# Data Shaping for Improving Endurance and Reliability in Sub-20nm NAND

Eran Sharon, Stella Achtenberg, Idan Alrod, Avi Klein, Alon Eyal
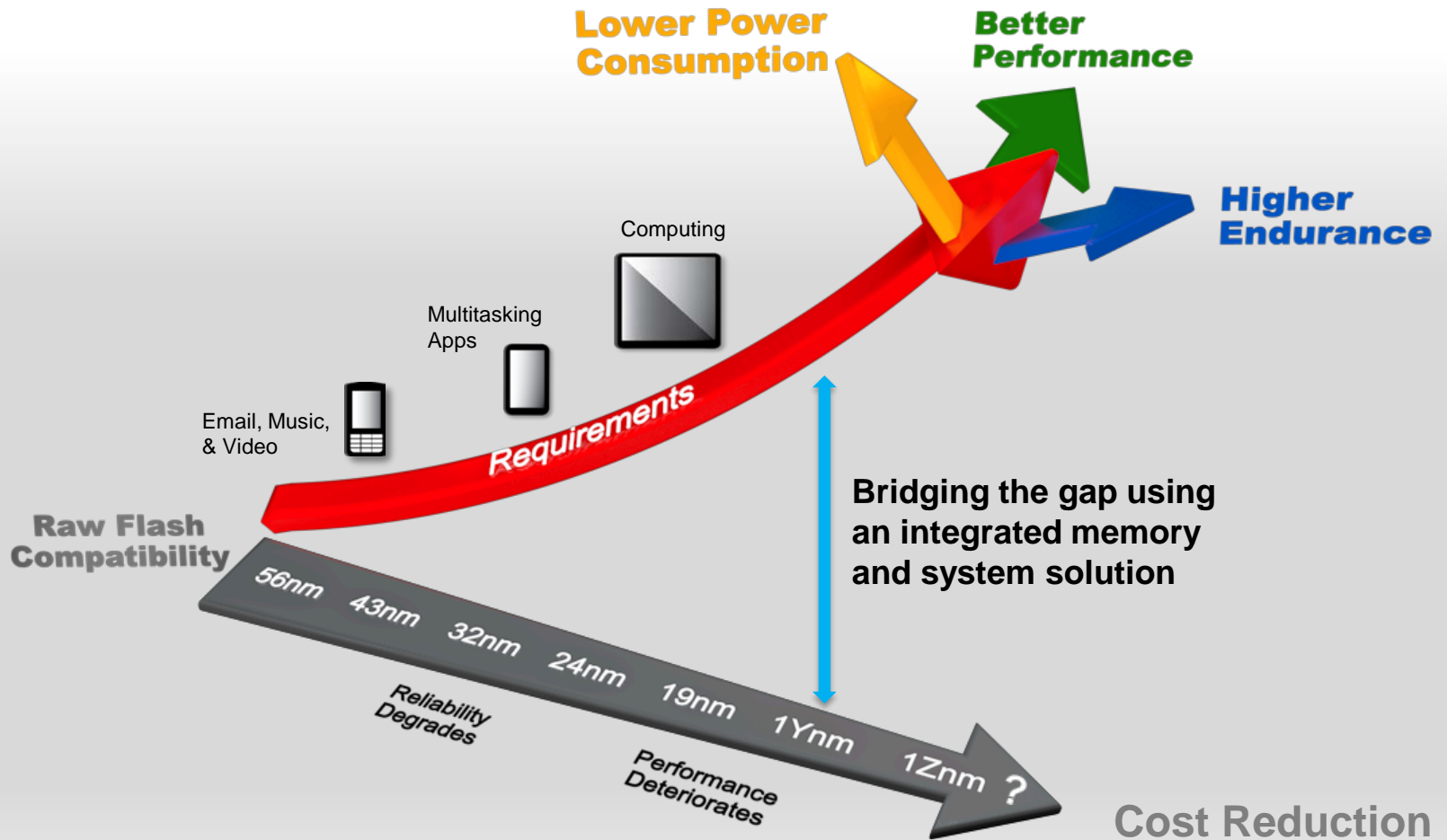*Intelligent Memory Systems, SanDisk Corp.*
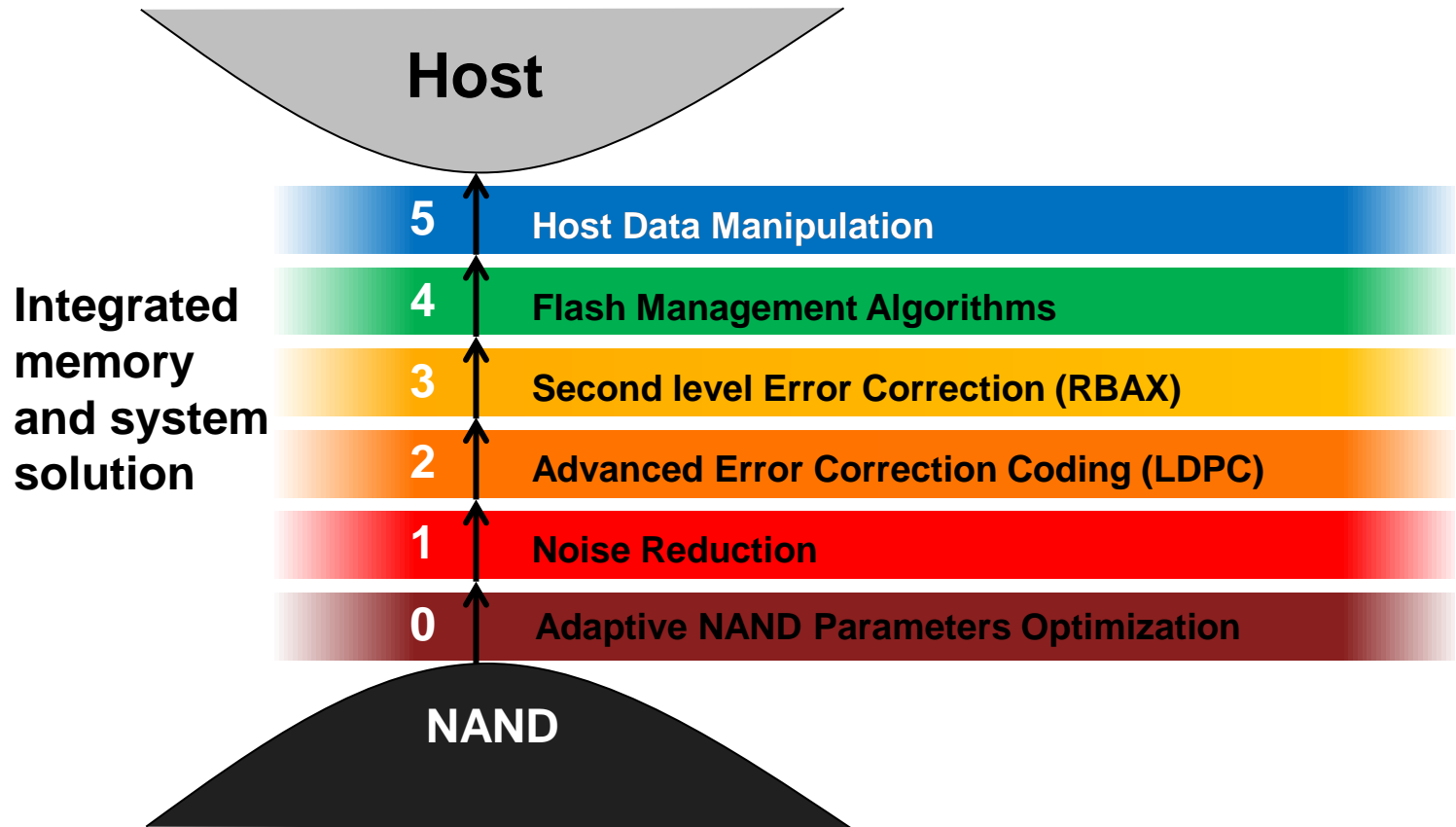
*August 2014*

# Outline

- Gap Between Product Requirements and Technology Capability

- Multi Tier Integrated System & Memory Solution - Recap

- Low Data Entropy in Typical Hosts – Unrealized Potential

- Leveraging Low Host Data Entropy for Data Shaping

- Data Shaping Effect on Memory Endurance and Reliability

- Practical Approaches for Data Shaping

- Eco-System Consideration

- Summary

*Disclaimer: This tutorial provides an overview of various techniques and concepts, some or all of which may not necessarily reflect what SanDisk is actually using in their products.*

# Gap Between Raw Memory Capability and Applications Requirements



**Lower Power Consumption**

**Better Performance**

**Higher Endurance**

Computing

Multitasking Apps

Email, Music, & Video

Requirements

Raw Flash Compatibility

56nm 43nm 32nm 24nm 19nm 1Ynm 1Znm ?

Reliability Degrades

Performance Deteriorates

**Bridging the gap using an integrated memory and system solution**

**Cost Reduction**

SanDisk®

# Integrated Memory & System Solution

# Low Data Entropy in Typical Hosts – Unrealized Potential

- Examination of typical hosts data traffic shows that significant fraction of the data is of low entropy, having many repetitive data patterns
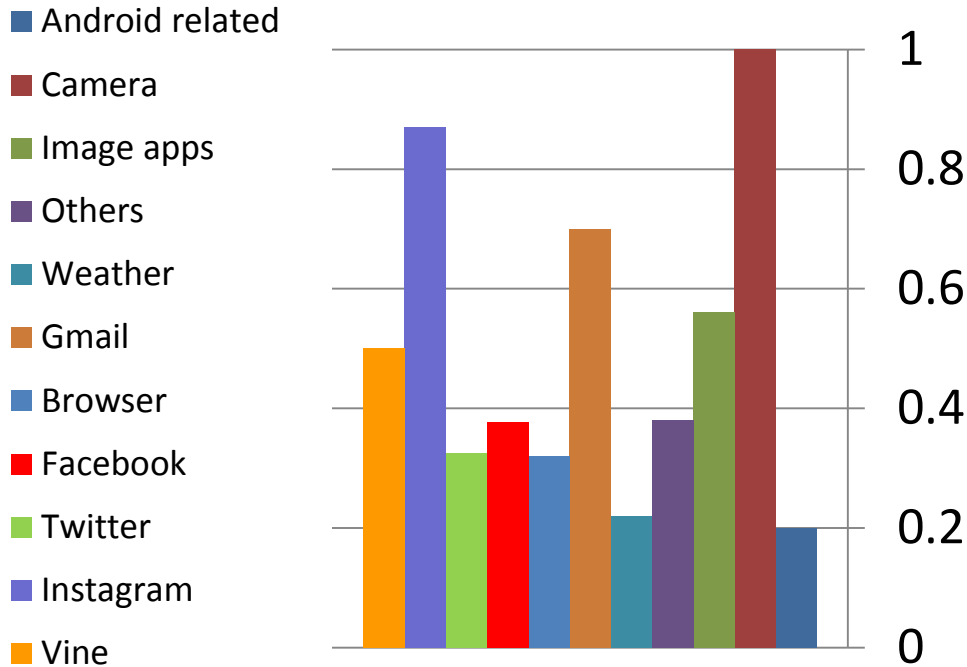
**Entropy** – measure of data randomness

- **Unrealized potential**: the inherent "redundancy" in the host data can be leveraged for improving endurance, reliability, performance and power, by manipulating host data:
  - **Deduplication**
  - **Compression**
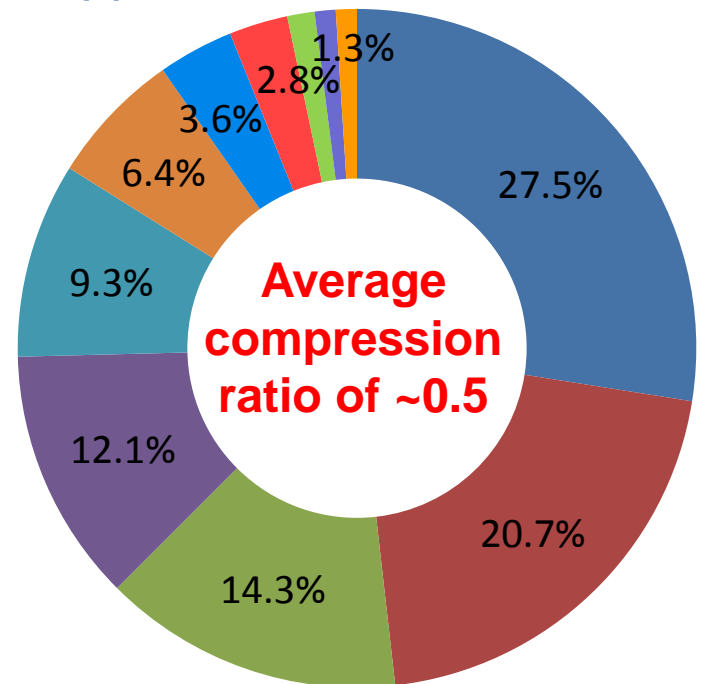  - **Data Shaping – a.k.a Endurance Coding**

# Analyzing mobile traffic of a sample user

- Record the traffic between the host and the controller during sample usage.
- Platform: Android 4.2.2. based Smart Phone
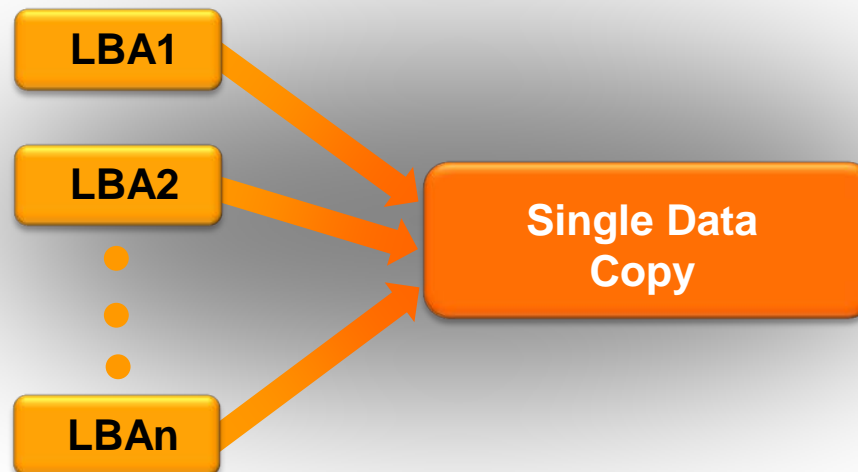- **Average compression ratio of ~50%**

## Compression ratio

- Android related
- Camera
- Image apps
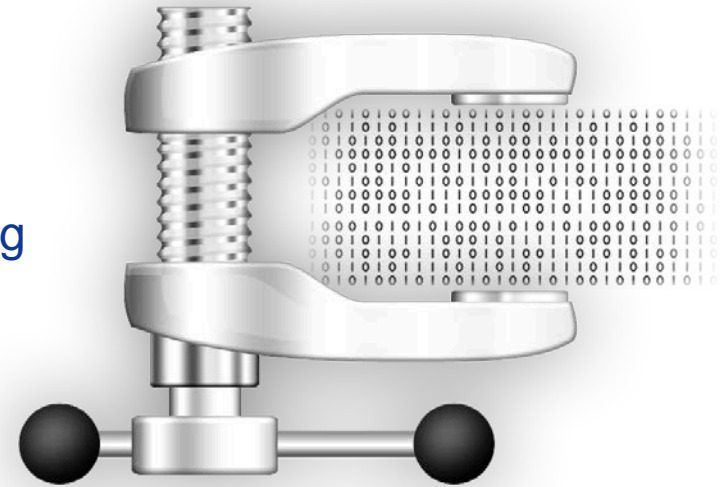- Others
- Weather
- Gmail
- Browser
- Facebook
- Twitter
- Instagram
- Vine

## Applications traffic distribution



Average compression ratio of ~0.5

- 27.5%
- 20.7%
- 14.3%
- 12.1%
- 9.3%
- 6.4%
- 3.6%
- 2.8%
- 1.3%

# Deduplication

- Specialized data compression technique for eliminating duplicate copies of repeating data
  - Manage multiple pointers to a single stored copy
- Operates on the file system level
  - Less suitable for eMMC level implementation (operates on 4KB sectors, unaware of files)
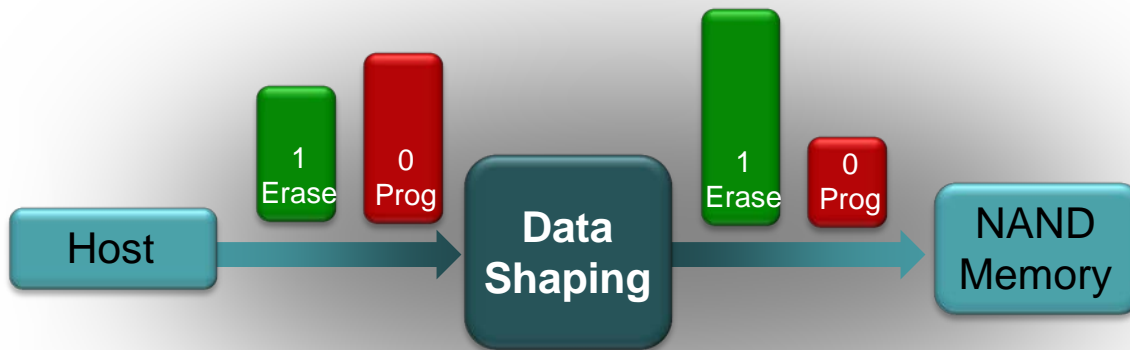- Highly suitable for enterprise backup applications

# Compression

- Typical traffic in Mobile applications is highly compressible

- Compression can provides significant endurance & performance gains
  - Less P/E cycles per GygaByte (GB) written
  - Increases the effective memory over provisioning
    - ➤ Improved garbage collection efficiency
    - ➤ Reduced write amplification
    - ➤ Performance stability

- System level considerations – impact on controller complexity power & cost
  - Requires significant changes in the Flash Management –
    mapping the logical address space into a variable physical space
  - High throughput, low power and cost compression engine design is challenging
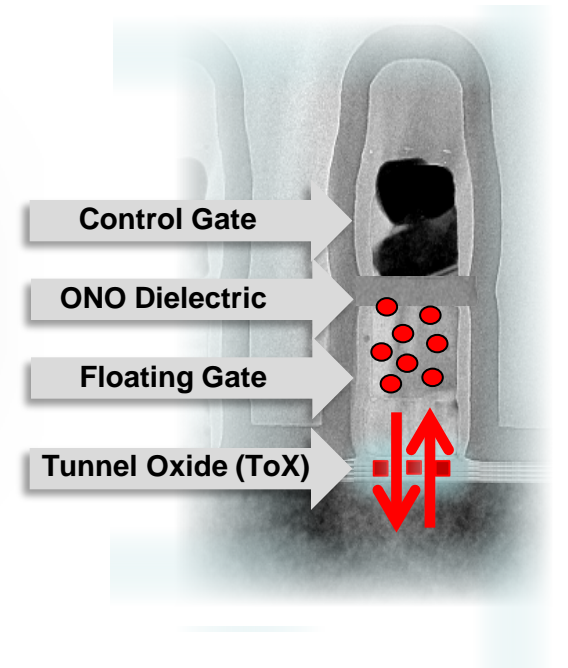
# Data Shaping ("Endurance Coding")

- **The Challenge**: Increasing Endurance

- **The Means**: Data Shaping – transform input data sequence into a "shaped" data sequence which induces less wearing when programmed to the NAND.

- **SLC Example**: transform the input data into shaped data having **less 0's**



**Minimize number of programmed cells per P/E cycle**

➤ Minimize average number of electrons tunneling in & out of the ToX per P/E cycle
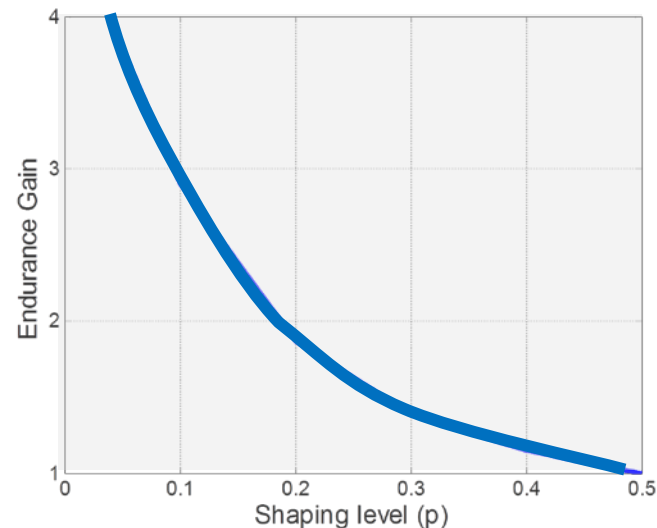
➤ Slow down ToX quality degradation

# Achievable Endurance Enhancement via Data Shaping

- Cell wearing is proportional to the probability $p$ of the cell to be programmed

- **Simplified model**:

  - $W_E$ – Wearing of an erased cell during a P/E cycle

  - $W_P$ – Wearing of a programmed cell during a P/E cycle

  - $W_P \gg W_E$   (Much more electrons passing through ToX for programmed cells)

  - Total wearing as a function of the shaping level $p$:  $W(p) = (1\text{-}p) \cdot W_E + p \cdot W_P$

  - Endurance gain due to using Shaped data ($p < 0.5$) vs. Scrambled data ($p = 0.5$):

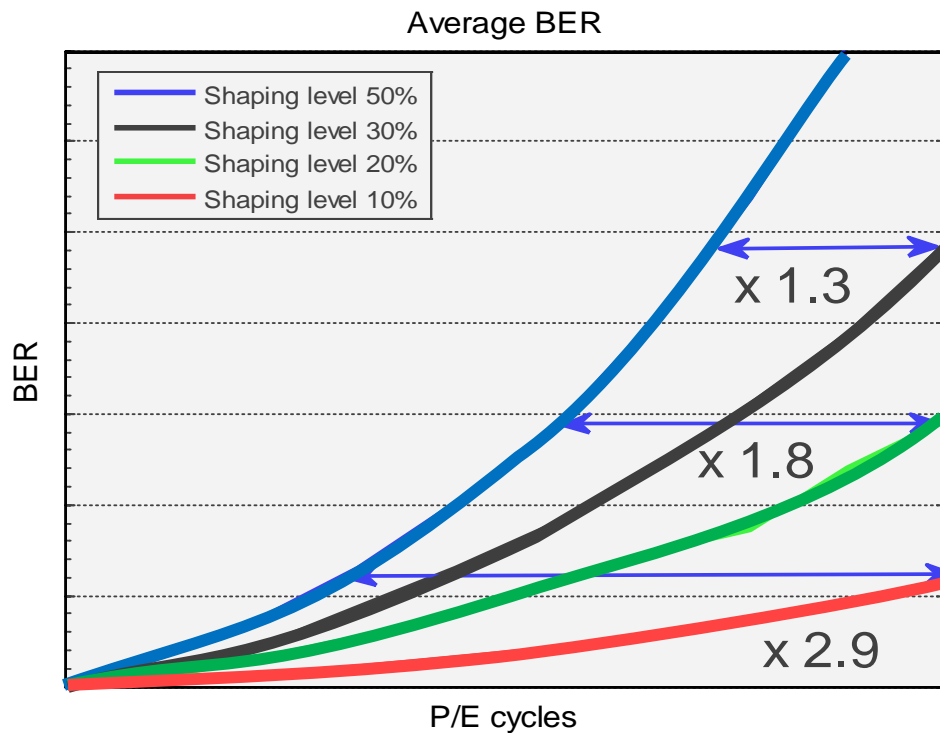$$Gain(p) = \frac{W(0.5)}{W(p)} = \frac{\alpha + 0.5}{\alpha + p},$$

where $\alpha = \dfrac{W_E}{W_P - W_E}$

($\alpha$ is specific per memory technology)



SanDisk®

- **Objective**: measure the cumulative wearing reduction effect of shaping

- **Experiment**:

  - Cycle the memory with shaped data (different shaping levels, up to different cycles)

  - At the last cycle, program with scrambled data ($p = 0.5$) and measure BER

  - Compare BER deterioration with cycling as a function of the average shaping level

Average BER



Measured endurance gain increases as the fraction of programmed cells ($p$) reduces

**SanDisk**®

# "Noise" Reduction due to Data Shaping - Empirical measurements

- **Objective**: measure the local BER reduction effect of shaping

- **Experiment**:
  - Cycle the memory with scrambled data
  - Program with shaped data at the last cycle and measure BER
  - Compare BER level at the last cycle as a function of the shaping level *p*.



Average BER vs. shaping level

# Effect of Shaping on Error Correction Capability

- Decoder correction capability can be significantly improved for shaped data
- Adjust decoder soft input metrics based on the estimated shaping level



Correction Capability for Different Shaping Levels

# Leveraging low host data entropy for shaping
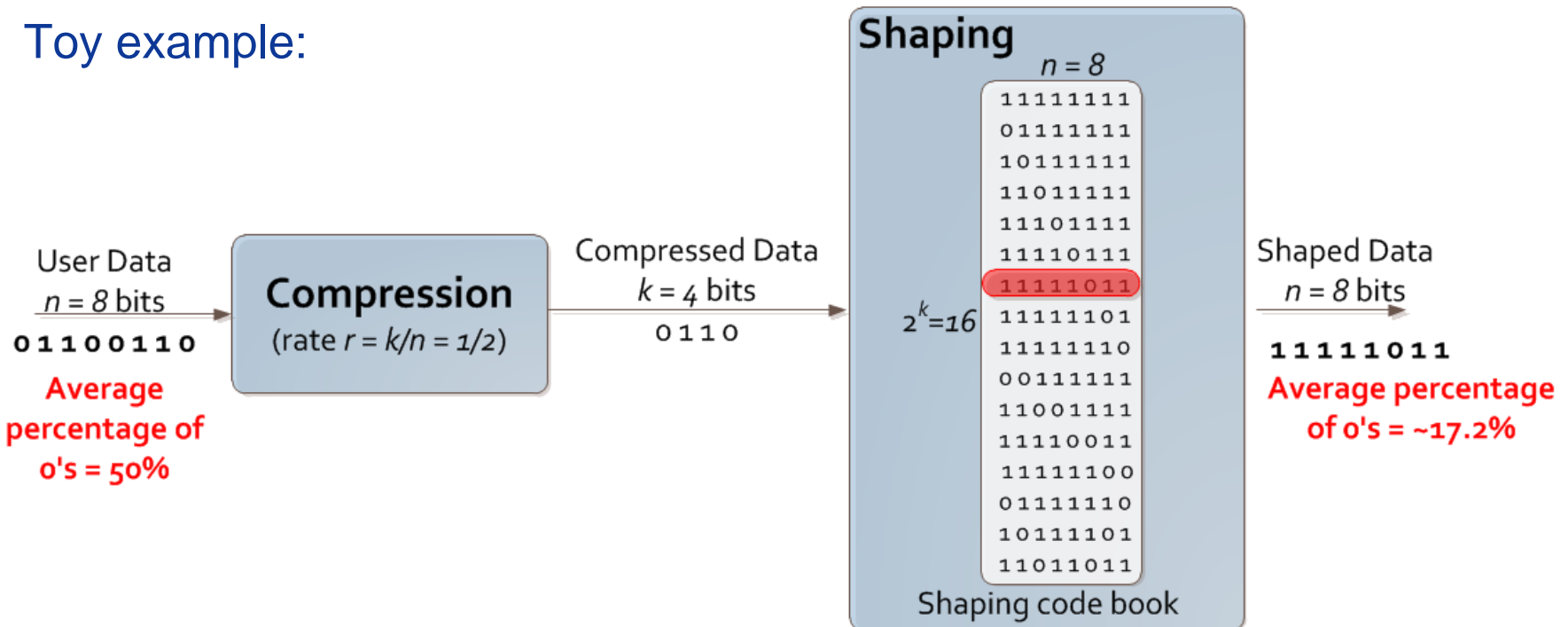
- Compression – Expansion approach
  - Compress: $n$ user bits into $k$ compressed bits
  - Expand: $k$ compressed bits into $n$ shaped bits

- Toy example:

# Achievable Shaping Level as a function of Data Compressibility

- What is the achievable shaping level of the Compression- Expansion approach?

- Assume data compression rate $r=k/n$, an optimal shaping code book will include all the $2^k$ length $n$ binary vectors having a minimal number of 0's $j$, up to at most $m$.

$$\Rightarrow 2^k = \sum_{j=0}^{m} \binom{n}{j} \underset{n \to \infty}{\cong} 2^{n \cdot Hb\left(\frac{m}{n}\right)} \quad \Rightarrow H_b\left(\frac{m}{n}\right) = \frac{k}{n} = r$$

where $H_b(p) = -p\log_2(p)-(1-p)\log_2(1-p)$ is the binary entropy function

$\Rightarrow$ The achievable shaping level of an optimal scheme is: $\quad \boldsymbol{p = \frac{m}{n} = H_b^{-1}(r)}$

# Endurance Enhancement Potential
# - Shaping Vs. Compression

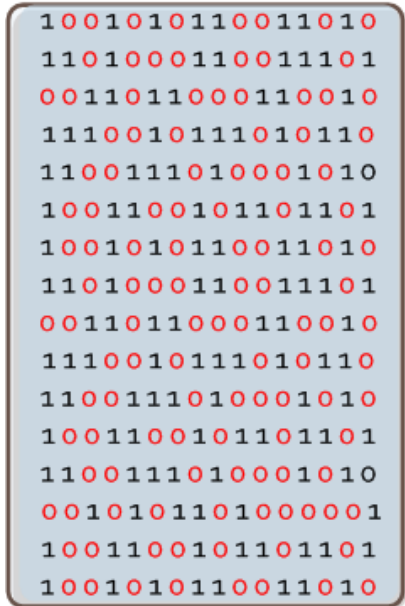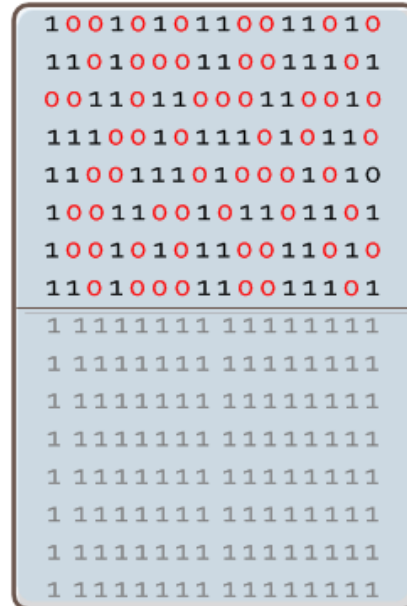- **Example:** Compression rate $r$ = ½ , achievable shaping level = $p = H_b^{-1}(r)$ = 0.11

**Reference system**
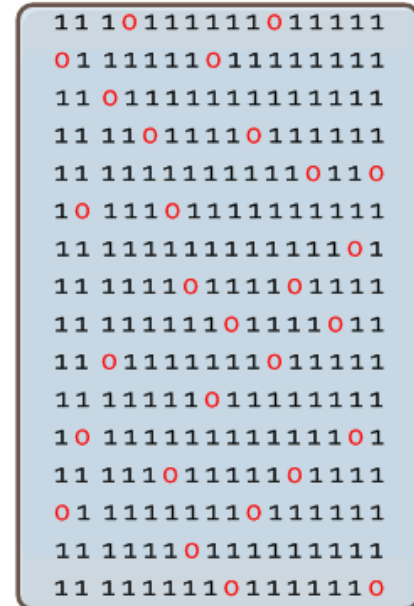Host data is scrambled
~50% of the cells are
programmed per GB

```
1001010110011010
1101000110011101
0011011000110010
1110010111010110
1100111010001010
1001100101101101
1001010110011010
1101000110011101
0011011000110010
1110010111010110
1100111010001010
1001100101101101
1100111010001010
0010101101000001
1001100101101101
1001010110011010
```

**System employing compression**
Host data is compressed
~25% of the cells are
programmed per GB

```
1001010110011010
1101000110011101
0011011000110010
1110010111010110
1100111010001010
1001100101101101
1001010110011010
1101000110011101
1 1111111 11111111
1 1111111 11111111
1 1111111 11111111
1 1111111 11111111
1 1111111 11111111
1 1111111 11111111
1 1111111 11111111
1 1111111 11111111
```

**System employing shaping**
Host data is shaped
~11% of the cells are
programmed per GB

```
11 101111111011111
01 11111011111111
11 01111111111111
11 11011110111111
11 11111111110110
10 11101111111111
11 1111111111101
11 11110111101111
11 1111110111011
11 01111111011111
11 11111011111111
10 11111111111101
11 11101111101111
01 11111110111111
11 11110111111111
11 11111101111110
```

**Potentially X2 the endurance**

(excluding indirect effects like write
amplification reduction)

**Potentially X2.8 the endurance**

(depending on the specific memory technology
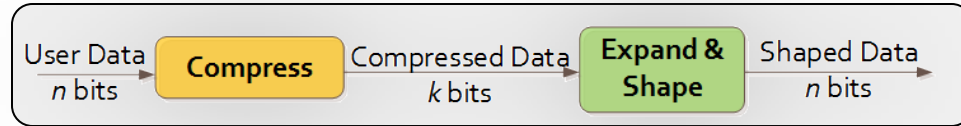and the shaping scheme optimality)

# Flash Management Implications - Shaping Vs. Compression

- **Impact on the Flash Management (Backend Firmware):**

  - **Compression** – Significant impact
    Converts $n$ bits to $k$ bits ($k < n$)
    Changes the logical to physical address management –
    Map logical sectors to variable physical sub-sectors

  - **Shaping** – Transparent to the Flash Management
    Converts $n$ bits to $n$ bits
    Logical to physical address mapping unchanged –
    Can be considered simply as a different type of scrambler…
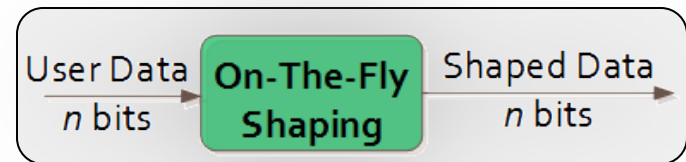
# Practical approaches for data shaping

- **Compression – Expansion approach:**



User Data n bits → Compress → Compressed Data k bits → Expand & Shape → Shaped Data n bits

- Two stage approach:
  - Compress using a lossless compression algorithm – e.g. LZ compression
  - Expand using a shaping code - e.g. Adaptive Reverse Huffman/Run-Length
- **Pros:** near optimal – can closely approach the theoretical shaping limit
- **Cons:** High complexity, High power consumption, Large latency (need to support variable compression-expansion rates)

- **Direct shaping approach:**



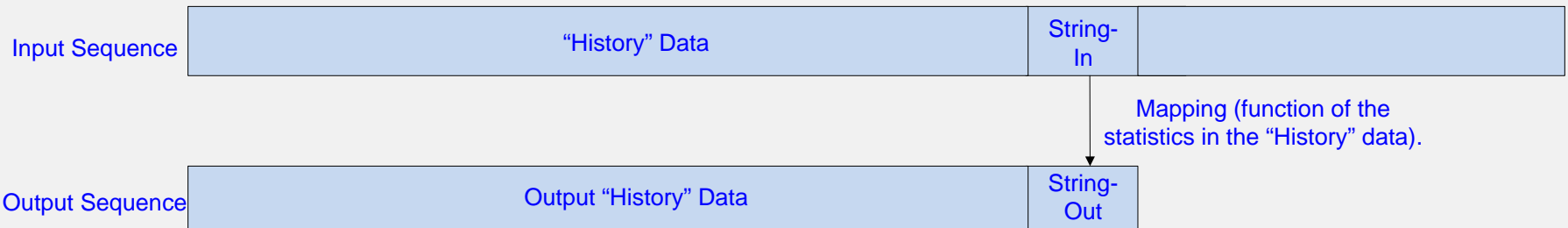User Data n bits → On-The-Fly Shaping → Shaped Data n bits

- Single stage approach: Direct shaping transformation from n (compressible) input bits to n(shaped) output bits
- **Pros:**
  - Negligible complexity
  - Negligible power consumption
  - Can be done On-The-Fly at very high throughputs
- **Cons:** Sub-optimal – achieves lower shaping level than theoretical limit

# Direct Shaping

- Transform *n* **compressible** bits into *n* **shaped** bits

- Convert each input string into a shaped output string using an adaptive mapping

- The mapping used for the current input string is a function of the statistics of previous strings, matching the most frequent "historic" strings to the most shaped strings

**Direct Shaping Scheme**

| Input Sequence | "History" Data | String-In | |
|---|---|---|---|

Mapping (function of the statistics in the "History" data).

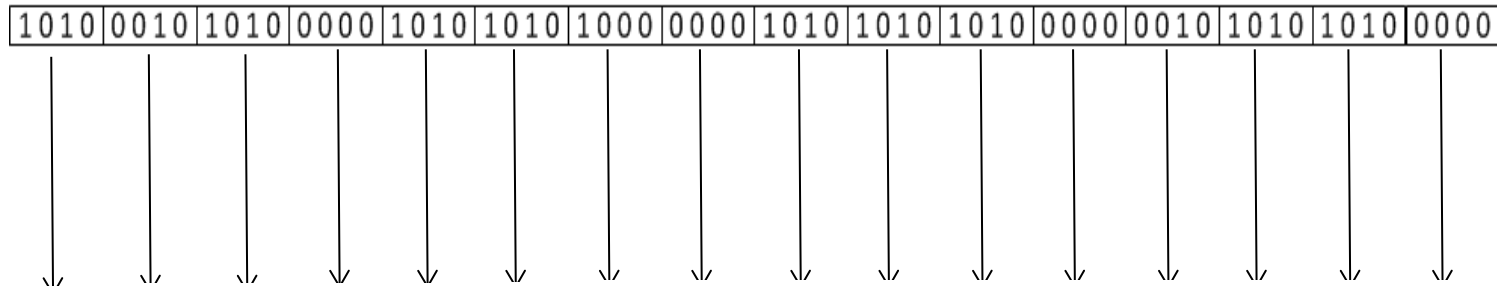| Output Sequence | Output "History" Data | String-Out |
|---|---|---|

- **Reversibility**: all mapping decisions are based on the "history" and hence can traced back by the De-Shaping algorithm → No need to store any side information

- **Amenable to an extremely slim design (few Kgates), negligible power consumption, OTF operation at high throughput**
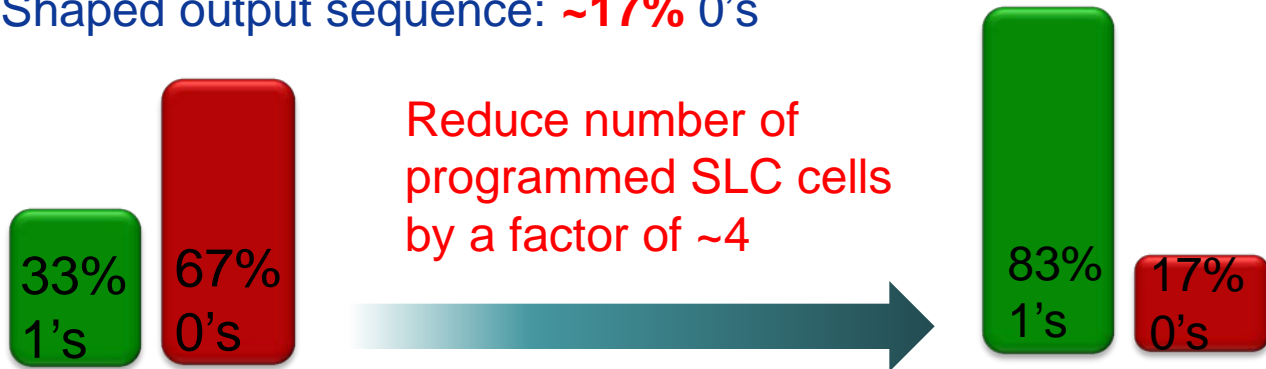
# Direct Data Shaping – How does it Work?

- **Toy example:**
- Convert a 64 bit compressible input sequence into a 64 bit shaped output sequence
- At step *j* map the most frequent 4 bit strings up to step *j*-1 to 4 bit strings with less 0's

| input | output | count |
|-------|--------|-------|
| 1 0 1 0 | 1 1 1 1 | 9 |
| 0 0 0 0 | 1 1 1 0 | 4 |
| 0 0 1 0 | 1 1 0 1 | 2 |
| 1 0 0 0 | 1 0 1 1 | 1 |
| 0 0 0 1 | 0 1 1 1 | 0 |
| 0 1 0 0 | 1 1 0 0 | 0 |
| 0 0 1 1 | 1 0 1 0 | 0 |
| 0 1 0 1 | 0 1 1 0 | 0 |
| 1 0 0 1 | 1 0 0 1 | 0 |
| 0 1 1 0 | 0 1 0 1 | 0 |
| 1 1 0 0 | 0 0 1 1 | 0 |
| 0 1 1 1 | 1 0 0 0 | 0 |
| 1 0 1 1 | 0 1 0 0 | 0 |
| 1 1 0 1 | 0 0 1 0 | 0 |
| 1 1 1 0 | 0 0 0 1 | 0 |
| 1 1 1 1 | 0 0 0 0 | 0 |

Compressible input sequence: **~67%** 0's

| 1010 | 0010 | 1010 | 0000 | 1010 | 1010 | 1000 | 0000 | 1010 | 1010 | 1010 | 0000 | 0010 | 1010 | 1010 | 0000 |

Shaped output sequence: **~17%** 0's

33% 1's    67% 0's

Reduce number of programmed SLC cells by a factor of ~4

83% 1's    17% 0's

# Shaping Advantages - Summary

1. **Reduced cell wearing** (main motivation)
2. **Less disturb effects**
3. **Higher ECC capability**

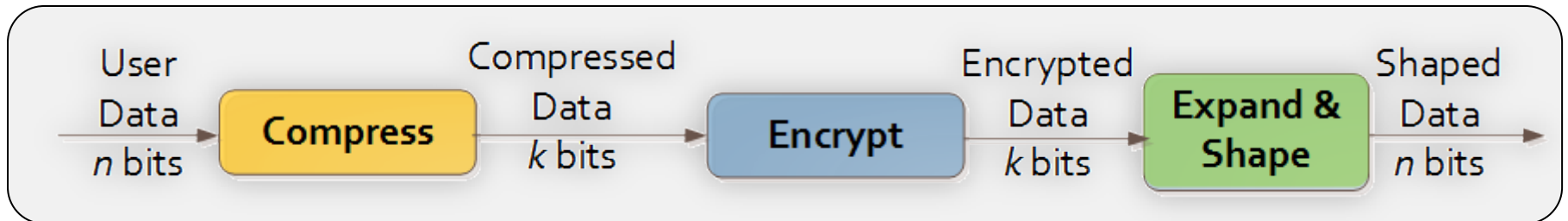**Orthogonal advantages**

**Shaping gain is threefold**

- First advantage is **cumulative** – cell wearing is a function of the entire history of shaped and non-shaped data that was programmed to it

- The second two advantages are **local** – observed only when currently programmed data is shaped - improve the average performance, power and reliability

4. **Negligible complexity & power, High throughput On-The-Fly operation**

5. **FW transparent – can be considered as a different type of scrambling**

# Leveraging Low Host Data Entropy - Ecosystem Considerations

- Data encryption results in high data entropy (randomizes the data)

- Data encryption at the host side should be avoided in order to take advantage of the low host data entropy via compression or shaping

- Encryption and Data Shaping can co-exist if they are performed at the memory controller level in the following order:
    - Compress
    - Encrypt
    - Expand via shaping

# Summary

- **Analyzing mobile traffic reveals low host data entropy**

  - ~0.5 average compression rate measured for sample usage on an Android based Smart Phone

- **Unrealized potential: the inherent "redundancy" in the host data can be leveraged for improving endurance, reliability, performance and power**

  - Apply methods of Deduplication, Compression and Shaping

- **Shaping provides a FW transparent low complexity & power approach for taking advantage of the low host data entropy**

  - Reduced cell wearing

  - Reduced error rates                    Improved endurance,
                                            performance and reliability
  - Increased error correction capability

- **Ecosystem cooperation is required in order to take advantage of the low host data entropy, under security and encryption requirements**

**SanDisk**

Thank you!

Questions?

**Contact:**
eran.sharon@sandisk.com or
stop by SanDisk booth # 204