

When it comes to Emerging MemoriesOne Size Doesn't Fit All!!

August 8, 2019



Overwhelming Deluge of Data

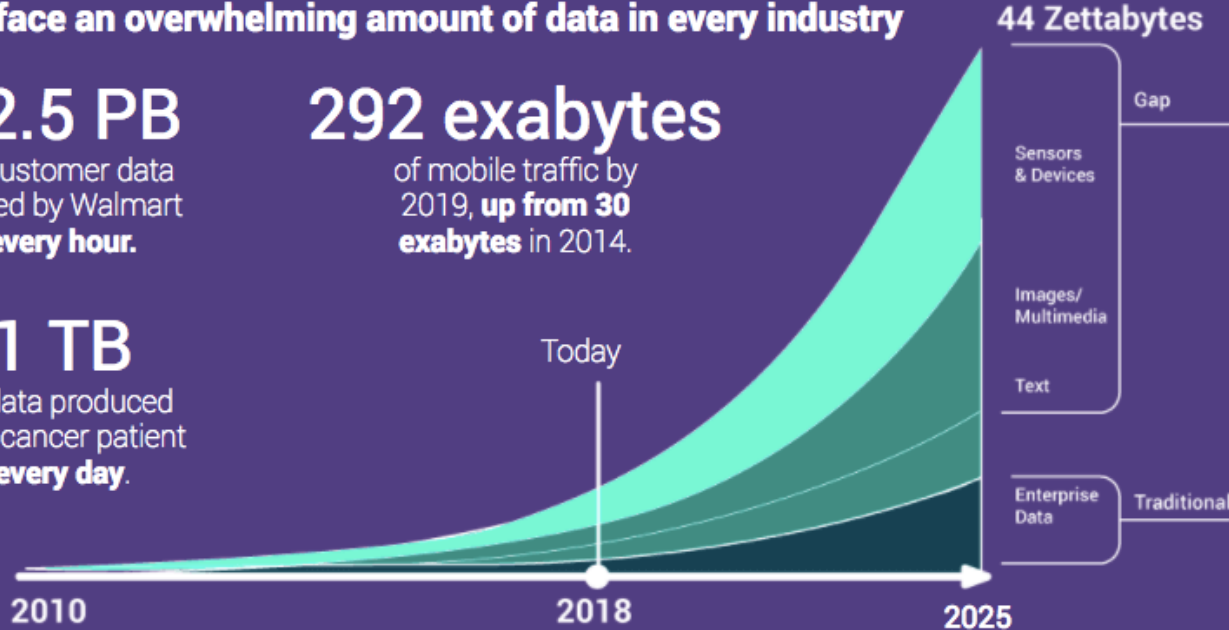
1 Zettabyte (ZB) = 1 Trillion Gigabytes (GB)

We face an overwhelming amount of data in every industry

>2.5 PB
of customer data
stored by Walmart
every hour.

292 exabytes
of mobile traffic by
2019, **up from 30**
exabytes in 2014.

1 TB
of data produced
by a cancer patient
every day.



Source © 2018 DVmobile Inc. All Rights Reserved.

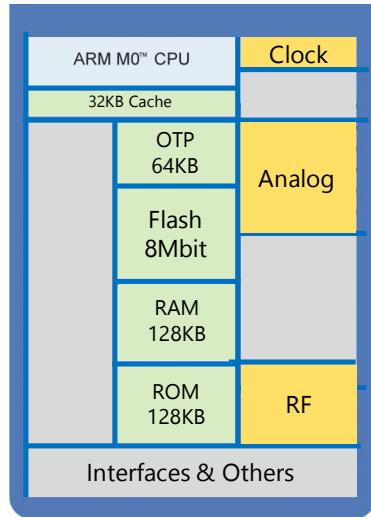
**...Driving the Need for more efficient Memory
and more distributed processing**

Edge Processing: Increased Edge Compute/Memory

- Due to the increased data collection, it is no longer practical and power efficient to move all the data to the Cloud
- Higher Processing at the Edge
- Higher Non-Volatile Memory to store
 - Program, Models/Coefficients, increased amount of Data collected
- Higher Volatile Memory for AI/Signal Processing
- Ultra-low Power to extend Battery Life and to reduce thermal issues to reduce cost and form factor

➔ More than ever, it is critical to have efficient memory in edge nodes

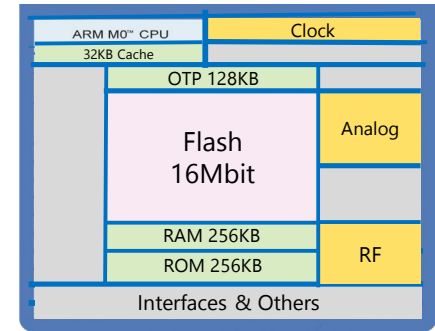
Resistive RAMs Provide a Simpler/Smaller Memory



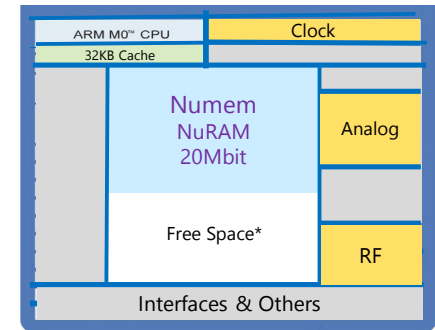
IoT SoC

Next Gen. Option-I --- > Current Memory Tech.

Next Gen. Option-II --- > Numem NuRAM



* Not to scale



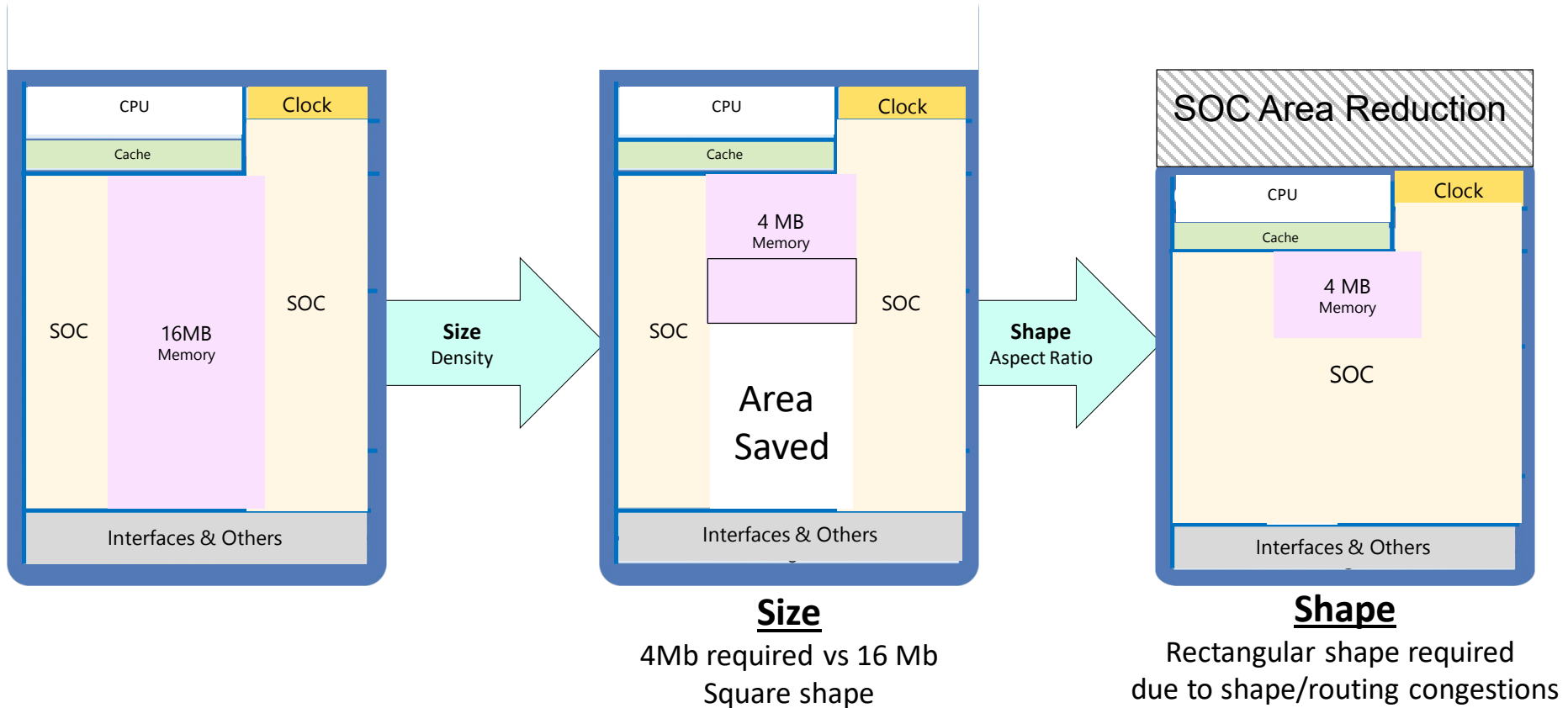
* Not to scale

Resistive RAMs such as MRAM, RRAM, PCRAM, enable Simplified & Smaller Unified Memory

Optimization: Significant Improvements in Speed/Power/Area

- Architecture
 - Memory Architecture: Read/Write Circuits, Decoders, Array
 - SOC Architecture and Distribution of Memory throughout the SOC
 - Memory Subsystem: BIST, Failure Correction (ECC, Repair, Calibration)
- Optimized Circuit Design affecting even more RRAM/PCRAM/MRAM (newer technologies)
- Design Customization including:
 - Size
 - Shape/Aspect Ratio (Column Muxing)
 - Array Subdivision/ Bit line length
 - Pipelining
 - Reflow/Non-Reflow Bitcells

Area: Exact Sizing and Shape Optimization



Not to Scale

Area: Exact Sizing and Shape Optimization

- Standard IP Cores are usually One-size-fits-all: e.g. 16Mb & 32Mb
- Example: you need a 4Mb & only 512K Reflow Capable (~20% larger)
- Tally up delta size $4 \times 1.2 - 0.5 \times (1 - 1.2)$ **4.7x larger**
- Shape/Routing Congestions also result in larger SOC: **~30% larger**
 - It is not always about the memory best area but about best SOC area
- Memory Design Architecture also affects size enormously: **>2x delta**

→ Tally up: **overall size delta could be in excess of ~13.4x**
in this example $(4.7/0.7 \times 2)$

Speed/Power Optimization

- The Memory design can significantly affect Speed/Power
 - Array Subdivision (Bitline length)
 - Pipelining
 - Optimized Circuit Design (Sense Amp, Decoder, Data Pipe, etc..)
- Shorter Bitline Integration time can improve performance by **>2x**
 - 5-30% larger area depending on optimization
- Design architecture is of paramount importance for best possible performance especially in newer technologies like RRAM/MRAM
 - Read Access times vary can widely from 2.5ns to 20+ns typ so by **8x**
- Pipelining achieves significantly higher bit rate (if acceptable), **2x** per stage
 - ➔ All and all optimization can reasonably **improve performance by >16x**

Optimization/Customization Improvements Summary

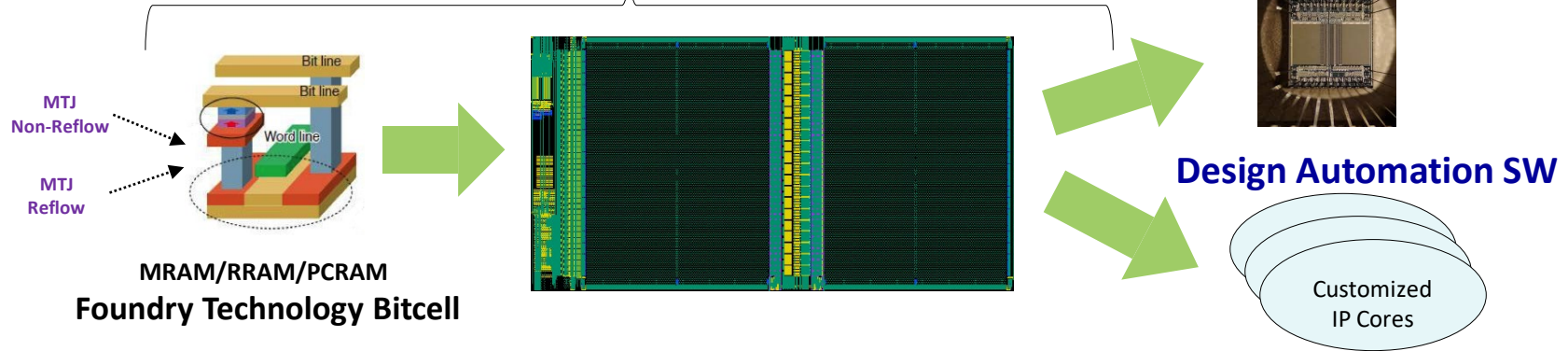
- > 13x Smaller Area

- > 16x better Speed/Power

- 50X Faster Time to Market through Design Automation SW

- More Reliable Product through Validation/Characterization SW

Patented Architecture and Design (Smaller size, higher Perf./Lower Power)



- **Optimized Standard & Custom Memory Subsystem IP Cores**
 - Validated on MRAM through tape-out of 19 devices to date and extensive testing
 - Our memory are non-volatile and 2-3x smaller than SRAM
- **Design Automation SW for Memory Customization & Validation/Characterization**
 - Yields much greater PPA improvements and higher reliability
- **Memory Chipllets** (die) for use in MCP SOCs

Thank You