



EFFICIENT DATA ENCODING FOR LOW ENERGY CONSUMPTION FOR NOC

1. ACHANTA NAGARJUNA, 2. K. SIRISHA

1. M. Tech Student, Dept. of ECE, BVC INSTITUTE OF TECHNOLOGY & SCIENCES, AMALAPURAM, A.P
2. Associate Professor, Dept. of ECE, BVC INSTITUTE OF TECHNOLOGY & SCIENCES, AMALAPURAM, A.P

ABSTRACT: The network on chip (NOC) is a widely discussed concept for handling the large on chip communication requirements of complex system on chip (SOC) design. The traditional bus based architecture does not communicate properly in very large SOCs. As a result the on chip communication uses the packet switching paradigm for routing information between the intellectual property (IP) blocks. The concept of code division multiple access (CDMA) is applied for on chip packet switch communication network. The technique of applying CDMA principle in NOC design is the point to be discussed in this project. A packet switched network on chip that applies the CDMA principle is realizable in a very common logic that is Register Transfer Logic (RTL) by using the VHDL coding technique. The globally asynchronous and locally synchronous (GALS) scheme is used for the realization of CDMA NOC by using both synchronous and asynchronous designing technology. Packet switched NOC is divided into two designing schemes which are named as CDMA NOC and POINT TO POINT NOC. The packet switch NOC which uses point to point design scheme, which is shown by the example of ring topology NOC, has a varying data transfer latency when the packets are transferred to different destination or to the same destination by different routes in the network. For the elimination of variation of data transfer logic CDMA NOC is used. The structure of the CDMA NOC is proposed and the process is coded and implemented by using ALTIUM software in this project. The model of CDMA NOC is described by the ALTIUM software. The comparative study of the characteristics of CDMA NOC and point to point NOC mainly ring topology are examined.

KEYWORDS: globally asynchronous and locally synchronous (GALS), network on chip (NOC), Register Transfer Logic (RTL)

INTRODUCTION:

Advanced silicon technology offers the possibility of integrating hundreds of millions of transistors into a single chip, which makes system-on-chip (SoC) design possible. With the continuous scaling of silicon technology, area and power dissipation of interconnects are one of the main bottlenecks for both on-chip and off-chip buses. Multiplexing parallel buses into a serial link enables an improvement in terms of reducing interconnect area, coupling capacitance, and crosstalk, but it may increase the overall switching activity factor (AF) and energy dissipation. Therefore, an efficient coding method that reduces the switching AF is important issues in serial interconnect design. The embedded transition inversion (ETI) coding scheme reduce the switching activity factor and solve the extra bit indication of TIC coding scheme by embedding the inversion information in the phase difference between the clock and the encoded data.

When there is an inversion in the data word, a phase difference is generated between the clock and data. Otherwise, the data word remains unchanged and there is no phase difference between the clock and the data. Need of phase difference is decided by the decision bit. The decision bit is decided by the number of transitions when the number of transitions is more than half of the word length, every even bit of the data word will be inverted and then decision bit sets to high. This operation is performed at transmitted section. The receiver side adopts a phase detector (PD) to detect whether the received data word has been encoded or not. Statistical analysis and experimental results show that the proposed coding scheme has low transitions for different kinds of data patterns. Low power design, in a system perspective, happens at all levels of the digital electronic system stack. It is being done from the lowermost device level design to the top most software

design. And there are the intermediate levels where a lot of effort is being expended to make systems run at low power.

Self-Transitions and Coupling Transitions: Self-transitions are defined as transitions on the capacitance between a bus line and the substrate (ground) while coupling transitions are defined as transitions on the capacitance between adjacent lines [4]. Figure 1 shows a simplified bus model with coupling (ignoring all the resistances). C_s is the self-capacitance from each bus line to ground; C_c is the coupling-capacitance between two adjacent lines. There has been some confusion in the literature about the difference between power consumption and power dissipation on buses with coupling, so here we explain the difference and give several examples

LITERATURE REVIEW

Parallel buses multiplexed into a serial link enables an improvement in terms of reducing interconnect area, coupling capacitance, and crosstalk, but it increases the overall switching Activity Factor (AF) and energy dissipation. Therefore, an efficient coding method needed to reduce the switching AF is an important issue in serial interconnect design. Many studies attempt to reduce the AF of parallel buses. Stan and Burleson introduced a bus-invert method that transmits the original or inverted pattern to minimize the switching activity. Researchers have proposed many techniques to improve the bus-invert coding method, such as the partial bus-invert coding and weight-based bus invert coding methods. The schemes mentioned above use an extra channel to send the inversion indication signal. Kuo et al. proposed the serial coding technique to solve the extra channel problem. They append extra information bits to the back of the original data word. Although this approach resolves the area overhead problem, it increases data latency. Three level differential encoding is proposed for parallel bus to enable multiple drivers at the transmitter and to recycle the same current and reduce power consumption. Joint crosstalk avoidance code and error correction code are proposed to reduce the power in parallel bus. Huang et al. further proposed combining serializing bus with the joint crosstalk avoidance code and error correction code to reduce the power. Serialized low-energy transmission (SILENT) is a coding method used in reducing the switching activity for serial links. This approach encodes every single bit in the parallel bus using the XOR gate, and multiplexes the encoded parallel buses into a serial link. The XOR operation sets an adjacent bit with the same value to zero. The greater the correlation is, the more zeros the encoder produces this method is designed for data with strong correlation. NoC Power and Energy Researchers have

recently begun focusing on the energy and power in NoCs, which have been shown to be significant contributors to overall chip power and energy consumption [3, 4, 10, 11]. One effective way to reduce NoC power consumption is to reduce the amount of data sent over the network. To that extent, recent work has focused on compression at the cache and network levels [12, 13] as an effective power-reduction technique. Compression is complementary to our approach.

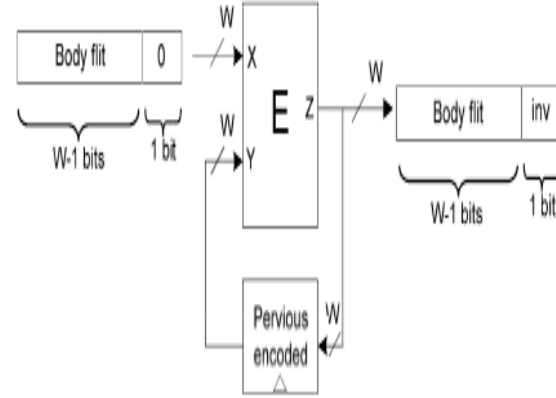


Fig: Encoder architecture

While our work seeks to reduce the amount of data transmitted through identification of useless words, compression could be used to more densely pack the remaining data. Researchers have also proposed a variety of techniques to reduce interconnect energy consumption through reduced voltage swing [14]. Schinkel et al. propose a scheme which uses a capacitive transmitter to lower the signal swing to 125mV without the use of an additional low-voltage power supply [15]. In this work we evaluate our prediction and packet encoding techniques for links composed of both full-signal swing as well as low-signal swing wires finally, static power consumption due to leakage currents is also a significant contributor to total system power. However, researchers have shown that power-gating techniques can be comprehensively applied at the NoC level and are highly selective at reducing leakage power at periods of low network activity. Our goal is to save dynamic energy in the memory system interconnect NoC. The proposed encoding architecture, which is based on the odd invert condition defined by (12), is shown in Fig. We consider a link width of w bits. If no encoding is used, the body flits are grouped in w bits by the NI and are transmitted via the link. In our approach, one bit of the link is used for the inversion bit, which indicates if the flit traversing the link has been inverted or not. More specifically, the

NI packs the body flits in $w-1$ bits. The encoding logic E , which is integrated into the NI, is responsible for deciding if the inversion should take place and performing the inversion if needed. The generic block diagram is the same for all three encoding schemes proposed in this paper and only the block E is different for the schemes. To make the decision, the previously encoded flit is compared with the current flit being transmitted. This latter, whose w bits are the concatenation of $w-1$ payload bits and a “0” bit, represents the first input of the encoder, while the previous encoded flit represents the second input of the encoder. The $w-1$ bits of the incoming (previous encoded) body flit are indicated by $X_i(Y_i)$, $i = 0, 1, \dots, w-2$. The w^{th} bit of the previously encoded body flit is indicated by inv which shows if it was inverted ($\text{inv} = 1$) or left as it was ($\text{inv} = 0$). In the encoding logic, each T_y block takes the two adjacent bits of the input flits (e.g., $X_1X_2Y_1Y_2$, $X_2X_3Y_2Y_3$, $X_3X_4Y_3Y_4$, etc.) and sets its output to “1” if any of the transition types of T_y is detected. This means that the odd inverting for this pair of bits leads to the reduction of the link power dissipation (Table I). The T_y block may be implemented using a simple circuit. The second stage of the encoder, which is a majority voter block, determines if the condition given in (12) is satisfied (a higher number of 1s in the input of the block compared to 0s). If this condition is satisfied, in the last stage, the inversion is performed on odd bits. The decoder circuit simply inverts the received flit when the inversion bit is high. To this end we developed a simple, low complexity, spatial locality predictor, which identifies the words expected to be used in each cache block. A used word prediction is made on a L1 cache miss, before the request packet to the L2 is generated. The basic idea of the proposed approach is encoding the flits before they are injected into the network with the goal of minimizing the self-switching activity and the coupling switching activity in the links traversed by the flits. In fact, self-switching activity and coupling switching activity are responsible for link power dissipation. In this paper, we refer to the end-to-end scheme. This end to end encoding technique takes advantage of the pipeline nature of the wormhole switching technique [4]. Note that since the same sequence of flits passes through all the links of the routing path, the encoding decision taken at the NI may provide the same power saving for all the links. For the proposed scheme, an encoder and a decoder block are added to the NI. Except for the header flit, the encoder encodes the outgoing flits of the packet such that the power dissipated by the inter-router point-to-point link is minimized [2].

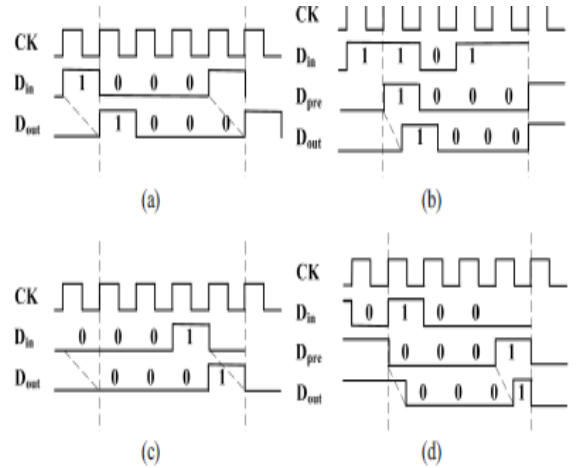
PROPOSED TECHNIQUE:

word exceeds the threshold N_{th} , the bits in the data word should be encoded. Otherwise, the data word remains the same. When an encoding is needed in a data word, this method checks every two-bit in the data word, as Fig. 3 shows. Every two bit in the serial stream is combined as a base to be encoded. In this case, the $b_{11}b_{21}$ is a base and the $b_{31}b_{41}$ is another base. The 2-bit in a base is denoted as b_1b_2 and the encoded output is denoted as $b_{e1}b_{e2}$. When the N_t in a data word is less than N_{th} , b_1b_2 remains unchanged. Otherwise, we perform the inversion coding and the phase coding. For the inversion coding, the bitstreams “01” and “10” are mapped to “00” and “11,” respectively. The bitstreams “00” and “11” are mapped to “01” and “10,” respectively. For the phase coding, we embed the inversion information in the phase difference between the clock and the encoded data.

$$b_{e1} = b_1$$

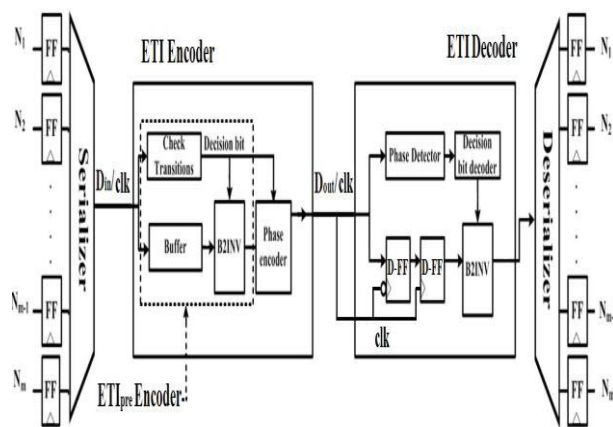
$$b_{e2} = \begin{cases} b_2, & \text{with } N_t < N_{\text{th}} \\ !b_2, & \text{with } N_t \geq N_{\text{th}}. \end{cases}$$

Since this operation is on a two-bit basis and only the second bit is inverted, it is called bit-two inversion (B2INV).



2) Phase Coding: The ETI coding uses the phase difference between the data and the clock to encode the indication information. The ETIpre has the same data word as the TIC, except that it removes the extra bit bex. Removing the bex leaves eight sets of data words that are exactly the same. For example, there are two “1000” data words after the ETIpre coding.

Within every data word duration, the phase difference between the data and the clock distinguishes these two data words, as Fig. 4 illustrates. Same Dout “1000” in If Din “1000” and “1101” without and with inversion. A half clock cycle difference between Dout and CK is shown in Fig. (b), indicating that Din has been encoded. The Dout and CK are aligned in Fig. (a), indicating that Din has not changed. Dout “0001” is the same in Fig. (c) and (d) from Din “0001” and “0100” without and with inversion. This approach is able to identify whether Dout has been encoded or not as long as there is a half cycle delay between the Dout and CK. Although the phase difference



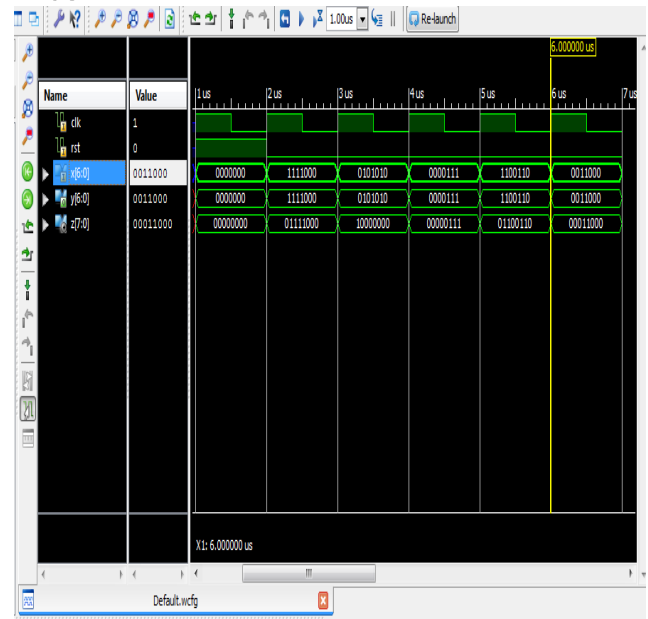
ETI Decoder

The ETI encoder generates the phase difference between the clock and the data word. Normally, a PD identifies an early or delayed phase. A variety of PDs could detect the phase difference. This paper adopts the commonly used Alexander PD [18]. The

Alexander PD architecture is shown in Fig. (a), which uses three consecutive clock edges to generate four sampling signals (S0, S1, S2, and S3). The PD is controlled by the clock CK and input data Din. When the clock CK and input data Din are valid, the PD is activated to identify the phase relation between the clock and the data. The PD can determine whether a data transition exists from the condition that the clock leads or lags the data. The basic waveform is shown in Fig. (b) to judge the un-inverted, inverted, no transition, or the special data word. If the clock leads the data (early conditions), the signal $S1 \oplus S2$ is high and the $S2 \oplus S3$ is low. Conversely, if the clock lags the data (late conditions), the signal $S1 \oplus S2$ is low and $S2 \oplus S3$ is high. Thus, $S1 \oplus S2$ and $S2 \oplus S3$ could provide the clock and data relation are related.

$S1 \oplus S2$	$S2 \oplus S3$	Clock	Coding state
High	Low	Early	Has not been encoded
Low	Low	No transition	
Low	High	Late	Has been encoded
High	High	Special case	

RESULT:



CONCLUSION

In this paper, we have presented a set of new data encoding schemes aimed at reducing the power dissipated by the links of a NoC. In fact, links are responsible for a significant fraction of the overall power dissipated by the communication system. In addition, their contribution is expected to increase in future technology nodes. As compared to the previous encoding schemes proposed in the literature, the rationale behind the proposed schemes is to minimize not only the switching activity, but also (and in particular) the coupling switching activity which is mainly responsible for link power dissipation in the deep sub-micron meter technology regime. The proposed encoding schemes are agnostic with respect to the underlying NoC architecture in the sense that their application does not require any modification neither in the routers nor in the links. An extensive evaluation has been carried out to assess the impact of the encoder and decoder logic in the NI. The encoders implementing the proposed schemes have been assessed in terms of power dissipation and silicon area. The impacts on the performance, power, and energy metrics have been studied using a cycle- and bit accurate NoC simulator under both synthetic and real traffic scenarios. Overall, the application of the proposed encoding schemes allows savings up to 21.6% of power dissipation and reduces the error during transmission without any significant performance degradation.

REFERENCES

- [1] "International technology roadmap for semiconductors - interconnect," Semiconductor Industry Association, 2006.
- [2] S. R. Vangal, J. Howard, G. Ruhl, S. Dighe, H. Wilson, J. Tschanz, D. Finan, A. Singh, T. Jacob, S. Debra A. Lelewer and Daniel S. Hirschberg, Data Compression. 1987.
- [11] Xilinx. [www.xilinx.com. \[Online\]. http://www.xilinx.com/support/documentation/data_sheets/ds100.pdf](http://www.xilinx.com/support/documentation/data_sheets/ds100.pdf)
- [12] Ian D. L. Anderson and Mohammed A. S. Khalid Jason G. Tong, "Soft-Core Processors for Embedded Systems," in 18th International Conference on Microelectronics, 2006.
- [13] xilinx., http://www.xilinx.com/ise/embedded/emb_rtf_guide.pdf.
- [14] www.xilinx.com. [Online]. <http://www.xilinx.com/microblaze>.
- [15] www.xilinx.com. [Online]. http://www.xilinx.com/support/documentation/sw_manuals/edk92i-ctt.pdf.
- Jain, V. Erraguntla, C. Roberts, Y. Hoskote, N. Borkar, and S. Borkar, "An 80-tile sub-100-W TeraFLOPS processor in 65-nm CMOS," *IEEE Journal of Solid-State Circuits*, vol. 43, no. 1, pp. 29–41, Jan. 2008.
- [3] M. B. Taylor, J. Kim, J. Miller, D. Wentzlaff, F. Ghodrat, B. Greenwald, H. Hoffman, P. Johnson, J.-W. Lee, W. Lee, A. Ma, A. Saraf, M. Seneski, N. Shnidman, V. Strumpfen, M. Frank, S. Amarasinghe, and A. Agarwal, "The raw microprocessor: a computational fabric for software circuits and general-purpose programs," *IEEE Micro*, vol. 22, no. 2, pp. 25–35, 2002.
- [4] J. C. S. Palma, L. S. Indrusiak, F. G. Moraes, A. G. Ortiz, M. Glesner, and R. A. L. Reis, "Inserting data encoding techniques into NoC-based systems," in *IEEE Computer Society Annual Symposium on VLSI*, Mar. 2007, pp. 299–304.
- [5] A. Jantsch, R. Lauter, and A. Vitkowski, "Power analysis of link level and end-to-end data protection in networks on chip," in *IEEE International Symposium on Circuits and Systems*, vol. 2, May 2005, pp. 1770–1773.
- [6] M. R. Stan and W. P. Burleson, "Bus invert coding for low power I/O," *IEEE Transactions on Very Large Scale Integration Systems*, vol. 3, pp. 49–58, Mar. 1995.
- [7] K. W. Kim, K. H. Baek, N. Shanbhag, C. L. Liu, and S. M. Kang, "Coupling-driven signal encoding scheme for low-power interface design," in *IEEE/ACM International Conference on Computer-aided Design*, 2000, pp. 318–321.
- [8] L. M. Ni and P. K. McKinley, "A survey of wormhole routing techniques in direct networks," *IEEE Computer*, vol. 26, pp. 62–76, Feb. 1993.
- [9] D. Bertozzi and L. Benini, "Xpipes: a network-on-chip architecture for gigascale systems-on-chip," *IEEE Circuits and Systems Magazine*, vol. 4, no. 2, pp. 18–31, 2004.
- [16] Khalid Sayood, Introduction to Data Compression.
- [17] International Technology Roadmap for Semiconductors (ITRS) Working Group, \International Technology Roadmap for Semiconductors (ITRS), 2009 Edition. "http://www.itrs.net/Links/2009ITRS/Home2009.htm.
- [18] W. J. Dally and B. Towles, \Route Packets, Not Wires: On-Chip Interconnection Networks," in The 38th International Design Automation Conference (DAC), 2001.
- [19] Y. Hoskote, S. Vangal, A. Singh, N. Borkar, and S. Borkar, \A5-GHz Mesh Interconnect for a Teraprocessor," *IEEE Micro*, vol. 27, 2007.