# Computing Evolution: Technology and Markets

'Beyond the LHCb Phase-1 Upgrade' Workshop, Elba (IT), 28 – 31 May 2017

Helge Meinhard / CERN

Presenting material prepared by Bernd Panzer-Steindel / CERN

# Outline

- Semiconductor market
- Device market
- Processors
- Hard Disk
- Solid-State Disks
- Memory
- Tapes
- Server

- Summary
- References

# General Market (1)



**Chart 12**
**World Electronic Equipment Production By Type**
+3.2%
Annual Rate $B (converted @ fluctuating exchange rates)
preliminary 1Q'17 estimate

- industr+instru
- gov/mil
- Computer & Peripherals
- automotive
- consumer
- Telecom/Datacom
- Business

## Few companies dominating the markets

| Server CPUs | Intel (99%) |
|---|---|
| FPGA | Xilinx (49%), Intel (38%) |
| GPU | Intel (68%), Nvidia (18%), AMD (14%) |
| Hard disks | Western Digital (41%), Seagate (37%), Toshiba (22%) |
| Tape drives | IBM |
| Tape media | Fujifilm, Sony |
| NAND | Samsung (35%), Toshiba (20%), Western Digital, Micron |
| DRAM | Samsung (50%), Hynix (25%), Micron/Intel (19%) |

Electronic equipment production
is essentially flat (market saturation)

# General Market (2)

**Worldwide Semiconductor Revenues**
Year-to-Year Percent Change

Revenues have increased due to high prices for NAND and DRAM

February '17 = 16.5% Y/Y

Source: WSTS

Computing market only small part
~18 B ARM processors,~ 0.3 B Intel/AMD.

Expect 1 trillion semiconductor units shipped in 2018.

2016 Semiconductor Unit Shipments
(868.8B)

Microcompoent 3%
Memory 5%
Sensor/Actuator 2%
Std Logic 6%
Opto 25%
Analog 15%
Discretes 44%

Source: IC Insights

Thyristors, power transistors, diodes, etc.

**Tracking Semiconductor Unit Growth**

1,002.6
868.8
623.7
103.7
32.6

Source: IC Insights

# Device Markets (1)



**World Device Shipment Units by Segment**

Chart 25

Units (000)

|  | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 |
|---|---|---|---|---|---|---|---|---|
| Mobile Phone | 1,746,177 | 1,806,964 | 1,879,000 | 1,917,000 | 1,893,000 | 1,910,000 | 1,920,000 | 1,954,000 |
| Ultramobiles (Basic & Utility) | 120,203 | 209,788 | 226,000 | 196,000 | 169,000 | 161,000 | 158,000 | 157,000 |
| Ultramobiles (premium) | 9,787 | 21,517 | 37,000 | 44,000 | 50,000 | 60,000 | 72,000 | 82,000 |
| PC (Desk & Notebook) | 341,273 | 296,131 | 277,000 | 244,000 | 220,000 | 205,000 | 196,000 | 191,000 |

Growth rates: +5.2%  +3.6%  -0.7%  -2.9%  0%  +0.4%  +1.6%

Gartner 4/17

Market saturation: minimal or negative growth rates
Longer product lifetimes

# Device Markets (2)



Smartphone Unit Shipments to End Users — World

Units (Millions)

Gartner Dataquest 2/17 & prior reports, 1Q'17 estimate based on Trendforce Q1'17/Q1'16 growth rate



New mobile subscriptions Q1 2016

63 million new mobile subscriptions globally in Q1 2016

Top 5 countries by net additions Q1 2016

1 India +21 million
2 Myanmar +5 million
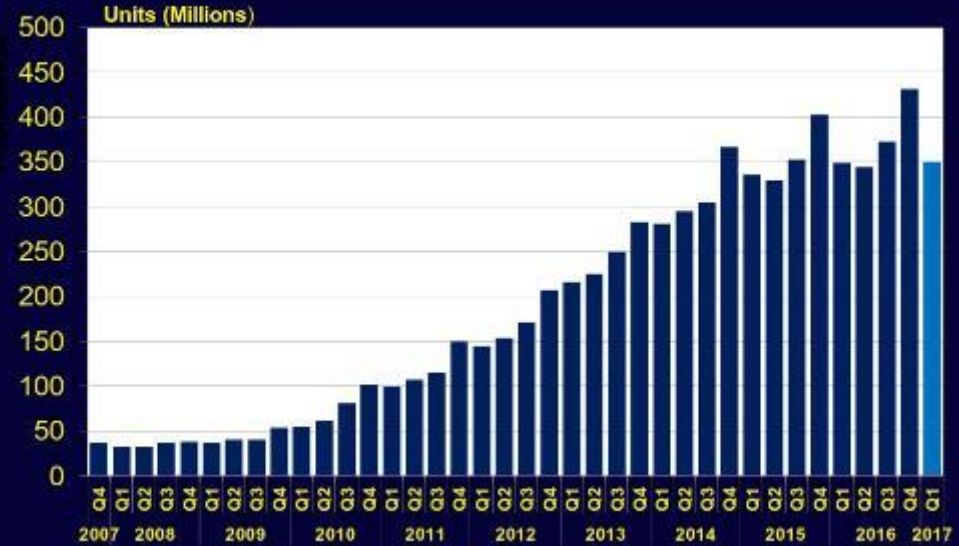3 Indonesia +5 million
4 USA +3 million
5 Pakistan +3 million

5 BILLION subscribers

The number of mobile subscriptions exceeds the population in many countries. This is largely due to inactive subscriptions, multiple device ownership or optimization of subscriptions for different types of calls. This means the number of subscribers is lower than the number of subscriptions. Today there are around 5 billion subscribers compared to 7.4 billion subscriptions.

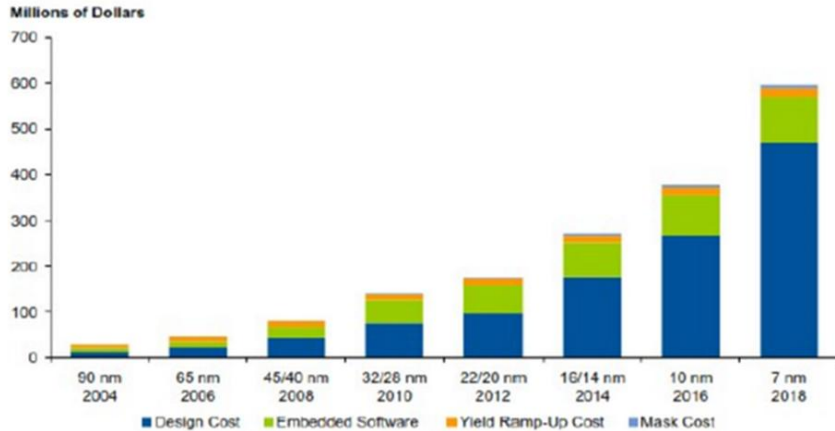Penetration (percent of population)

Saturation:
7.3 B phone subscriptions world-wide – more than the population

Replacement bump expected in 2018

# Processors (1)

**Estimated Cost of Developing Lower Node Chips**

Millions of Dollars



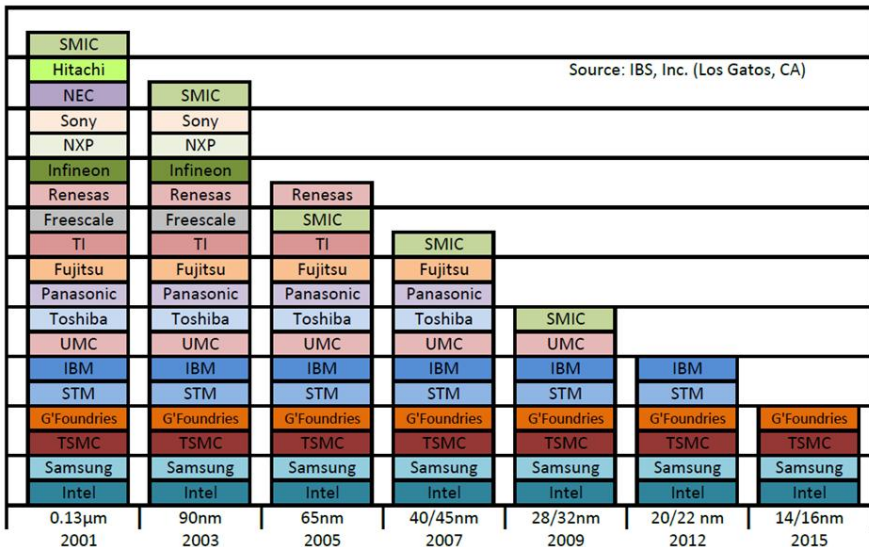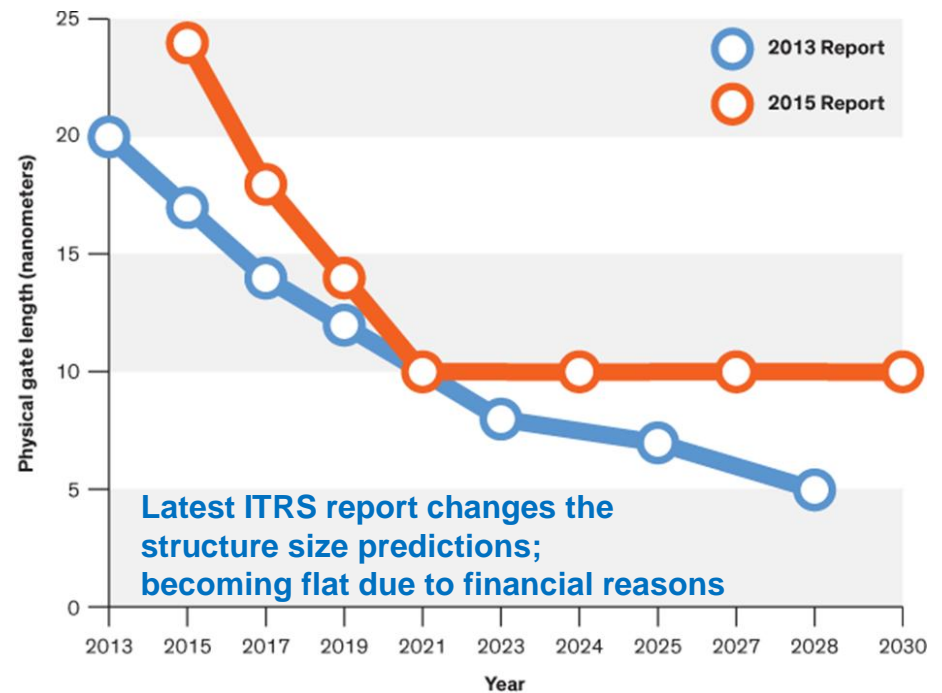Source: Gartner

Market Realist

| | SMIC | | | | | | |
| | Hitachi | | | | Source: IBS, Inc. (Los Gatos, CA) | | |
| | NEC | SMIC | | | | | |
| | Sony | Sony | | | | | |
| | NXP | NXP | | | | | |
| | Infineon | Infineon | | | | | |
| | Renesas | Renesas | Renesas | | | | |
| | Freescale | Freescale | SMIC | | | | |
| | TI | TI | TI | SMIC | | | |
| | Fujitsu | Fujitsu | Fujitsu | Fujitsu | | | |
| | Panasonic | Panasonic | Panasonic | Panasonic | | | |
| | Toshiba | Toshiba | Toshiba | Toshiba | SMIC | | |
| | UMC | UMC | UMC | UMC | UMC | | |
| | IBM | IBM | IBM | IBM | IBM | IBM | |
| | STM | STM | STM | STM | STM | STM | |
| | G'Foundries | G'Foundries | G'Foundries | G'Foundries | G'Foundries | G'Foundries | G'Foundries |
| | TSMC | TSMC | TSMC | TSMC | TSMC | TSMC | TSMC |
| | Samsung | Samsung | Samsung | Samsung | Samsung | Samsung | Samsung |
| | Intel | Intel | Intel | Intel | Intel | Intel | Intel |
| | 0.13µm | 90nm | 65nm | 40/45nm | 28/32nm | 20/22 nm | 14/16nm |
| | 2001 | 2003 | 2005 | 2007 | 2009 | 2012 | 2015 |

Figure 4. Dramatic Consolidation of state of the art CMOS Fabs. Source: IBS , Inc. (Los Gatos, CA).

## Non-linear costs for development

- Only four companies able to fabricate 14 nm chips
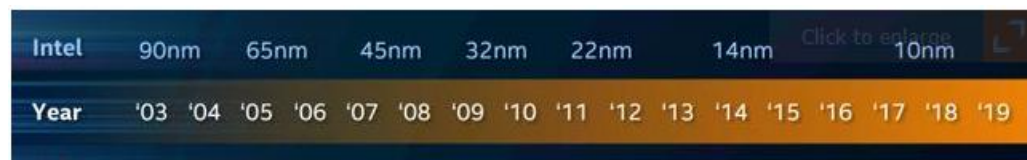- 10 nm Samsung fab costs $14 B



**Latest ITRS report changes the structure size predictions; becoming flat due to financial reasons**
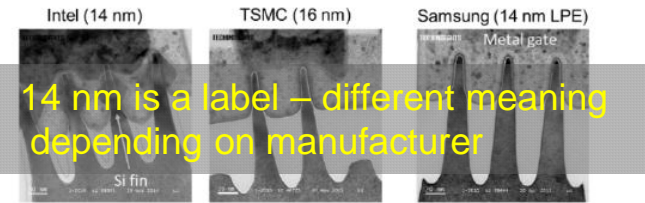
# Processors (2)



YESTERDAY PROCESS TECHNOLOGY — TICK (PROCESS) → TOCK (ARCHITECTURE)

TODAY PROCESS TECHNOLOGY — PROCESS → ARCHITECTURE → OPTIMIZATION

**Intel moved from 2-year cycle to 3 years or more**

| Intel | 90nm | 65nm | 45nm | 32nm | 22nm | 14nm | 10nm |
|-------|------|------|------|------|------|------|------|
| Year | '03 '04 | '05 '06 | '07 '08 | '09 '10 | '11 '12 '13 | '14 '15 '16 '17 | '18 '19 |

**16/14 nm finFET Comparison**

TECHINSIGHTS | Proving Patent Val

Intel (14 nm) — TSMC (16 nm) — Samsung (14 nm LPE) — Metal gate — Si fin

14 nm is a label – different meaning depending on manufacturer

| Feature | Intel | TSMC | Samsung |
|---------|-------|------|---------|
| Gate length (nm) | 24 | 33 | 30 |
| Min contacted gate pitch (nm) | 70 | 90 | 78 |
| Fin height under gate (nm) | 42 | 37 | 37 |
| Fin pitch (nm) | 43 | 45 | 49 |
| Min metal pitch (nm) | 52 | 70 | 67 |

• Intel transistors are smaller than TSMC or Samsung

The ConFab.

#TheConFab2016

## Incubation Time

- **Strained Silicon**
  - 1992->**2003**
- **HKMG**
  - 1996->2007
- **Raised S/D**
  - 1993->2009
- **MultiGates**
  - 1997->2011

Source — Metal — Gate Insulator — Drain — 1998

~ 12-15 years

Decrease of feature size goes along with new material technologies

R&D à production needs 12-15 years



Currently in production — FinFET — Fully depleted silicon-on-insulator — Lateral nanowire 2019 — Vertical nanowire 2021 — Monolithic 3D 2024

Insulator | Source or drain | Gate | Silicon

**7nm structures need new technologies: nanowires and non-silicon material**

# Accelerators: GPU (1)

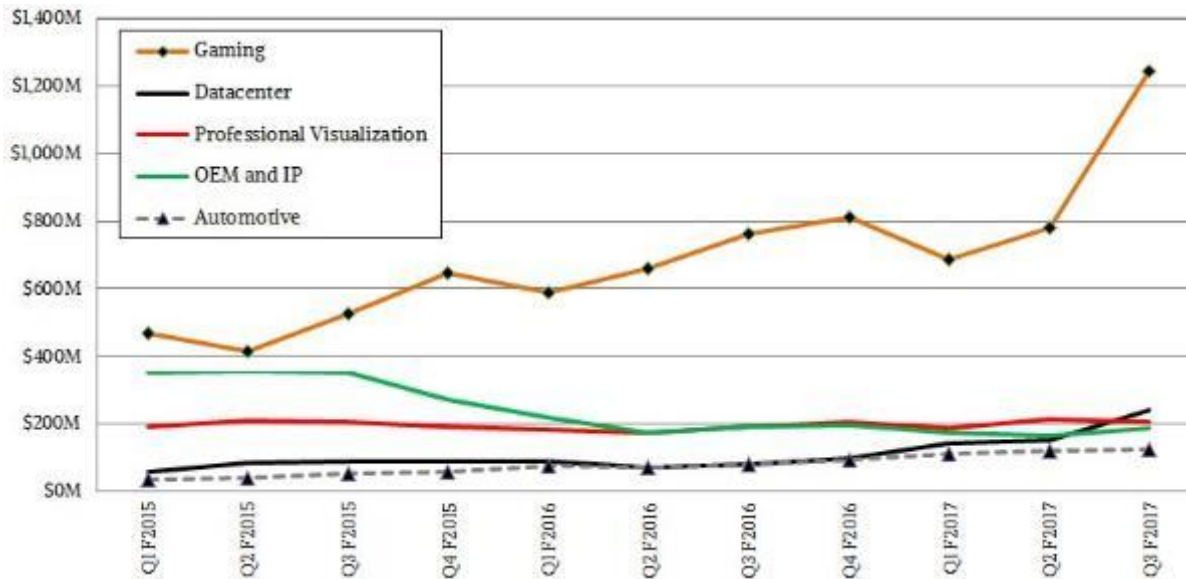Embedded market shares (CPU+GPU):
Intel 68%, Nvidia 18%, AMD 14%
Discrete GPU cards: Nvidia 77%, AMD 23%

New products announced:
- Nvidia Volta: 12 nm, 21 B transistors, 15 TFlops SP (Q1 2018)
- AMD Vega:  14 nm,  12 B transistors, 12.5 Tflops SP (Q3 2017)

Focus: high-end Gamer (DP and FP16 artificially reduced)

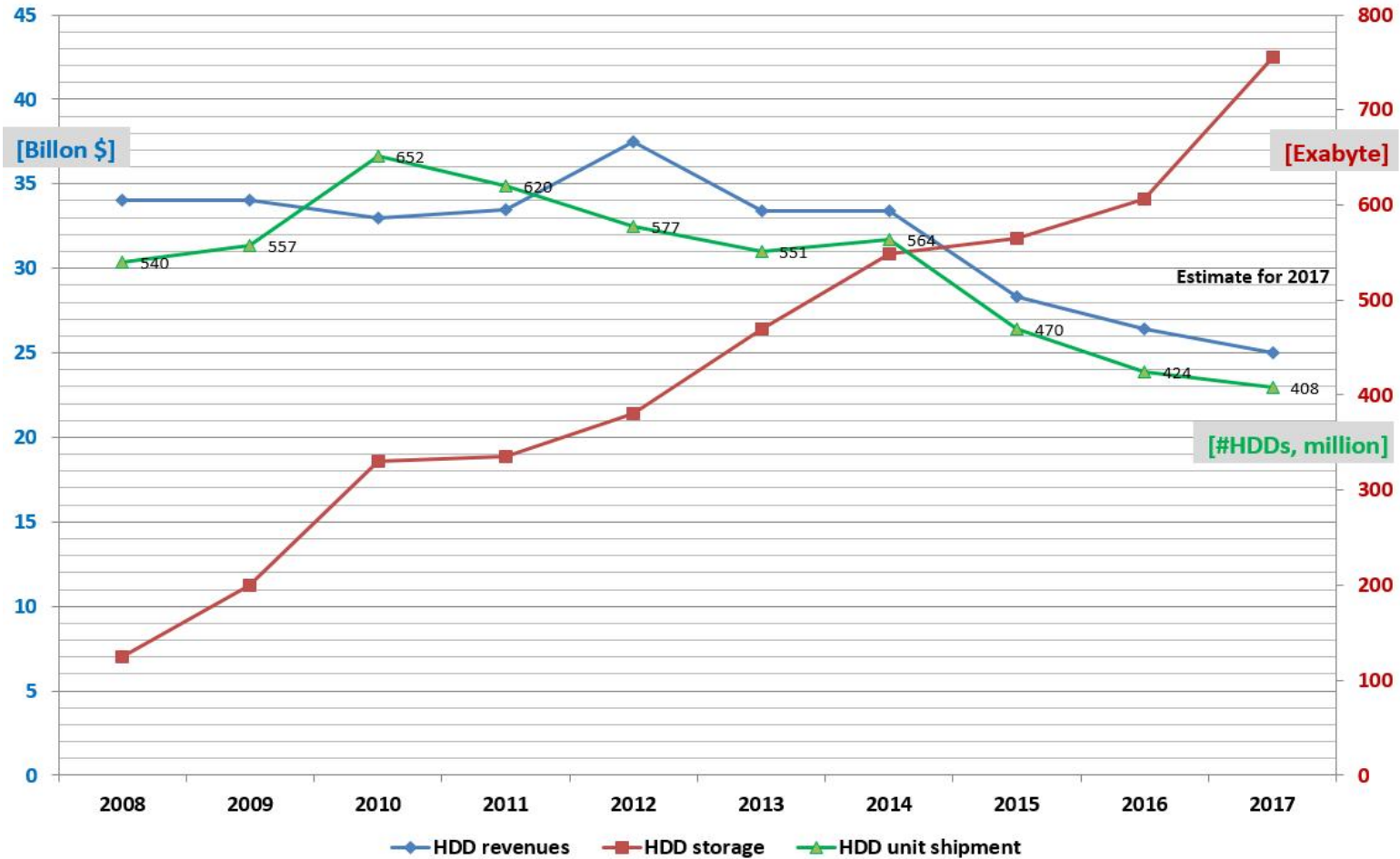Professional workstation cards and HPC: small niche,  ~2 million cards per year (compared to 350 million total GPUs)



Nvidia Revenues

# Accelerators: GPU (2)

- New focus for graphic cards: machine learning

- Move to FP16 and even INT8 architectures, less precision à 8 bit processing !

- Google TPU Tensor Processing Unit

- New start-ups with special processor designs:
  e.g. KnuEdge, Nervana (just bought by Intel), SpiNNaker, Eyeriss, P-Neuro, NeuRAM3
  - Essentially not usable as general purpose processors
  - Online?!

- Intel changing strategies also for their KnightsXX processors, 'forking' models (increase FP16 and decrease DP)
  - ~100k units per year, very small market

- Qualcomm plans to add neuromorphic chips into the smartphone

# Hard Disks (1)



Continuous decrease in revenues
Shipments decreasing, Seagate just closed one of their major production fabs
Pressure from SSDs in the high end enterprise drive market

# Hard Disks (2)

Combining bit density (30% annual growth rate) and volume density (number of platters, helium) → 100 TB in 2025 conceivable

Areal density improvement dropped from ~40% to 16% per year



## ASTC Technology Roadmap



## Areal Density Trends

Chart provided courtesy of the Information Storage Industry Consortium (INSIC)



PMR  limit at 1 TbPSI
SMR adds ~25%,  market small
HAMR should provide 5 TbPSI

HAMR delayed, production in 2018

# Hard Disks (3)

Focus on infrastructure cost reduction
i.e. Helium drives,  more platters per drive

**TrendFocus Nearline HDD Forecast by Disk-per-Drive Ratio**



Source: TrendFocus; Stifel

To reduce costs, at CERN, we are looking into
- Desktop drives
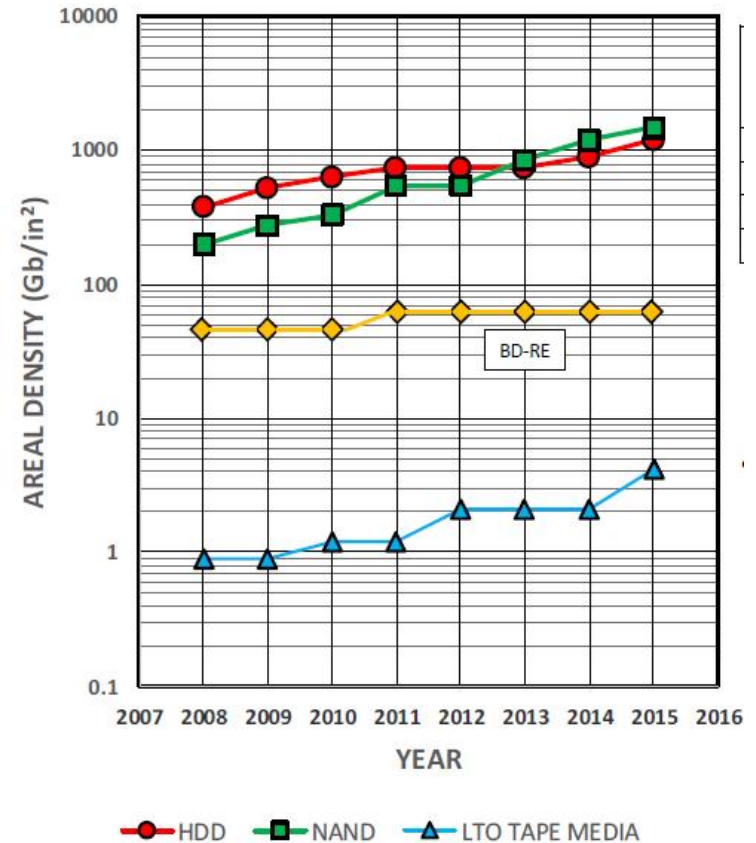- Multi PB disk server
- Erasure encoding

Google paper:   interest in much bigger drives (> 3.5")
Amazon patent:  separate mechanics from electronics
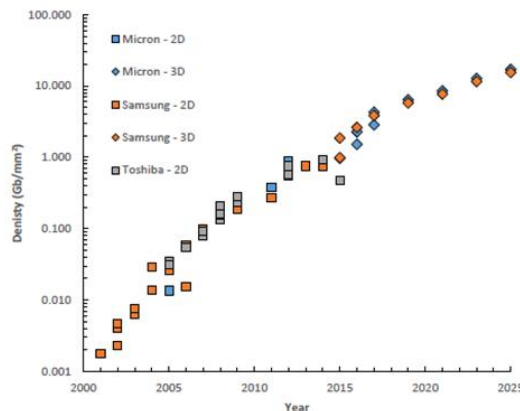                (steering boards outside of the drive)

# Solid-State Disks (1)

NAND density has surpassed HDD density

## 3D NAND – scaling in the third dimension

- 2D NAND scaling beyond 16nm/15nm is uneconomical.

- 3D NAND adds additional layers for scaling in place of 2D lithographic scaling.

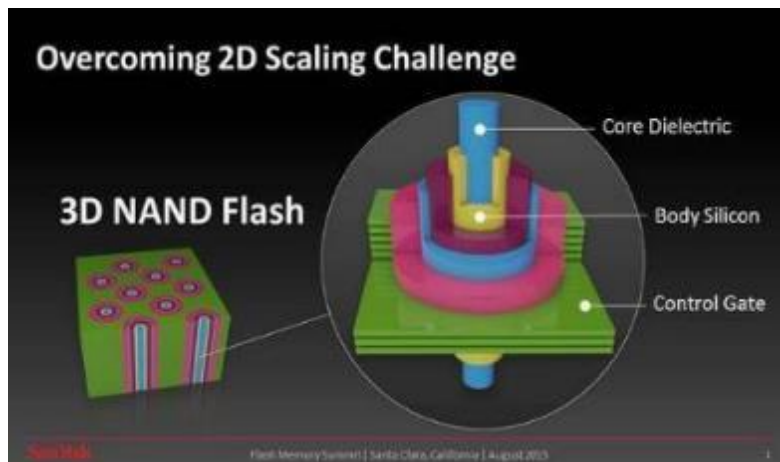- Bit density is continuing to scale with the potential for terabit NAND die.
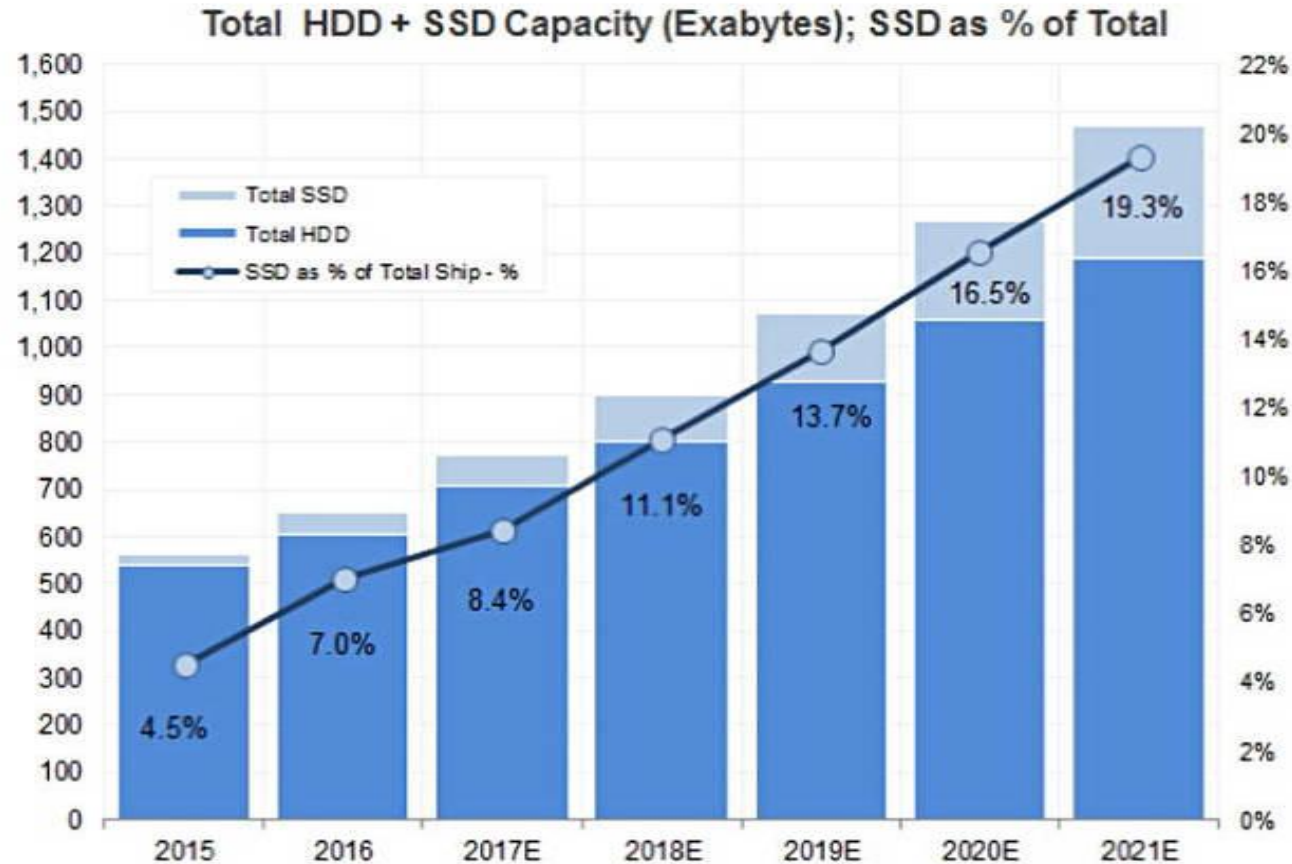
# Solid-State Disks (2)



| | 2016 | | | | 2017 | | | | 2018 | | | | 2019 | | | | 2020 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1Q | 2Q | 3Q | 4Q | 1Q | 2Q | 3Q | 4Q | 1Q | 2Q | 3Q | 4Q | 1Q | 2Q | 3Q | 4Q | 1Q | 2Q | 3Q | 4Q |
| SEC | 48L | | | 64L | | | | 96L | | | | | | | >=128L | | | | | |
| SK Hynix | 48L | | 64L | | | 72L | | | | 96L | | | | | | >=128L | | | | |
| Toshiba | 48L | | | | 64L | | | | 96L | | | | | | >=128L | | | | | |
| Micron | 32L | | 64L | | | | 96L | | | | | | >=128L | | | | | | | |
| YMTC | | | | | | | | 32L | | | 64L | | | >=96/128L | | | | | | |

Source: Yangtze Memory, Samsung Securities

NAND:

- 2D scaling came to an end 2 years ago
- Feature size manufacturing lags slightly behind processor structure sizes, <20nm today
- 3D: 64-layer 3D NAND in production; 72-layer expected end of 2017
- >30% price increases since last year, expect to last until end 2017
- Yield of 3D NAND improving, expect >50% of all shipments in Q4 2017 to be 3D instead of 2D

# Solid-State Disks vs. Hard Disks (1)



Total HDD + SSD Capacity (Exabytes); SSD as % of Total
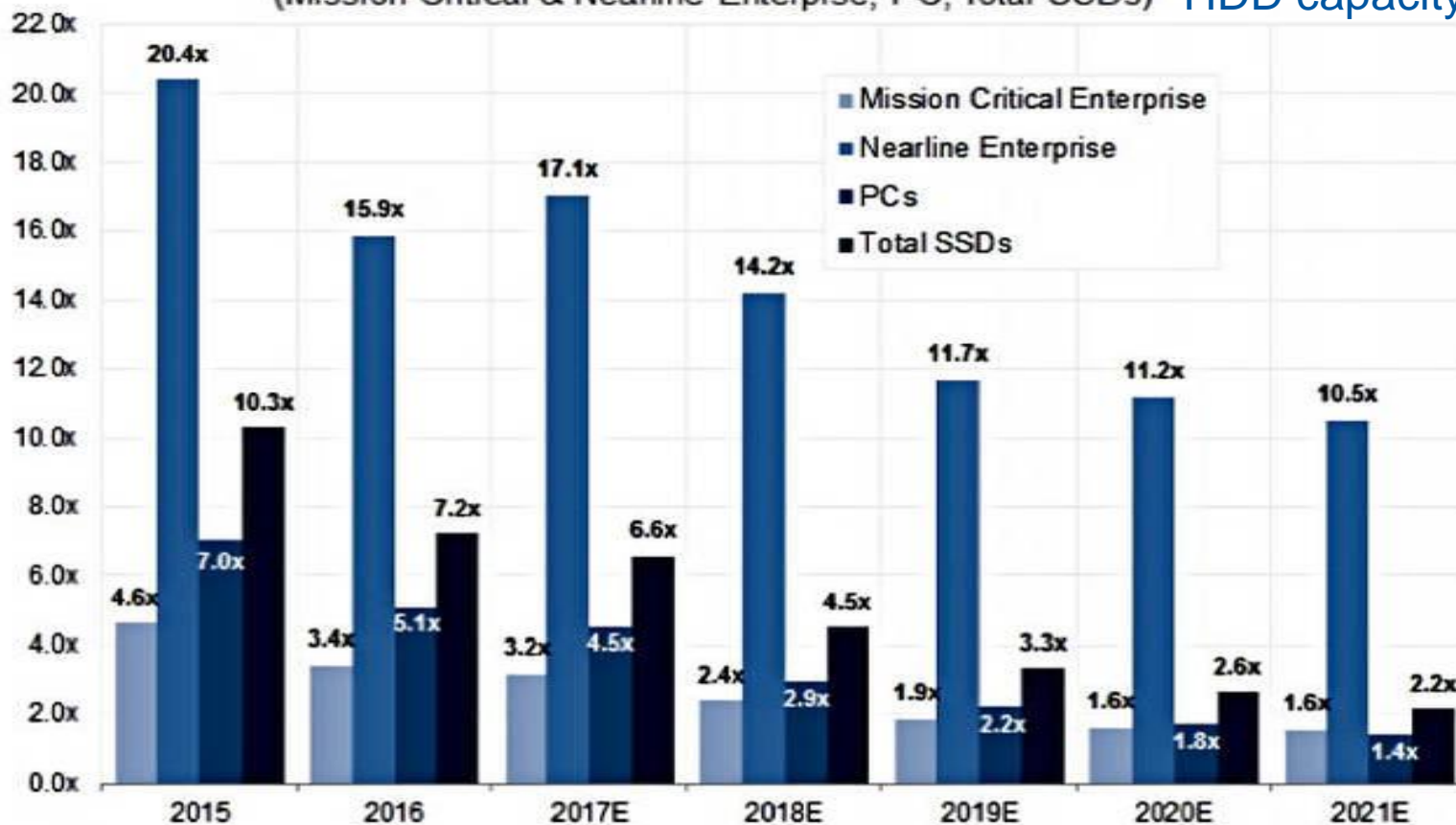
Source: IDC; Stifel

**16 times more HDD capacity than SSD was shipped in 2016**

**NAND Fab investment of 100-200 B$ necessary to achieve HDD ExaByte deliveries**

# Solid-State Disks vs. Hard Disks (2)

**SSDs not cost effective for HDD capacity replacement**



**Total SSD $/TB Premium vs. HDDs**
(Mission-Critical & Nearline Enterpise, PC, Total SSDs)

Legend:
- Mission Critical Enterprise
- Nearline Enterprise
- PCs
- Total SSDs

Source: IDC; Stifel

**Currently supply shortage of NAND
à  SSD price increases**

**Nearline = capacity drives
à  SSDs not foreseen for large scale storage**

# Memory: DRAM

Limited future improvements on performance and energy efficiency



Memory technology trend
- GDDR6 with over 14Gbps, beyond 10Gbps GDDR5
- LP5, 20% more power-efficient than LP4X
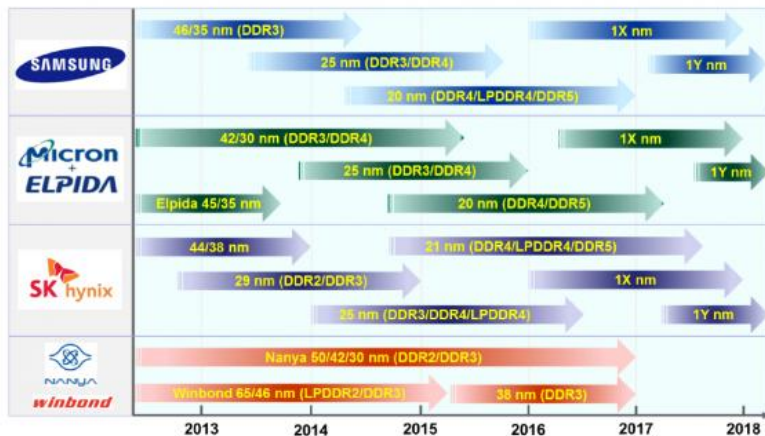


**Branded DRAM Revenue**

DRAMeXchange 2/2017 & earlier reports



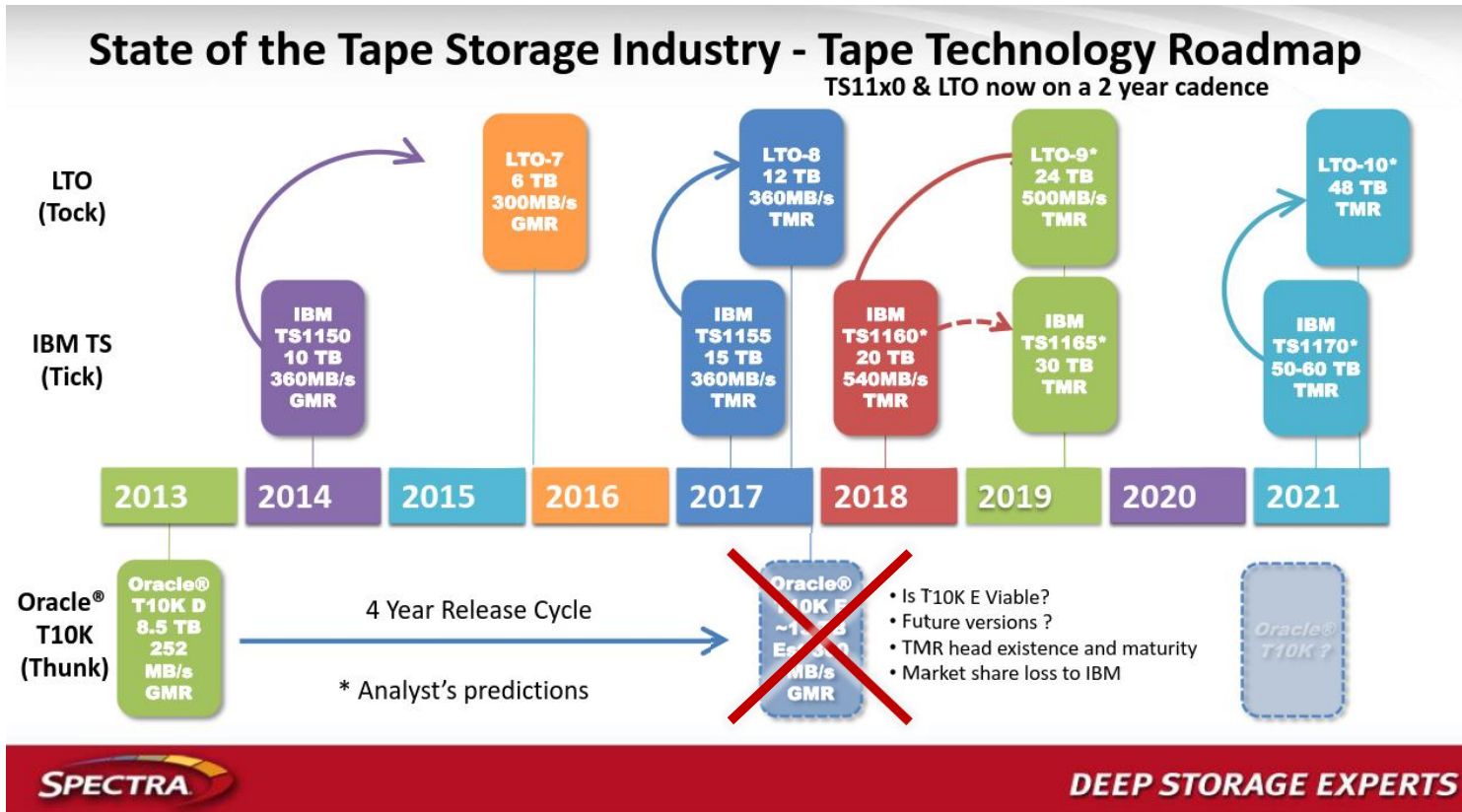DRAM Technology Review — TECHINSIGHTS

■ DRAM Process Node Roadmap (Manufacturers)

Volatile market, supply shortage started second half 2016, > 50% price increase until now, will be ongoing until end 2017

2 Chinese companies will enter the DRAM market in 2017

# New Memory Technologies

- 3D Xpoint: new technology from Intel and Micron, presumably a variant of Phase Change Memory
  Specs are changing:
  Announcement 2015:  1000x  faster, 1000x  endurance,  10x  denser than NAND
  IDF 2016:                         10x  faster,       3x endurance,    4x  denser than NAND
  first products (Optane) announced, usage currently limited to high end HDD caching

- Memristors: developed since 2008; HPE now collaborating with SanDisk (ReRAM)

- Spin torque MRAM in larger production units available (Everquest + Globalfoundries)
  Low density and high price

- Nantero just received extra funding for their Carbon Nanotube NVRAM, exists since 17 years, no product yet

- RRAM or ReRAM , various new categories being developed: Oxide RAM (OxRAM), Conductive-Bridge RAM (CBRAM) or Self-Rectifying Cells (SRC)

à    But… NAND fab investments are high, extended technology lifetime with 3D, hard to replace in the short term

# Magnetic Tapes (1)



## State of the Tape Storage Industry - Tape Technology Roadmap
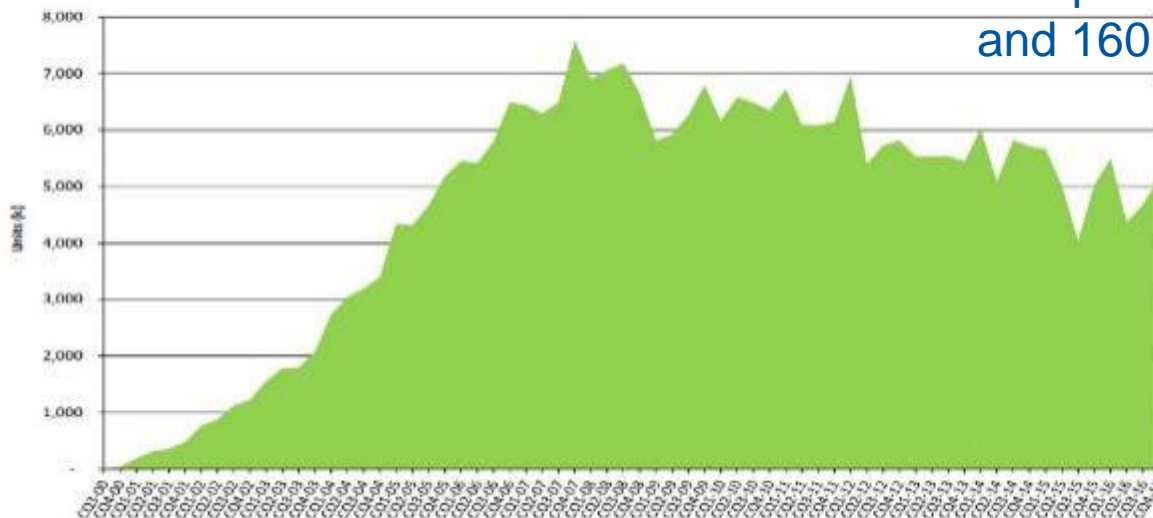TS11x0 & LTO now on a 2 year cadence

- Enterprise drives:
  Oracle 2017: stopped Enterprise production !!
  IBM end 2017: 10 TB à 15 TB  (just announced, TS1150)
  LTO-8 end 2017:  12 TB
- Technology in the lab:
  Fujifilm 154 TB, Sony 185 TB, IBM 220 TB

Technology change to
Tunnel Magnetoresistive heads
(used already in HDDs)

# Magnetic Tapes (2)

**Unit Shipments: Calendar Quarter**

~10 EB LTO capacity shipped per quarter
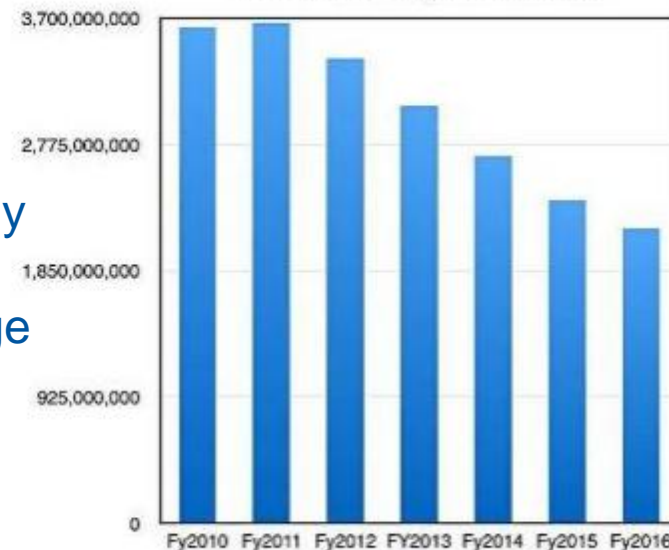à Comparted to 12 EB SSDs, 36 EB NAND
and 160 EB HDDs

Steady decrease of LTO tapes
shipped and revenues

Future of tapes ?

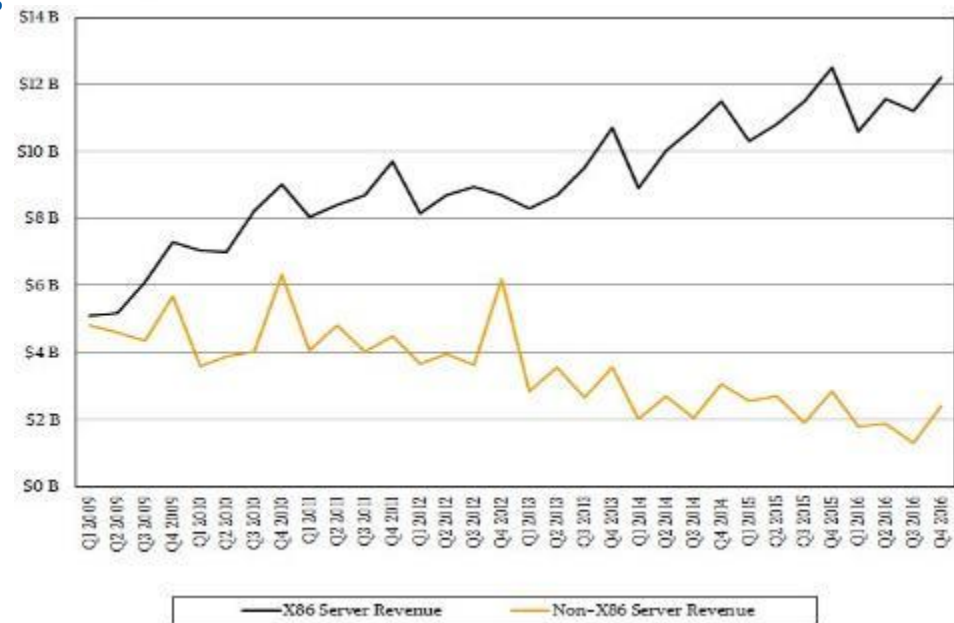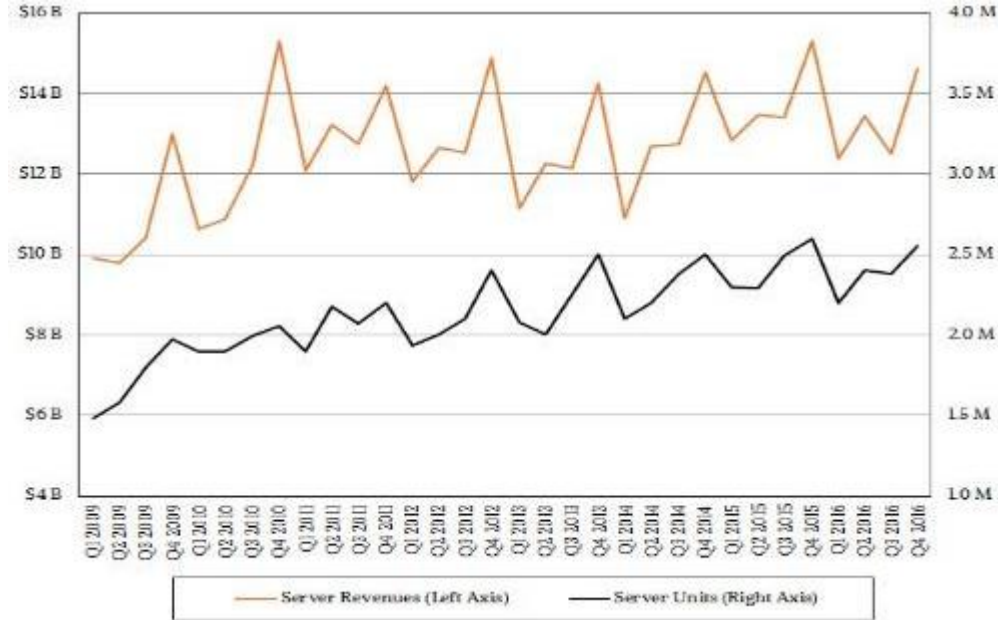All tape drive heads are now produced by IBM only

à Steady decrease of revenues in the IBM storage
group

**IBM Annual Storage HW Revenues**

# Servers (1)

- Server market is saturated: flat revenues and unit shipments during the last 2 years
- High profit market

- Single vendor: Intel, 99% market share

- Several initiatives to change that:

- OpenPower (IBM): consortium with many members
  - But revenues still going down, little impact so far
  - Announcement of POWER9 might help

- ARM server:
  - AppliedMicro , Qualcomm, Cavium:  new high end products  Announcements for 2H2017  (third ARMv8 Wave 2017-2018),   First two waves had little impact
  - Phytium (China), "Mars" processor

- AMD with new processor design (Zen), Naples for servers ~Q3 2017

- Fujitsu ARM-powered supercomputer
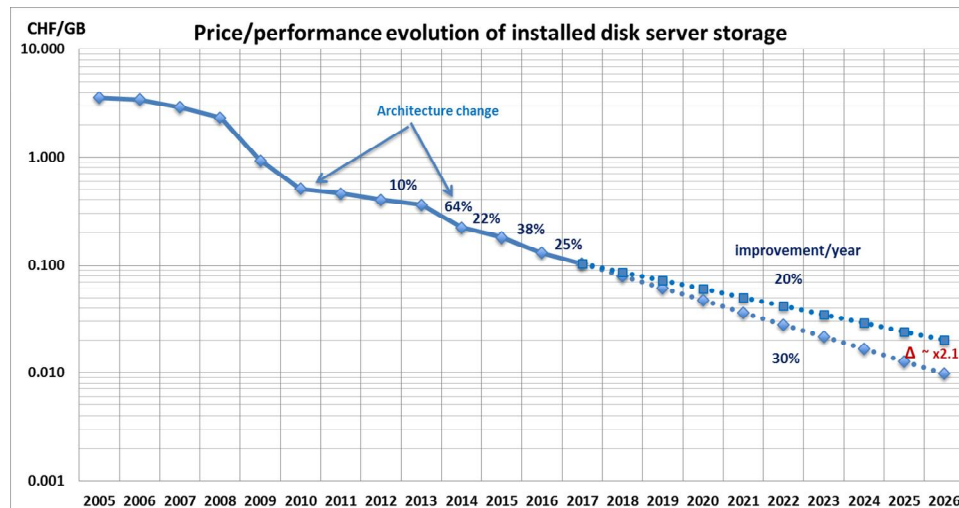  - Add large vector instructions to the ARM design
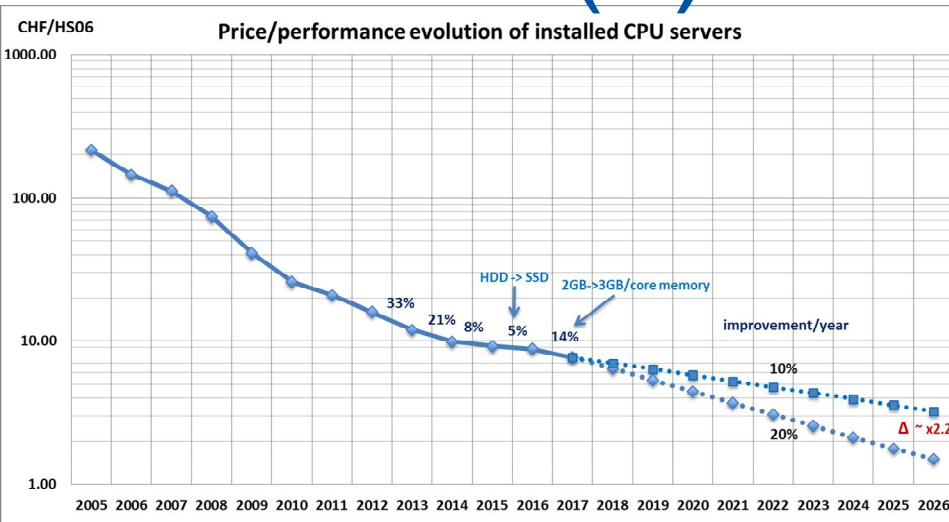  - Aimed for 2020, now ~2022



Server Revenues (Left Axis) — Server Units (Right Axis)



X86 Server Revenue — Non-X86 Server Revenue

# Servers (2)



Price/performance evolution of installed CPU servers

**Preliminary extrapolation of CPU and disk server costs (based on CERN procurements)**

Pessimistic and reasonable improvement extrapolations

Influence of changing software and hardware architecture requirements to be taken into account (programs, data model, data centre, …)

e.g. CERN moves from 2 to 3 GB/core (+8% cost), driven by experiment usage AND technology boundary conditions

- **Moore's Law and Kryder's Law are slowing down**
  - **18 months à  >= 3 years**

- **Real cost/performance evolution driven by financial and market aspects rather than technology**



Price/performance evolution of installed disk server storage

# Summary (1)

- Device markets (smartphones, tablets, PCs, notebooks, servers, HPC) saturated or even negative growth
  - Replacement market

- Moore's Law in trouble, financial issues
  - Not clear how this effects price/performance evolution
  - So far okay for CPU and disk servers

- Technology improvements still continuing, but requires high CAPEX
  End-product price tag evolution more complicated

- Market dominance of few companies increases, competition diminishing

# Summary (2)

- Technology alone unlikely to solve the computing problem at HL-LHC and beyond
  - Not much more to be expected than minor contributions

# References

More details and a (long) list references:

https://twiki.cern.ch/twiki/bin/view/Main/TechMarketPerf