potential (0V) from a transmitter that acts as a bidirectional current source. At the receiver, the signal is compared to the ground potential for signal amplification and data recovery. GRS requires a negative supply (below ground potential), and charge pumps are usually employed for this purpose. The GRS transmitter combines this charge pump with the data multiplexer and line driver. It has two operating phases: *pre-charge* and *drive*. During the *pre-charge* phase, storage capacitor ( $C_S$ ) is charged to the supply voltage. On the next phase of the clock (*drive*), the charge stored in  $C_S$  is forced into the line by connecting the capacitor terminals between the line and ground nodes. The polarity of the connection drives either a positive or negative current impulse into the line, representing logic 1 or 0, respectively. Two charge-pumps (Pump1 & Pump0) drive the line on opposite clock phases to provide a 2:1 output multiplexer.

The charge-pump transmitter avoids SSN since the current drawn from the supply each half-cycle is used solely to pre-charge  $C_S$ . This stored charge is injected into the line irrespective of data polarity, resulting in near constant supply current regardless of data patterns.

## IV. CHANNEL INTERCONNECTS

GRS can be used across RC-dominated channels as well as low-attenuation LC transmission lines, making it practical for both on-chip and off-chip interconnects. The single-ended nature of this approach requires careful co-design of the link circuitry and interconnect to achieve an energy efficient link.

### A. RC-Dominated Channels

RC-dominated channels are common in on-chip interconnects and silicon interposer packages, where small wire dimensions result in high trace resistances. Link/channel co-design must deal with insertion loss that increases, while bandwidth decreases, with  $l^2$  ( $l$  is the channel length). The high attenuation of RC lines mitigates the problem of reflections, while cross-talk remains a serious concern.

Various circuit techniques can be employed to overcome the high-frequency attenuation and resulting intersymbol interference (ISI) of these channels by trading voltage swing for

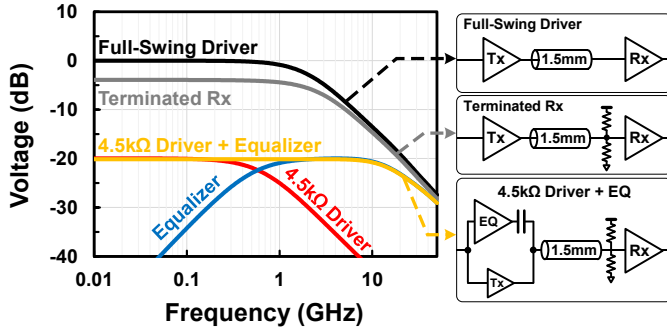


Fig. 3. Simulated signaling performance across an RC-dominated channel.

signaling bandwidth. Fig. 3 shows the frequency response of various circuits signaling across a 1.5mm RC-dominated channel ( $R=130\Omega/\text{mm}$ ,  $C=305\text{fF}/\text{mm}$ ). A full-swing 100 $\Omega$  driver achieves a -1dB bandwidth of 1GHz, a conservative estimate to ensure 2Gb/s NRZ signaling. The addition of a 525 $\Omega$  Rx terminator reduces the signal swing by -4dB while extending the bandwidth to 1.6GHz. Combining a 4.5k $\Omega$  driver (BW = 320MHz) with a high-pass transmitter composed of the 100 $\Omega$  driver in series with a 60fF coupling capacitor (BW = 1GHz to 11GHz) reduces signal amplitude to -20dB while yielding a flat frequency response up to 11GHz to support maximum signal rates of 22Gb/s.

Typical approaches to overcoming cross-talk from adjacent wires, such as spatially and/or temporally offsetting signal transitions to minimize interference, generally do not work as well when operating links near the bandwidth limit of the channel. For on-chip and silicon interposer channels, the most effective way to mitigate cross-talk is to insert shielding wires (connected to power or ground) between channels, which practically eliminates cross-talk.

### B. LC Transmission Line Channels

Channels that behave as LC transmission lines are common in interconnects with larger conductors, such as those in organic packages and PC boards. High-speed SE signaling across these channels must deal with frequency-dependent attenuation and crosstalk, while reflections become a potential problem. Stripline routes significantly reduce cross-talk through insertion of reference planes (preferably ground) above and below all signal traces. Additionally, the bump and ball patterns should have sufficient ground returns interspersed. The most effective approach is a signal-ground checkerboard pattern, where ground bumps behave as shields for vertical transitions. Ground shields can also be inserted between traces in longer-channel applications, such as the PCB link reported here.

Impedance discontinuities cause reflections in LC channels, where on-chip circuitry, ESD protection, and pad parasitic capacitances present a very low impedance to fast signal edges. A portion of the signal energy is reflected onto the channel, which results in loss of voltage and timing margin while also producing frequency dependent attenuation, leading to ISI. An effective method for reducing these effects is through series inductors or t-coils, which compensate for the parasitic capacitances by inserting series positive reactance in between.

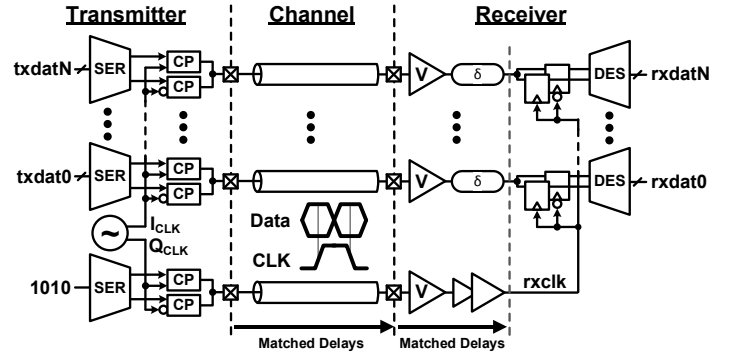


Fig. 4. Bundled-data clock forwarding architecture.

Since this technique uses passive components, power consumption does not increase.

### V. BUNDLED-DATA CLOCK FORWARDING SERIAL LINKS

Clock generation, distribution, and recovery are important features in any high-speed serial link and typically consume a large percentage of link power. Long-reach links employ per-lane trimming of the clock-to-data phase relationships, which requires very precise low-jitter clock generation to maintain the timing correlation over many bits, and per-lane measurement and control of data-to-clock phase. Most short-reach links use a bundled-data, clock-forwarding approach as the clocking architecture since it is almost always practical to match the delays across a bundled link to a tolerance  $\ll 1\text{UI}$ . By forwarding the reference clock on an identical channel, it is possible to cancel nearly all sources of jitter, while avoiding the complexity and power consumption of clock recovery circuits. A simplified clock-forwarded approach is implemented here.

Fig. 4 shows a *matched-delay* forwarded clock serial link, where the clock signal is transmitted across the channel in quadrature with the data. The clock receiver is identical to the data receivers and directly drives the input of the Rx clock distribution buffer. A programmable delay cell is included within each data signal path with an insertion delay ( $\delta$ ) and power-supply delay variation designed to match that of the Rx clock distribution network. Since signal paths are closely matched, this method essentially cancels all jitter with bandwidth nearly equal to the bit-rate. The data delay element operates on a data stream amplified to CMOS levels; it introduces some timing variation, most of which can be trimmed during link calibration, and consumes negligible energy.

A matched-delay forwarded clock architecture, while simple, is extremely effective in cancelling power-supply induced timing variations since the data and clock paths reliably track across supply voltage and temperature. Additionally, the Tx clock generator design is relaxed since clock jitter is exported along with the data and cancelled at the receiver. The generated clock is not required to have low-jitter performance or even operate at a precisely set frequency. This contributes to a very low-power clocking system relative to conventional clock-recovery schemes.

The matched-delay forwarded clock architecture can be adapted for on-chip serial links transmitting across lossy high-resistance channels. Both in-phase ( $Iclk$ ) and quadrature phase ( $Qclk$ ) clocks are forwarded with the data across RC-dominated channels comprised of on-chip metal layers. To overcome the quadratic-wire-length attenuation, intermediate re-timing stages sample and buffer the data signals along the channel to increase the reach of the serial link.  $Iclk$  and  $Qclk$  are alternated as a simple way to reset timing at each successive repeater.

## VI. GRS TRANSMITTER DESIGN

The GRS transmitter must generate positive and negative voltages symmetrically about the ground potential. This section describes the circuit details of a GRS transmitter and corresponding equalization, which can be adapted for signaling across RC-dominated channels and LC transmission lines.

### A. Charge-Pump Transmitter

Section III describes the basic idea of a charge-pump implemented as a ground-referenced data transmitter. In a 2:1 output multiplexed GRS transmitter, two pumps drive current into the line, as shown in Fig. 2. Each pump pre-charges capacitor  $C_S$  on one clock phase, and drives the line on the other phase. Charge-pump output current is  $fC_S(V_{DD} - V_{LINE})$ , where  $f$  is the switching frequency, and ideally behaves as a voltage source in series with a resistance of  $R = 1/fC_S$ . The switches introduce additional series resistance  $R_S$ .

Fig. 5 shows the transistor-level schematic of a GRS charge-pump transmitter. Transistors M0, M5 pre-charge the storage capacitor ( $C_S$ ) to the supply voltage. During the drive phase, transistor pairs M1, M4 or M2, M3 discharge the capacitor into the line to generate the positive or negative signal currents, respectively. For a given line characteristic impedance,  $C_S$  is sized to provide the desired line amplitude at the target bit-rate. The signaling rate dictates the drive-switch transistor sizes, and thus the series resistance  $R_S$ . The current developed in the transmitter is the sum of two exponentials, one rising and one decaying, whose time constants depend on  $C_S$  and  $R_S$  of the pump, the line impedance  $R_0$ , and the shunt capacitance on the transmitter output  $C_0$ . At design time, transistor sizing (and thus  $R_S$ ) is optimized so that the maximum current is driven into the line at the center of the data eye.

There is inherent asymmetry between the positive and negative signal amplitudes generated by the charge-pump transmitter. Parasitic capacitances referred to ground at  $C_{TOP}$  and

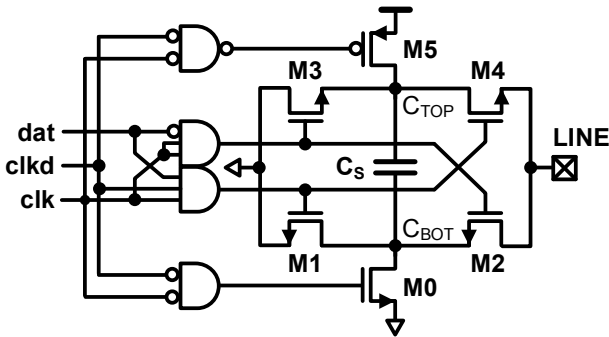


Fig. 5. Charge-pump transmitter schematic.

$C_{BOT}$  tend to provide additional positive current drive, but divert part of the current for negative drive. Overlap of pre-charge and drive waveforms also contribute to asymmetry by allowing pre-charge current to flow into the line; this effect advances rising-edge timing. This problem is avoided by using a slightly delayed clock ( $clk_d$ ) to control pre-charge and positive drive. Negative drive has an inherent delay because M2's gate voltage must rise above the threshold voltage ( $V_{TH}$ ) before conduction begins. When M2 starts to turn on, the source terminal ( $C_{BOT}$ ) begins to fall below ground, which pumps the gate voltage and boot-straps M2 conduction. To compensate for these effects, negative drive is activated on the  $clk$  phase.

Asymmetry also results from differing voltage trajectories when driving the line positive and negative. This results from the time-varying operating regimes of the drive transistors. Transistors M4 and M3 begin in saturation during positive and negative drive, respectively, and gradually transition into the triode region, while M1 and M2 spend most of the drive phase in the triode region, thus the shapes of the current impulses differ. Most of the asymmetry, however, can be tuned out at design time by sizing the negative drive transistors larger to adjust their relative drive strengths and compensate for the inherently smaller negative drive amplitude.

During negative drive, positive gate-to-source voltages are developed across M0, M1, and M4 (the sources fall below ground while the gates are at 0V), causing these transistors to enter weak to moderate inversion regimes and sink charge from the negative signal drive path. In this state, the devices operate in a diffusion-dominated regime where induced channel currents are correlated to the device threshold voltage across operating temperature. This decreases the negative signal amplitude when operating at higher temperatures, resulting in a time varying offset voltage that can affect link performance. The most straight-forward solution is to use higher-threshold devices, especially for pre-charge transistor M0, as it is the largest contributor to this temperature sensitivity. The transistors can be sized larger to compensate for the increased channel resistance of the higher- $V_{TH}$  devices while ensuring the transistors remain in the weak inversion regime during negative drive.

### B. Transmit Equalization

The GRS transmitter can be modified to perform transmit equalization (EQ), important for reducing intersymbol interference due to high-frequency signal attenuation. In our experimental links, we use the equalizing structure shown in Fig. 6 to implement additive pre-emphasis. This auxiliary EQ transmitter consists of a CMOS stage that drives charge into the line during edge transitions through a coupling capacitor ( $C_{EQ}$ ).

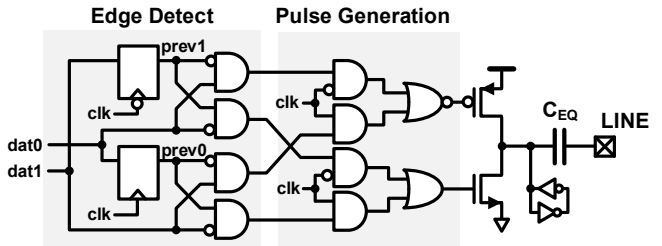


Fig. 6. Pre-emphasis transmit equalizer.

The capacitor size is chosen to boost signal transitions using a current impulse of duration  $< 1UI$ . Edge-detection and pulse-generation circuitry enable pull-up or pull-down of the output driver whenever subsequent data bits differ. The output driver is only active during edge transitions, since constantly driving the inboard capacitor node introduces a finite resistance in series with  $C_{EQ}$ , which can modulate the transmitter return impedance during long periods of idle data. To prevent the inboard node from being left in a high-Z state, small cross-coupled inverters maintain the voltage potential between transitions. This ensures full EQ strength after long periods of inactivity since the current impulse magnitude is proportional to the change in voltage at the drive node. While transmit EQ introduces data-dependent current consumption (SSN), the EQ current is about an order of magnitude smaller than the line drive circuitry.

## VII. GRS RECEIVER DESIGN

The receiver input amplifier within a GRS link requires three basic functions to recover a line signal with a common-mode voltage near ground: 1) Provide sufficient gain for signal pre-amplification, 2) Level-shift the input waveform to the CMOS inverter threshold (nominally  $V_{DD}/2$ ) to allow for inverters to amplify the signal down-stream, and 3) Compensate for high-frequency channel attenuation. Common-gate amplifiers for near-ground signaling were first described in [9] and are capable of providing the three functions required in the GRS receiver. The Rx design has evolved over the development of GRS signaling, and three such structures are summarized here.

### A. Single-Ended Common-Gate Amplifier

Fig. 7 is the first-generation GRS amplifier, a resistively loaded ( $R_{LOAD}$ ) common-gate (CG) amplifier input stage followed by a 2<sup>nd</sup> stage CMOS inverter. The CG amplifier levels-shifts the input signal while providing small-signal gain (A) set by the M0 transconductance ( $g_{m0}$ ) and  $R_L$ , such that:

$$A = \left( g_{m0} + \frac{1}{r_o} \right) \frac{r_o R_L}{R_L + r_o}$$

where  $r_o$  is the drain-to-source transconductance  $= 1/g_{ds0}$ .  $R_{LOAD}$  is chosen at design time to achieve a specific small-signal gain, and to place the output pole of the first stage above Nyquist, while M0 is subsequently sized to adjust the amplifier bias current and set the output common-mode voltage to the high-gain region of the 2<sup>nd</sup> stage (nominally  $V_{DD}/2$ ). INV2 operates

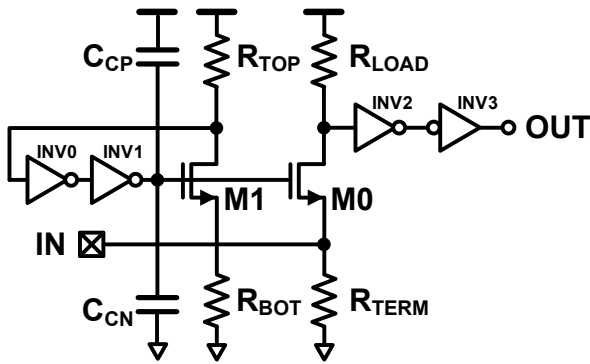


Fig. 7. Single-ended common-gate amplifier schematic.

as a linear, continuous-time amplifier to further amplify the signal to near CMOS levels.

The DC operating point is controlled by a replica, which mimics the amplifier structure with the addition of inverters INV0 and INV1 to bias M1's drain at the switching threshold of a CMOS inverter. Decoupling capacitors  $C_{CP}$  and  $C_{CN}$  stabilize the feedback loop, while tunable resistors  $R_{TOP}$  and  $R_{BOT}$  replicate the voltage drops associated with  $R_{LOAD}$  and  $R_{TERM}$ , respectively. DC operation can be trimmed for offset calibration through the P/N ratio of INV0 and the ratio of  $R_{TOP}$ -to- $R_{BOT}$ . Continuous-time linear equalization (CTLE) can be implemented by replacing  $R_{LOAD}$  with an active inductor (PMOS transistor with feedback resistance). The replica bias feedback loop has poor temperature sensitivity, which results in input-referred offset drift across operating temperatures.

### B. Complimentary Common-Gate Amplifier

To compensate for the temperature sensitivity of the CG amplifier in Fig. 7, a more symmetric amplifier using complementary transistors can be used, as shown in Fig. 8. This complimentary common-gate (CCG) amplifier is composed of two input branches: a CCG amplifier (M0A & M1A) in parallel with a CCG equalizer (M0B, M1B,  $R_{FB}$ ,  $C_{EQP}$ ,  $C_{EQN}$ ). The CCG amplifier replaces the resistive load of Fig. 7 with a PMOS transistor M1A. A voltage drop is introduced at the source of M1A using resistor  $R_{VDD}$ , which is trimmed to replicate the DC operating conditions of M0A. The CCG amplifier is biased using a replica formed by a shorted inverter (M2 & M3), that is source-degenerated to ensure the bias voltage ( $V_{BIAS}$ ) is near the switching threshold of M0A/M1A. The complementary nature of the structure inherently matches the behavior of the 2<sup>nd</sup> stage inverter (INV0), allowing for better tracking across supply and temperature variations.

The CCG equalizer is implemented in parallel to overcome a portion of the high-frequency channel attenuation and compensate for the amplifier output pole. This branch is a gyrator-C network using feedback resistor ( $R_{FB}$ ) to close the loop around transistors M0B and M1B. Active inductors are routinely used for spectrum shaping and equalization, but only recently have inverter-based gyrators been exploited [7]. Assuming  $R_{FB} \ll 1/g_{ds}$  of M0B and M1B, the feedback loop-gain ( $T_o$ ) can be approximated as  $T_o = g_m (R_{FB} + 1/sC_{EQ}) / (1 + sR_{FB}C_{EQ})$ , where  $C_{EQ}$  is the sum of

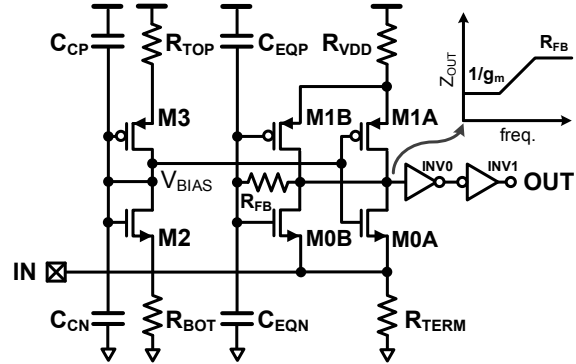


Fig. 8. Complimentary common-gate amplifier schematic.

the feedback capacitors  $C_{EQP}$  and  $C_{EQN}$  and  $g_m$  is taken as the sum of the transconductances of M0B and M1B. The EQ branch impedance is reduced by  $(1 + T_o)$  as a result of the shunt-feedback configuration, such that:

$$Z_{OUT,EQ} = \frac{1}{g_m + sC_{EQ}} + \frac{R_{FB}}{1 + g_m/sC_{EQ}}$$

Assuming  $g_m \gg 1/R_{FB}$ , the output impedance of the EQ branch behaves as an inductor with equivalent inductance of  $L \cong R_{FB}C_{EQ}/g_m$  and series resistance  $R_S \cong 1/g_m$ . At low frequencies, the EQ branch behaves as a diode-connected MOSFET ( $Z_{OUT,EQ} = 1/g_m$ ), and thus reduces gain for low-frequency signals. At high frequencies, the MOSFET gate terminals are shorted to AC ground through the equalizing capacitors, resulting in an equivalent output impedance of  $R_{FB}$ . Since, by assumption  $1/g_m \ll R_{FB}$ , a frequency response similar to that shown in Fig. 8 results.

### C. Single-Ended to Pseudo-Differential CG Amplifier

Fig. 9 shows a single-ended to pseudo-differential CG amplifier [10], consisting of CG amplifiers (M0 & M1) made differential through cross-connecting the biasing arrangement.  $R_{BOT}$  is trimmed to match the parallel combination of the Tx and Rx terminators and balance the DC operation of the two branches. Since M0 and M1 are in common-gate and common-source configurations, respectively, the small-signal gain and insertion delays differ between the two output nodes. The signal swing of  $V_{ON}$  is smaller than  $V_{OP}$  due to the asymmetric signal amplification paths, which can be approximated by:

$$A = V_{OP} - V_{ON} = g_{m0}Z_{OUT} + g_{m1}Z_{OUT}/(1 + g_{m1}R_{BOT})$$

RC networks are added to the diode-connected transistors to equalize the signal path to  $V_{ON}$  and thus reduce the delay mismatch from 0.1UI to 0.01UI between the output nodes.

CTLE is implemented through active inductors (M6 & M7) at the output nodes, which can be adjusted based on the channel loss conditions. Cross-coupled PMOS devices (M8 & M9) introduce positive feedback to effectively increase the active-inductor quality factor. The amplifier includes a common-mode (CM) feedback loop that compares the CM voltages of the 1<sup>st</sup> and 2<sup>nd</sup> stage outputs and regulates the input bias branches so the 1<sup>st</sup> stage outputs are centered on the succeeding inverter threshold ( $\sim V_{DD}/2$ ).

This structure has near constant performance across global corners and tolerates mismatches in the PMOS-to-NMOS beta

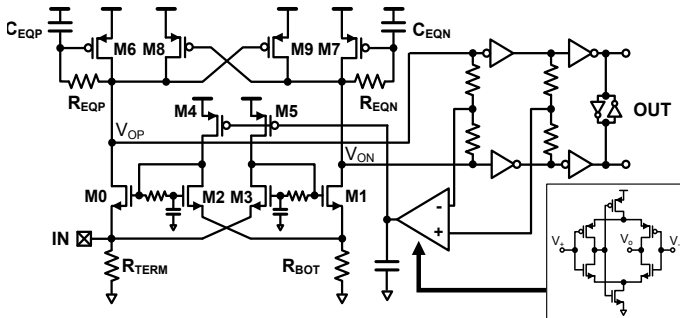


Fig. 9. Pseudo-differential common-gate amplifier schematic.

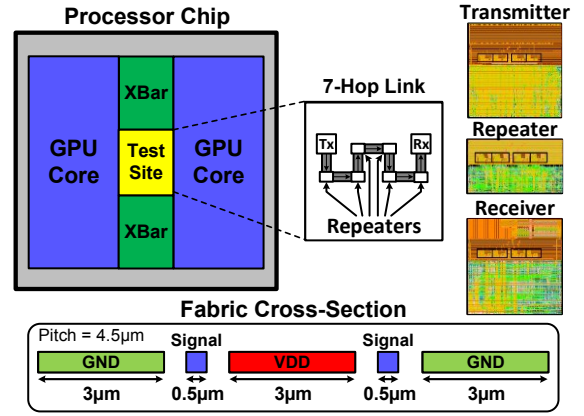


Fig. 10. On-chip serial link diagram, die images, and fabric cross-section.

ratios. The amplifier achieves large DC gains due to the SE-to-differential signal conversion, which attenuates input referred offset and removes the need for a SE-to-DIFF converter within the clock buffer path. The input referred offset can be trimmed through slewing of the bias currents in current sources M4, M5 (each an array of switched PMOS devices) and by tuning  $R_{BOT}$ .

## VIII. EXPERIMENTAL GRS LINKS

GRS is a general method of transmitting a single-ended signal across a channel for both intra-chip and inter-chip serial links. This section summarizes GRS test articles that were developed and experimentally verified to showcase the operation of GRS within serial links communicating over an on-chip communication network, chip-to-chip communication for chips mounted on the same organic package, and an off-package link for systems communicating over a short PC-board channel. All links utilize the clock-forwarding scheme described in Section V.

### A. 16Gbps On-Chip Serial Link

An on-chip communication link was implemented in a 28nm CMOS process using low-swing GRS between re-timing stages (hops). To experimentally validate the voltage sensitivity of the link, the test-site was implemented in a 600mm<sup>2</sup> production graphics-processing unit (GPU), as shown in Fig. 10. Operation of the GPU cores induces 155mV peak-to-peak noise on the supply, which is shared with the GRS-link. The link consists of 7 re-timing hops, where inclusion of PRBS checking-logic at each hop allows for link characterization across a programmable number of hops.

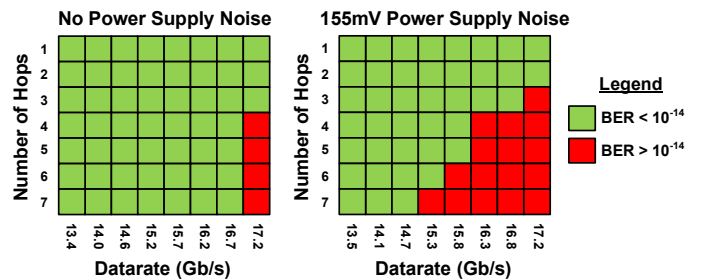


Fig. 11. Number of hops vs. datarate Shmoo plots.



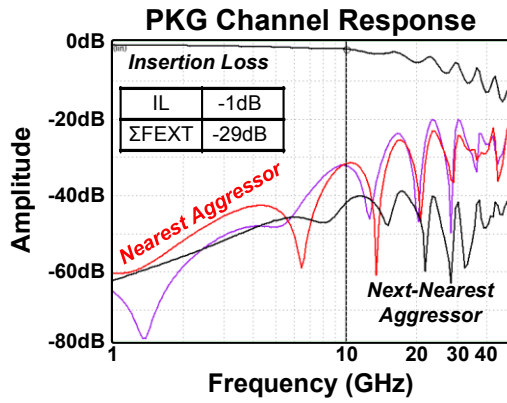


Fig. 12. Simulated HFSS organic package channel model.

The signal fabric was implemented using the 0.5 $\mu$ m-thick, on-chip upper metal layers, allowing for signal distribution without disturbing the underlying GPU compute logic. Signaling wires were interleaved between the existing GPU power-grid, thereby utilizing the VDD and GND distribution metals as shields between adjacent signals to mitigate cross-talk within the single-ended link. The signal wires were sized with 0.5 $\mu$ m width (1 $\mu$ m spacing) to implement a square wire cross-section. The characteristics of the on-chip metallization are described in Section IV(A). A target 16Gb/s signaling rate was selected for 1.5mm trace lengths, where trimmable transmit equalization up to 14dB was included to overcome line attenuation. Since the GRS link operates in a high-Z environment, 40fF storage capacitance was sufficient to generate 115mV DC amplitude at 16Gb/s. The single-ended common-gate structure, described in Section VII(A), was used as the receiver amplifier.

Fig. 11 shows measured signaling distance vs. data-rate shmoo plots of the link with and without GPU-induced supply noise. The link signaling rate was swept across a 13.5Gb/s to 17.2Gb/s range for 1-7 hops (1.5mm to 10.5mm transmission distance). The link achieves a 170fJ/b/mm energy efficiency, where the GPU-induced noise on the supply shortens the link reach by 4 hops (6mm distance). The limiting factor of the on-chip link results from quadrature timing distortion between the forwarded I & Q clocks that accumulates across the fabric due to cross-talk. The experimental results for the 1-hop case in Fig. 11 can also serve as a proxy for signaling over silicon interposer

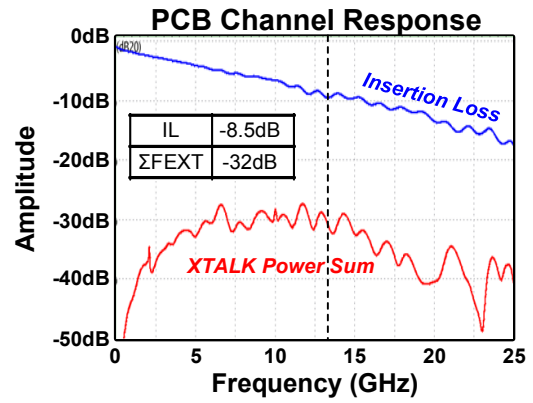


Fig. 13. Simulated HFSS printed circuit board channel model.

interconnects, which have similar RC-dominated line characteristics.

### B. 20Gbps On-Package Serial Link

An on-package, chip-to-chip GRS link was implemented on a conventional organic package for 20Gb/s operation [8]. The link was implemented in a 28nm planar CMOS process ( $V_{DD} = 0.92V$ ) and operates across LC-dominated channels with 50 $\Omega$  characteristic impedance. The charge-pump transmitter was sized with 200fF storage capacitance for 75mV DC output amplitude (150mV<sub>p-p</sub>). The receiver amplifier utilizes the single-ended common-gate structure described in Section VII(A). Clock source generation utilizes an LC oscillator, tunable over an 18GHz-22GHz range, with a quadrature clock divider to produce a 9GHz-11GHz half-rate clock.

The chips were attached to a custom organic package on a 4-2-4 stackup. Signal routes were restricted to the top build-up layers to use blind vias for vertical connections and thereby reduce stub impedances resulting from the plated-through hole (PTH) vias otherwise required to go below the PCB core layers. Power/Ground planes are interleaved between the 2 signal routing layers to implement stripline routes for reduced crosstalk between neighboring channels. Signals were routed on-chip using the top-layer RDL with equal 830 $\mu$ m trace lengths to ensure matched channel delays and RDL shunt capacitances. 6mm package routes (22 $\mu$ m width, 66.7 $\mu$ m height) were used to connect the separate die. A 1:1 signal-to-ground bump ratio was used in the bump pad array to ensure a robust signal return path.

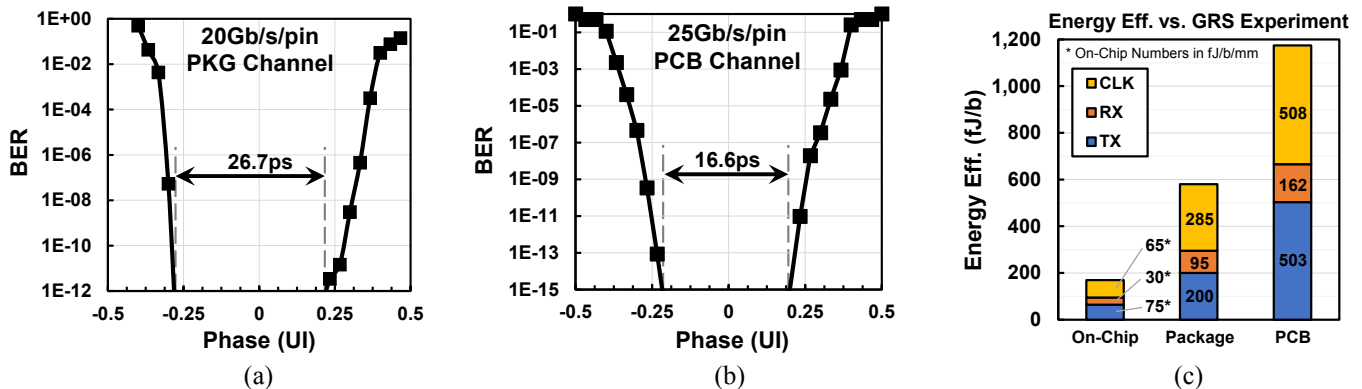


Fig. 14. Measured BER bathtub curves for (a) 20Gb/s on-package GRS and (b) 25Gb/s off-package GRS, and (c) Energy efficiency breakdown.

3-Dimensional channel modeling shows -1dB signal attenuation at Nyquist while the power sum of the crosstalk aggressors was -28dB below the signal insertion loss (Fig. 12). Most of the signal attenuation results from the  $\sim 100$ fF shunt capacitance at the transceiver outputs, comprising RDL trace, bump-pad, and circuit parasitics.

Fig. 14(a) shows the measured bathtub curve for the on-package link, which achieves a 0.53UI opening at a bit-error rate (BER) less than  $10^{-12}$  operating at 20Gb/s with a 0.58pJ/b energy efficiency. The 0.58pJ/b reported here uses a higher transmit EQ setting than the 0.54pJ/b reported in [8]. Fig. 14(c) shows the link efficiency breakdown using 4 segments of transmit EQ.

### C. 25Gbps Off-Package Serial Link

An off-package link was implemented using GRS for transmission across a conventional PCB channel [11]. The charge-pump transmitter uses a 300fF storage cap to generate 100mV DC output amplitudes ( $200\text{mV}_{\text{p-p}}$ ) at 25Gb/s. The link utilizes the single-ended to pseudo-differential receiver amplifier, described in Section VII(C), along with 5.8dB of transmit equalization. The link utilizes a PLL-regulated ring-oscillator (RO), locked to a 1.56GHz reference signal, for clock generation. The supply voltage is regulated so the RO, and thus the I/O circuitry, operates at a fixed rate independent of PVT. Adjusting the supply voltage as a function of process corner and operating temperature flattens circuit performance and current consumption across PVT.

The link consists of two GRS transceiver chips mounted on an organic package and connected by 54mm PCB channels. The routed channels include RDL traces and 13mm of on-package traces. Crosstalk was minimized by routing signals as striplines between ground planes for package and board traces. The signals were routed on the bottom-most stripline routing layer to minimize PTH via stub length and thus reduce signal reflections. 3D channel modeling showed -8.5dB insertion loss at Nyquist while the power sum of the crosstalk aggressors was -31.6dB, as shown in Fig. 13.

Fig. 14(b) shows the aggregate bathtub curve, which achieves 0.42UI opening at 25Gb/s. Also shown in Fig. 14(c) is the off-package link efficiency breakdown, which achieves 1.17pJ/b energy-efficiency.

## IX. CONCLUSION

Ground-referenced signaling enables short-reach, single-ended links to operate with high energy- and pin-efficiency. A series of experiments exploring the use of GRS has been demonstrated for both intra-chip and chip-to-chip interconnects. The single-ended links can operate on the same die or between chips mounted on the same organic package. The methodology also supports reliable high-speed communication on short PCB package-to-package links. In addition to the basics of the charge-pump GRS transmitter and ground-referenced level shifting receivers, simple delay-matched forwarded clocking architectures were demonstrated that avoid the need for a clock recovery system at the receiver. A comparison with prior work is provided in Table 1.

## REFERENCES

- [1] S.-K. Lee *et al.*, "A 95fJ/b Current-Mode Transceiver for 10mm On-Chip Interconnect," *Proc. ISSCC*, vol. 14.8 pp. 262-263, 2013.
- [2] E. Mensink *et al.*, "Power Efficient Communication Over Capacitively Driven RC-Limited On-Chip Interconnects," *IEEE JSSC*, vol. 45, no. 2, pp. 447-457, 2010.
- [3] J. Bae, J.-Y. Kim, and H.-J. Yoo, "A 0.6pJ/b 3Gb/s/ch Transceiver in 0.18um CMOS for 10mm On-chip Interconnects," *Proc. ISCAS*, pp. 2861-2864, 2008.
- [4] A. Shokrollahi *et al.*, "A Pin-Efficient 20.83Gb/s/wire 0.94pJ/bit Forwarded Clock CNRZ-5-Coded SerDes up to 12mm for MCM Packages in 28nm CMOS," *Proc ISSCC*, vol. 10.1, pp. 182-183, 2016.
- [5] W.-S. Choi *et al.*, "A 0.45-0.7V 1-6Gb/s 0.29-to-0.58pJ/b Source-Synchronous Transceiver Using Automatic Phase Calibration in 65nm CMOS," *Proc ISSCC*, vol. 3.8, pp. 66-67, 2015.
- [6] M. Mansuri *et al.*, "A Scalable 0.128-to-1Tb/s 0.8-to-2.6pJ/b 64-Lane Parallel I/O in 32nm CMOS," *Proc ISSCC*, vol. 23.2, pp. 402-403, 2013.
- [7] J. Song, S. Hwang, H.-W. Lee, and C. Kim, "A 1-V 10-Gb/s/pin Single-Ended Transceiver with Controllable Active-Inductor-Based Driver and Adaptively Calibrated Cascaded-Equalizer for Post-LPDDR4 Interfaces," *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. PP, no. 99, pp. 1-12, 2017.
- [8] J. Poulton *et al.*, "A 0.54pJ/b 20Gb/s Ground-Referenced Single-Ended Short-Haul Serial Link in 28nm CMOS for Advanced Packaging Applications," *IEEE JSSC* vol. 48, no. 12, pp. 3206-3218, 2013.
- [9] T. Knight and A. Krymm, "A Self-Terminating Low-Voltage Swing CMOS Output Driver," *IEEE JSSC* vol. 23, no. 2, pp. 457-464, 1988.
- [10] M. Johnson, "MOSFET Sense Amplifier Circuit," US Patent 4523110, Mostek Corporation, 1983.
- [11] J. Wilson *et al.*, "A 1.15pJ/b 25Gb/s/pin Ground-Referenced Single-Ended Serial Link for Off- and On-Package Communication in 16nm CMOS using a Process- and Temperature-Adaptive Voltage Regulator," *Proc ISSCC*, 2018.

Reference	[1]	[2]	[3]	GRS	[4]	GRS[8]	[5]	[6]	[7]	GRS[11]
Channel	On-Chip				On-Package		Off-Package			
	RC-Dominated Metal				Organic Substrate		SMA Cable	Ribbon Cable	Printed-Circuit Board	
Signaling	Diff	Diff	Diff	SE	CNRZ5	SE	SE	Diff	SE	SE
Tech.	65nm	90nm	180nm	28nm	28nm	28nm	65nm	32nm	65nm	16nm
Signal Rate	4Gb/s	2Gb/s	3Gb/s	16Gb/s	21Gb/s	20Gb/s	6Gb/s	16Gb/s	10Gb/s	25Gb/s
Reach	10mm	10mm	10mm	4.5mm	12mm	6mm	953mm	500mm	100mm	80mm
Energy Eff. [pJ/b]	0.01*	0.03*	0.06*	0.17*	0.94	0.58	0.58	2.58	4.18	1.17

\* On-chip numbers in pJ/b/mm

Table 1: Performance Summary and Comparison