## 13.4  A 48GB 16-High 1280GB/s HBM3E DRAM with All-Around Power TSV and a 6-Phase RDQS Scheme for TSV Area Optimization

Jinhyung Lee*, Kyungjun Cho*, Chang Kwon Lee, Yeonho Lee, Jae-Hyung Park, Su-Hyun Oh, Yucheon Ju, Chunseok Jeong, Ho Sung Cho, Jaeseung Lee, Tae-Sik Yun, Jin Hee Cho, Sangmuk Oh, Junil Moon, Young-Jun Park, Hong-Seok Choi, In-Keun Kim, Seung Min Yang, Sun-Yeol Kim, Jaemin Jang, Jinwook Kim, Seong-Hee Lee, Younghyun Jeon, Juhyung Park, Tae-Kyun Kim, Dongyoon Ka, Sanghoon Oh, Jinse Kim, Junyeol Jeon, Seonhong Kim, Kyeong Tae Kim, Taeho Kim, Hyeonjin Yang, Dongju Yang, Minseop Lee, Heewoong Song, Dongwook Jang, Junghyun Shin, Hyunsik Kim, Changki Baek, Hajun Jeong, Jongchan Yoon, Seung-Kyun Lim, Kyo Yun Lee, Young Jun Koo, Myeong-Jae Park, Joohwan Cho, Jonghwan Kim

SK hynix Semiconductor, Icheon, Korea
*Equally Credited Authors (ECAs)

With the emergence of large-language models (LLM) and generative AI, which require an enormous amount of model parameters, the required memory bandwidth and capacity for high-end systems is on an unprecedented increase. To meet this need, we present an extended version of the high-bandwidth memory-3 (HBM3 DRAM), HBM3E, which achieves a 1280GB/s bandwidth with a cube density of 48GB. New design schemes and features, such as all-around power through-silicon via (TSV), a 6-phase read-data-strobe (RDQS) scheme, a byte-mapping swap scheme, and a voltage-drift compensator for write data strobe (WDQS), are implemented to achieve extended bandwidth and capacity with enhanced reliability. The overall architecture and specifications, such as bump map footprint, the number of channel and I/Os, and the operation voltage, are identical to the latest HBM3 [1,2]; therefore, backward compatibility is provided, avoiding system modification.

HBM3E DRAM is a 3D-stacked structure, which has a base die at the bottom and a stack of core dies, up to 16 slices, by using TSVs. In addition, all slices can operate simultaneously, such as the refresh burst command, which consumes a huge amount of power. However, the core die receives its power from TSVs, which are mostly concentrated in the peripheral region; thus, the IR drop, due to the cell array, has a significant impact on performance degradation. Therefore, the power-distribution-network (PDN) design for the core die is of particular importance. To improve power integrity, power and ground TSVs are placed in all possible places: including the edge, center, and intermediate control region of the core die. As a result, the dynamic-voltage drop improves by up to 75%, with a 475% increase in the number of power TSVs, as shown in Fig. 13.4.1.

The past generations of HBM3 adopt a 4-phase RDQS scheme, and there are multiple sets of RDQS TSVs to support a rank-to-rank interleave operation under a $t_{CCDR}$ (required delay for consecutive read commands to different ranks) timing constraint. A conceptual view of a read operation and the timing diagram of RDQS TSV control for each rank are depicted in Fig. 13.4.2. In a 4-phase RDQS scheme, stack-ID (SID) signal and FIFO-out data strobes (FDQS), which generate RDQS and FIFO-out (F-out) signals for the read FIFO of a core die, are required for each data path. A 4-phase RDQS scheme, with multiple sets of RDQS TSVs, has a short signal path from FDQS to RDQS and a sufficient margin for the TSV-RDQS driver enable signal (RDQS_EN). However, multiple sets of TSVs inevitably cause an increase of the peripheral area, which is dominated by the number of signal TSVs.

To dramatically reduce the peripheral area for the memory capacity extension, a 6-phase RDQS scheme is proposed to reduce the number of TSVs. An area reduction could be realized by reducing the number of FDQS and RDQS TSVs in half. Based on the proposed RDQS scheme, only one set of FDQS TSVs is designed for a command path with an internal signal named PC, which is utilized to control pseudo-channel (PCH) 0 and 1. Moreover, RDQS TSVs, for each rank, are no longer required. Consequently, the peripheral height and the number of signal TSV can be reduced by 31% and 8%, compared to a 4-phase RDQS scheme. On the other hand, RDQS design difficulty increases due to the required 6-phase RDQS skew matching and RDQS_EN timing-margins.

The main reason to generate a 6-phase RDQS, from FDQS, is to control the enable timing of TSV drivers for RDQS. Figure 13.4.3 shows how RDQS_EN is generated with a 6-phase RDQS during the rank-to-rank interleave operation: RDQS_EN[n] is generated using RDQS[n] and RDQS[n+2] to guarantee the exact timing of the falling edge. Even with the proposed scheme, if core dies with different process corners are stacked in the same channel, the original enable timing margin $t_{M0}$ can be reduced to $t_{M1}$. Therefore, the calibration scheme to control the RDQS delay, due to the process variation of every core die, is also designed to ensure $t_{M0}$ as much as possible [1].

A read-enable (RD_EN) signal must be applied to generate a 6-phase RDQS. RD_EN is generated with the SID and PC signals at the command path. Then, 6-phase RDQS signals are generated from FDQS and RD_EN. Afterwards, 6-phase RDQSs are transferred to the data path and become the source signals to control RDQS_EN and F-out signal to read FIFO. At the RDQS TSV driver, it is important to evaluate the enable setup margin ($t_S$) and hold margin ($t_H$) because the relationship between RDQS and its enable signal is asynchronous.

Figure 13.4.4 shows proposed WDQS clock-distribution network (CDN) inside the DWORD of HBM3E. Each DWORD contains a WDQS CDN consisting of 6-skew groups, with 8 DQ I/Os per skew group. To deliver WDQS inputs through a total of 48 I/Os, several stages such as a WDQS buffer, internal driver, and repeater are involved, which can cause a delay drift due to supply noise. This delay degrades the sampling margin of DQ receiver (Rx) by compromising the pre-trained optimal DQ sampling time. The proposed voltage-drifted delay compensator (VDDC) is designed to reduce the impact of internal drift in WDQS. VDDC consists of a voltage-controlled bias-generator (VCBG) circuit that generates a bias current by tracking the supply voltage change in the reverse direction, and a VDDC circuit that effectively compensates for the delay without introducing WDQS duty degradation or a delay increase [3]. According to the resistance value ratio ($R_2/R_1$), VCBG adjusts the slope of the bias current that is reduced when the supply voltage rises. The ratio can adjust the required amount of compensation that varies depending on the characteristics of the WDQS CDN transistor type or the layout parasitic components. In addition, the VDDC circuit is implemented with a feedback equalizer form and the enable (EN) and feedback signals (FS_t, FS_c) are configured with logic gates to minimize the transistor stacking. The proposed scheme improves the voltage-drifted delay which represent the delay difference within the JEDEC I/O voltage ($V_{DDQ}$) range by 5.7×, on average, over various PVT corners.

As shown in Figure 13.4.5, the PHY bump map has a mirror symmetry between left and right channels [4], while the cell array has a shift symmetry to improve the DRAM yield. This mismatched structure causes the skew and the line length to increase between PHY-to-TSV interconnects in the right channel. To avoid the performance degradation caused by the above issues, the HBM3E employs a structure in which the byte order of the TSV is matched with that of the PHY and the byte order in the cell is swapped. Even though the location of the data accessing the memory cell is changed, the data being written and read from the system's perspective is identical. In addition, not only is the ECC correction capability within the symbol unit the same as before, but the data pattern within the octet unit (OCT) used for memory cell screening also remains unchanged.

Figure 13.4.6(top) shows the measured $t_{CK}$ Shmoo plot: the results show that the fabricated HBM3E DRAM achieves 10Gb/s/pin at 1.1V, which corresponds to the total bandwidth of 1280GB/s with 16-channel operation. A comparison table of key features between HBM3 and HBM3E is summarized in Figure 13.4.6(bottom). The proposed HBM3E DRAM accomplishes a 25% bandwidth improvement and a 2× maximum density increase compared to the previous HBM3 DRAM.

*References:*
[1] M.-J. Park et al., "A 192-Gb 12-High 896-GB/s HBM3 DRAM With a TSV Auto-Calibration Scheme and Machine-Learning-Based Layout Optimization," *IEEE JSSC*, vol. 58, no. 1, pp.256-269, Jan. 2023.
[2] Y. Ryu et al., "A 16 GB 1024 GB/s HBM3 DRAM With Source-Synchronized Bus Design and On-Die Error Control Scheme for Enhanced RAS Features," *IEEE JSSC*, vol. 58, no. 4, pp.1051-1061, Apr. 2023.
[3] M. Mansuri and C. -K. K. Yang, "A low-power adaptive bandwidth PLL and clock buffer with supply-noise compensation," *IEEE JSSC*, vol. 38, no. 11, pp. 1804-1812, Nov. 2003.
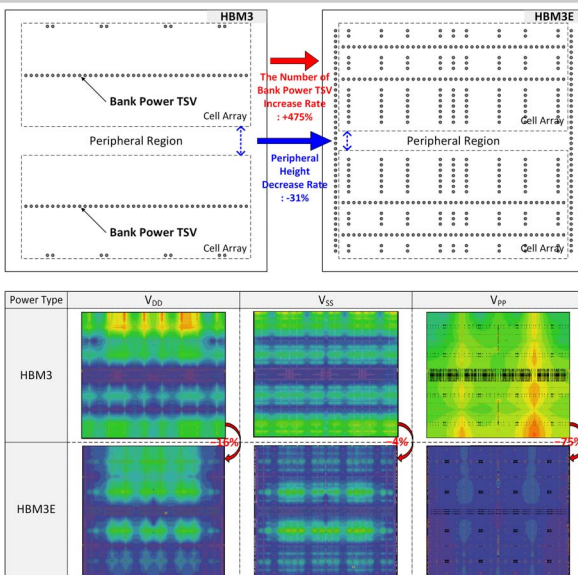[4] JESD238A: JEDEC Standard High Bandwidth Memory (HBM) DRAM Specification, Jan. 2023.

**Figure 13.4.1: Comparison of conceptual power TSV map (top) and dynamic voltage drop color map (bottom) between HBM3 and HBM3E.**
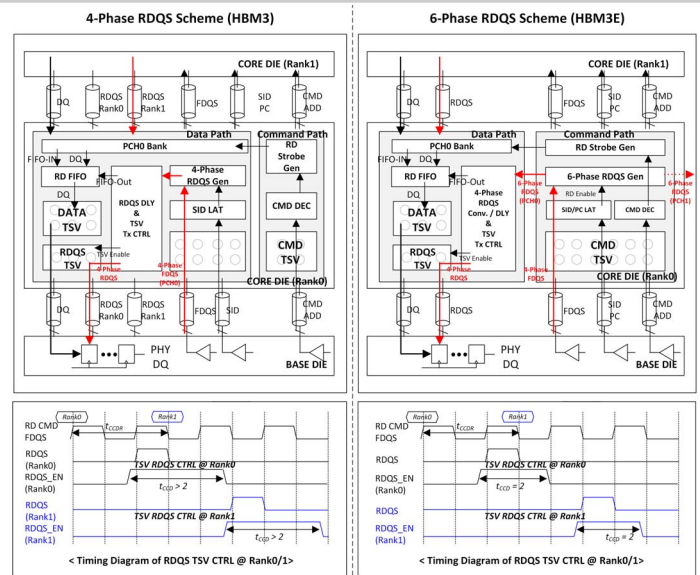


**Figure 13.4.2: Block diagram with signal flow at read operation and timing diagram in relation to 4-phase and 6-phase RDQS scheme.**
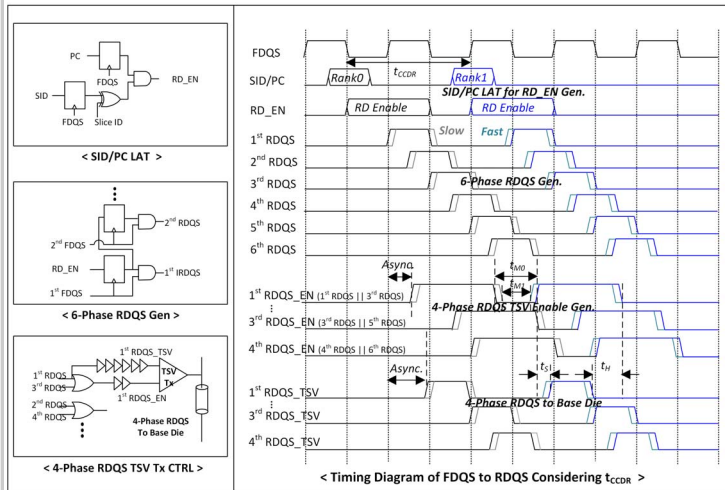
**13**



**Figure 13.4.3: Block and timing diagram of 6-phase generation to 4-phase RDQS transmission to base die.**
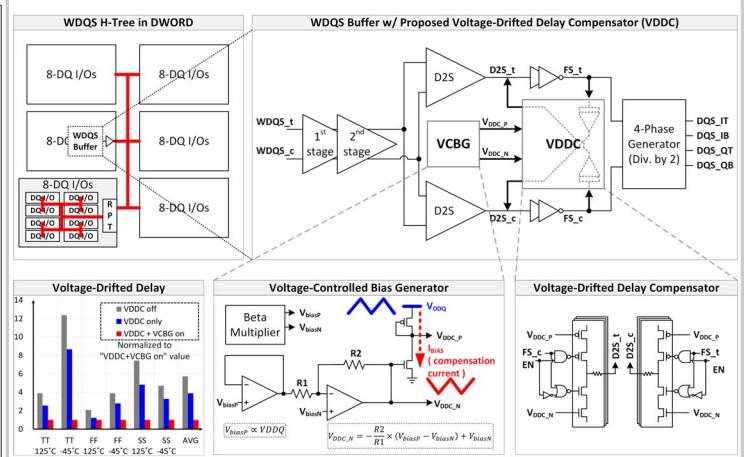


**Figure 13.4.4: Block diagram of voltage-drifted delay compensator (VDDC) and voltage-controlled bias generator (VCBG) for WDQS clock distribution network scheme.**
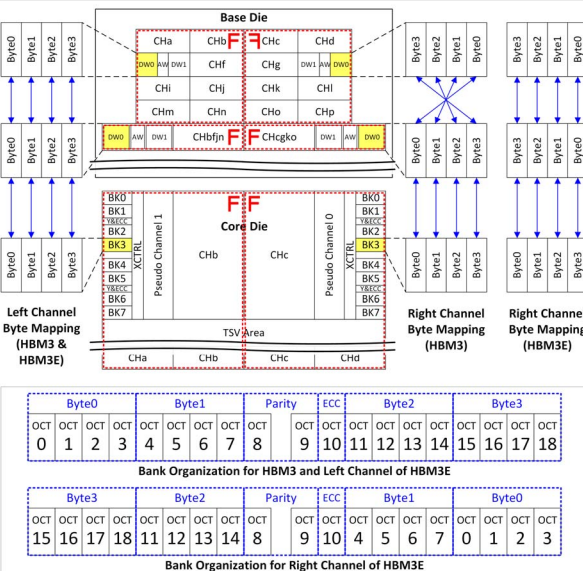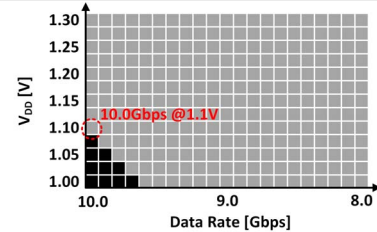


**Figure 13.4.5: Conceptual block diagram of byte-swap scheme and corresponding bank organization.**



**Figure 13.4.6: Measured $t_{CK}$ Shmoo plot and a table summarizing key features of HBM3 and HBM3E.**

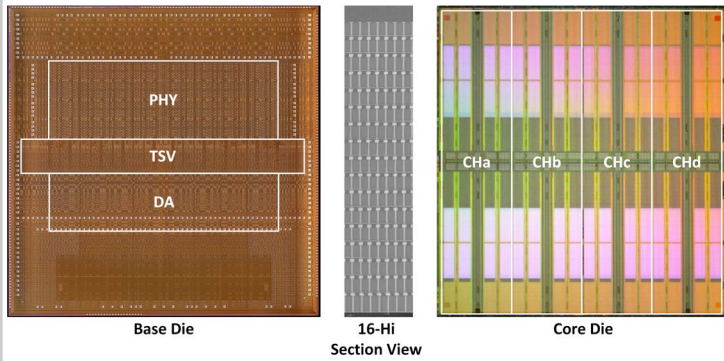| Generation | HBM3 | HBM3E |
|---|---|---|
| Supply Voltage | $V_{DD}$=1.1V, $V_{DDQ}$=1.1V, $V_{DDQL}$=0.4V, $V_{PP}$=1.8V | $V_{DD}$=1.1V, $V_{DDQ}$=1.1V, $V_{DDQL}$=0.4V, $V_{PP}$=1.8V |
| Data Rate | 7.0 Gb/s/pin*, 8.0 Gb/s/pin[†] | 10.0 Gb/s/pin |
| Bandwidth | 896 GB/s*, 1024 GB/s[†] | 1280 GB/s |
| Max. Density | 16 Gb × 12-High = 24 GB | 24 Gb × 16-High = 48 GB |
| Addresses | RA<0:13>, CA<0:4>, BA<0:3>, SID<0:1> | RA<0:14>, CA<0:4>, BA<0:3>, SID<0:1> |
| Organization | 16 channel × 2 PCH × 32 I/O | 16 channel × 2 PCH × 32 I/O |
| Microbump ballmap | 7.08 mm × 8.82 mm | 7.08 mm × 8.82 mm |
| Microbump pitch | 96 μm × 110 μm | 96 μm × 110 μm |
| Chip Size | 11 mm × 11 mm | 11 mm × 11 mm |

\* Ref.[1]  † Ref.[2]

**Figure 13.4.7: Chip micrograph of base die (left), a 16-high section view (middle), and core die (right).**