# Avoiding the Dark Ages with Memristors

## Shahar Kvatinsky
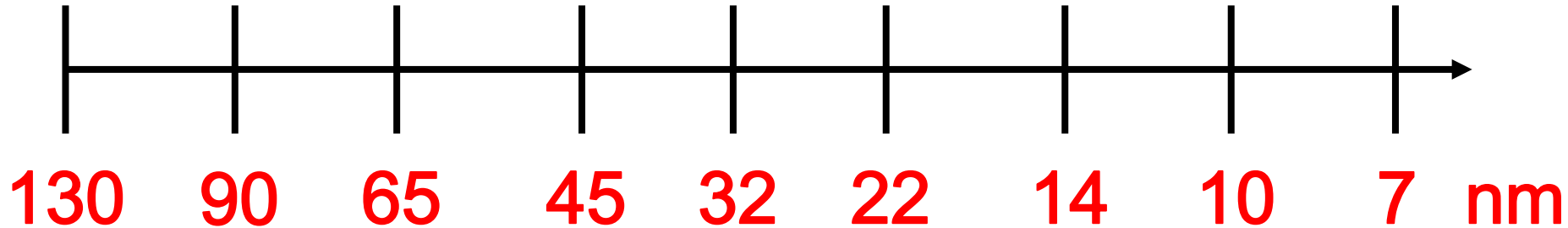
**Department of Electrical Engineering
Technion – Israel Institute of Technology**
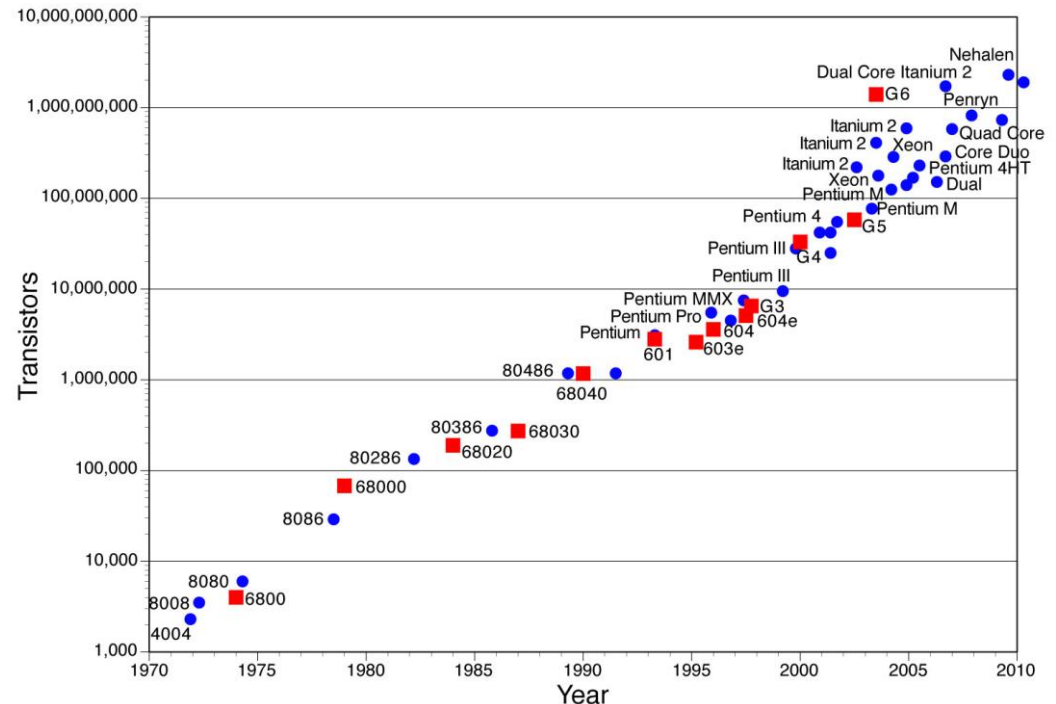March 2016

# Scaling 101 – Moore's Law

| 2001 | 2004 | 2006 | 2008 | 2010 | 2012 | 2014 | 2016 | 2018 |
|------|------|------|------|------|------|------|------|------|

130   90   65   45   32   22   14   10   7   nm

$$S = \frac{45}{32} = {\sim}1.4X$$

$$S^2 = {\sim}2$$

# Scaling 101 – Dennard Scaling

------------------------------------------- $S^3$

------------------------------------------- $S^2$

------------------------------------------- $S$

------------------------------------------- 1

# Scaling 101 – Dennard Scaling

------------------------------------------------------- $S^3$

------------------------------------------------------- $S^2$

$S^2 = 2X$
More transistors

------------------------------------------------------- $S$

------------------------------------------------------- $1$

# Scaling 101 – Dennard Scaling

S = 1.4X
Faster transistors
(Frequency scaling)

$S^2$ = 2X
More transistors

$S^3$

$S^2$

$S$

$1$

# Scaling 101 – Dennard Scaling

$S = 1.4X$
Faster transistors
(Frequency scaling)

$S^2 = 2X$
More transistors

Computing capabilities increased by $S^3 = 2.8X$

$S^3$

$S^2$

$S$

$1$

# Scaling 101 – Dennard Scaling

$S = 1.4X$
Faster transistors
(Frequency scaling)

$S^2 = 2X$
More transistors

Computing capabilities increased by $S^3=2.8X$

2.8X more transistors switches per second
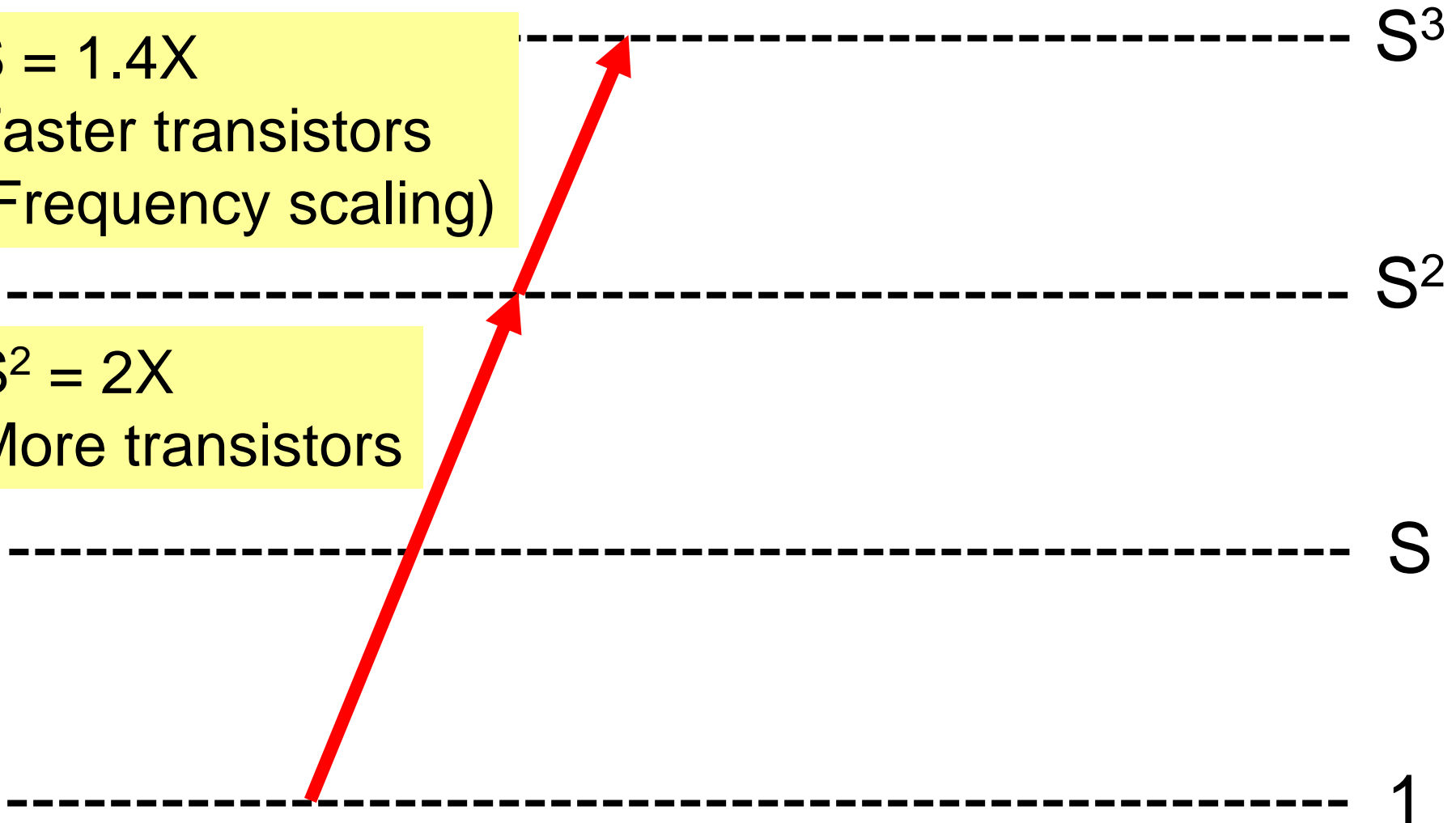Power increased by 2.8X

$S^3$

$S^2$

$S$

$1$

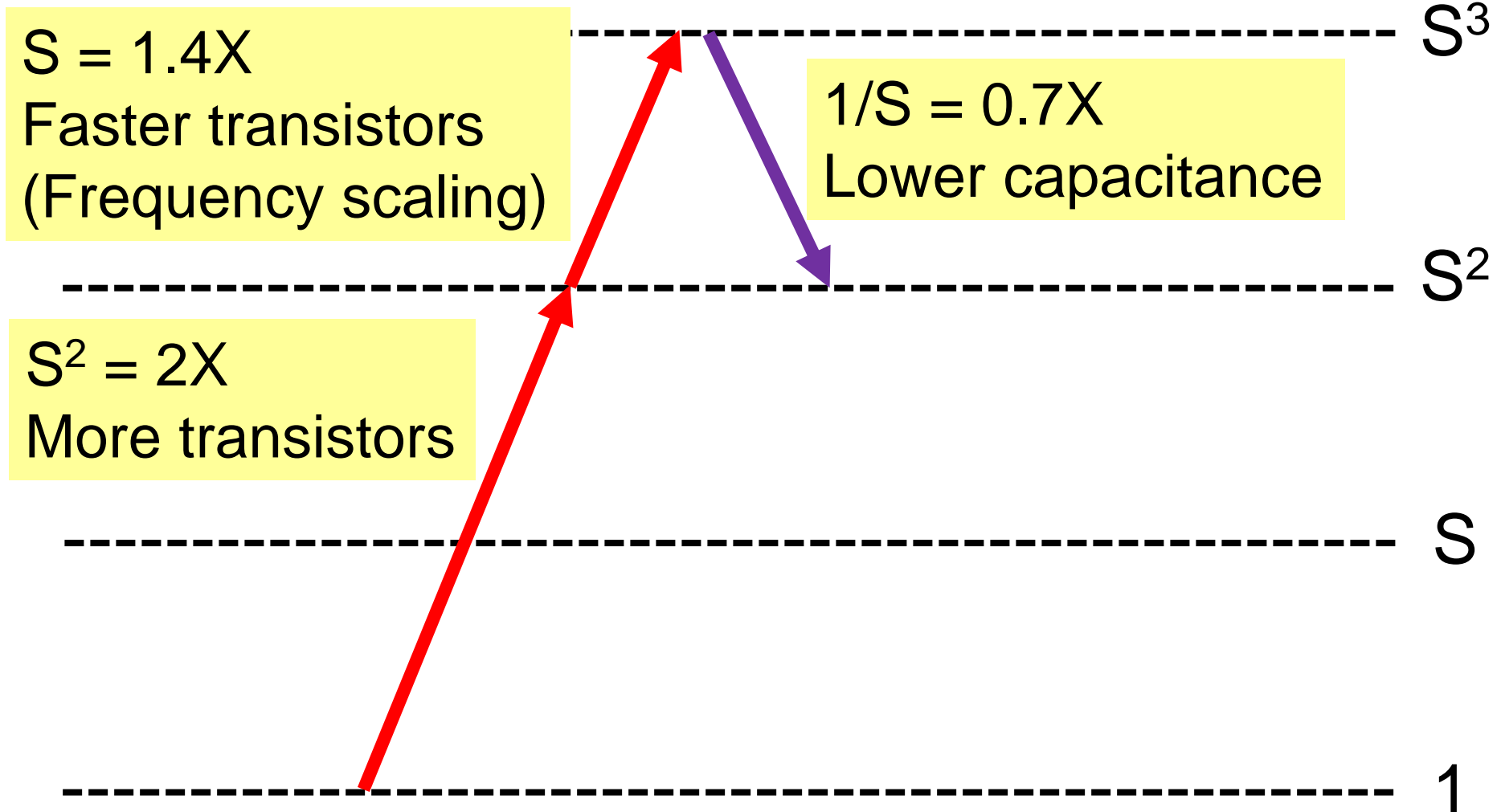# Scaling 101 – Dennard Scaling

S = 1.4X
Faster transistors
(Frequency scaling)
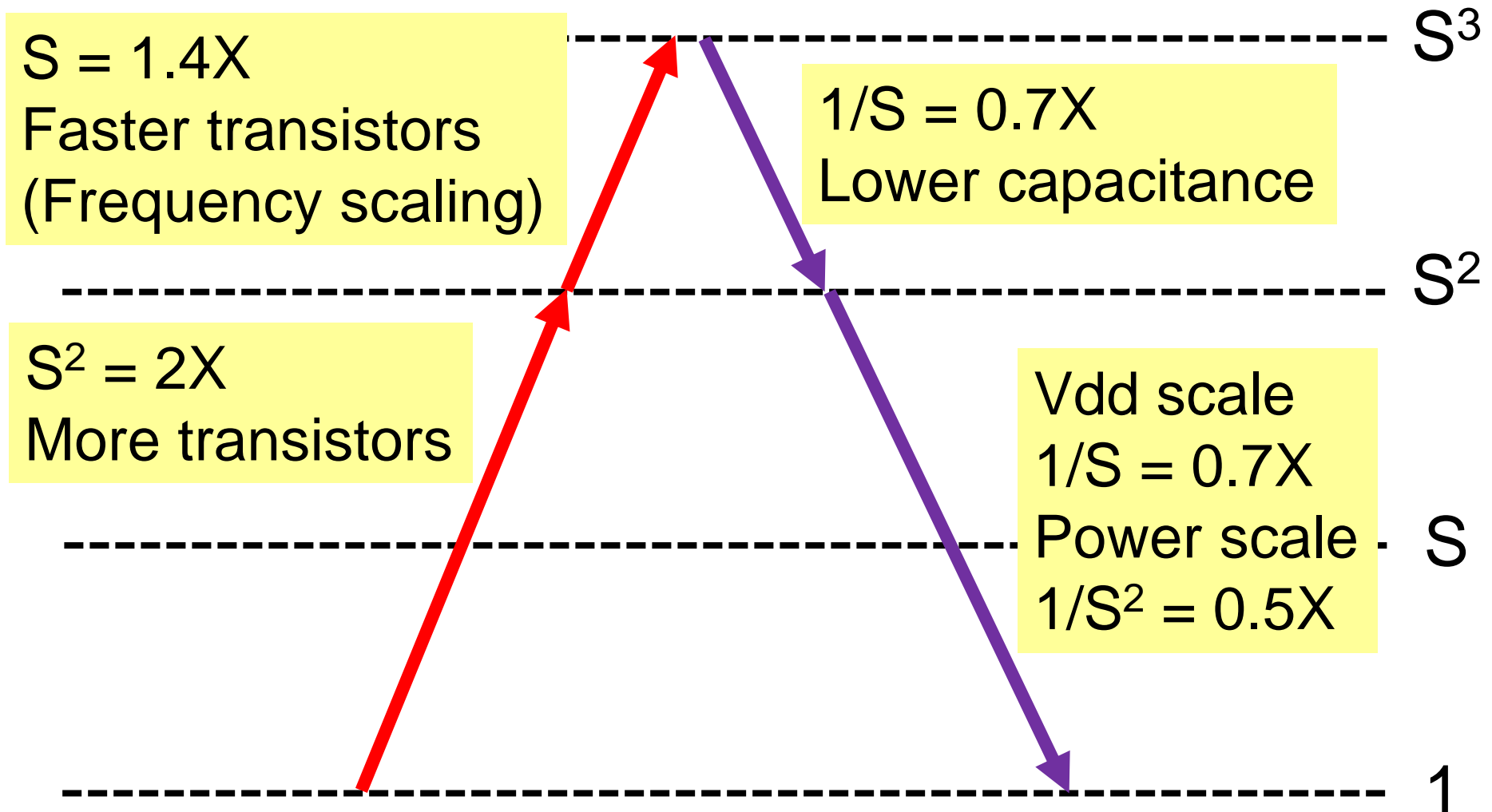
$S^2$ = 2X
More transistors

$S^3$

$S^2$

$S$

$1$

# Scaling 101 – Dennard Scaling

$S = 1.4X$
Faster transistors
(Frequency scaling)

$1/S = 0.7X$
Lower capacitance

$S^2 = 2X$
More transistors

$S^3$

$S^2$

$S$

$1$

# Scaling 101 – Dennard Scaling

$S^3$

$S^2$

$S$

1

S = 1.4X
Faster transistors
(Frequency scaling)

$S^2$ = 2X
More transistors

1/S = 0.7X
Lower capacitance

Vdd scale
1/S = 0.7X
Power scale
$1/S^2$ = 0.5X

# 2005 The End of Dennard Scaling

## Threshold Scaling and Leakage

$S = 1.4X$
Faster transistors
(Frequency scaling)

$1/S = 0.7X$
Lower capacitance

$S^2 = 2X$
More transistors

$S^3$

$S^2$

...cale
$1/S = ...7X$
...owe...cale
...5X

S

1

# The End of Frequency Scaling

S = 1.4X
Faster transistor
(Frequency scaling)

$S^2$ = 2X
More transistors

1/S = 0.7X
Lower capacitance

$S^3$

$S^2$

S

1

# Moving to Multicore

$S^3$

$S^2$

$S^2 = 2X$
More transistors

S

1

# Dark Silicon



$S^3$

$S^2$

$S^2 = 2X$
More transistors

S

**Power limitation**

1

14

# Dark Silicon



$S^3$

$S^2$

$S^2$ = 2X
More transistors

$S$

**Power limitation**

$1$

# The Four Horsemen of Dark Silicon

**Taylor DAC 2012**

- Shrink

- Dim

- Specialize

- Technology magic

  (*Deus Ex Machina*)

# The Four Horsemen of Dark Silicon
## Taylor DAC 2012

- Shrink

- Dim

- **Specialize**

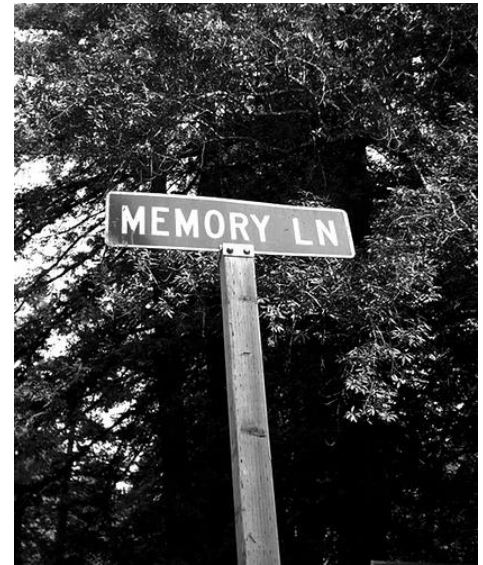- Technology magic

  (*Deus Ex Machina*)

# Sources of Energy Inefficiency

| Operation (16-bit operand) | Energy/Op (45 nm) | Cost (vs. Add) |
|---|---|---|
| Add operation | 0.18 pJ | 1X |
| Load from on-chip SRAM | 11 pJ | 61X |
| **Send to off-chip DRAM** | **640 pJ** | **3,556X** |

I-Cache Access   Register File Access   Control   Add

M. Horowitz, "Computing's Energy Problem (and what we can do about it)," ISSCC Keynote 2014
A. Pedram, S. Galal, S. Richardson, S. Kvatinsky, and M. Horowitz, "Dark Memory and Accelerator-Rich System Optimization in the Dark Silicon Era," IEEE Design & Test (submitted)
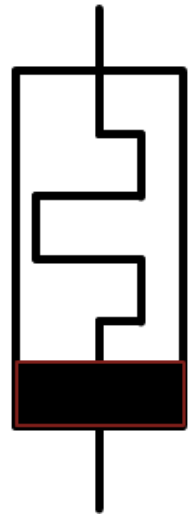
# Dark Memory and Specialization

- Memory system contributes >50% system power

- Memory hierarchy does not solve everything, SRAM is never completely dark

- Specialization increases memory power portion



- Amdahl's law - need to dim memory

A. Pedram, S. Galal, S. Richardson, S. Kvatinsky, and M. Horowitz, "Dark Memory and Accelerator-Rich System Optimization in the Dark Silicon Era," IEEE Design & Test (submitted)
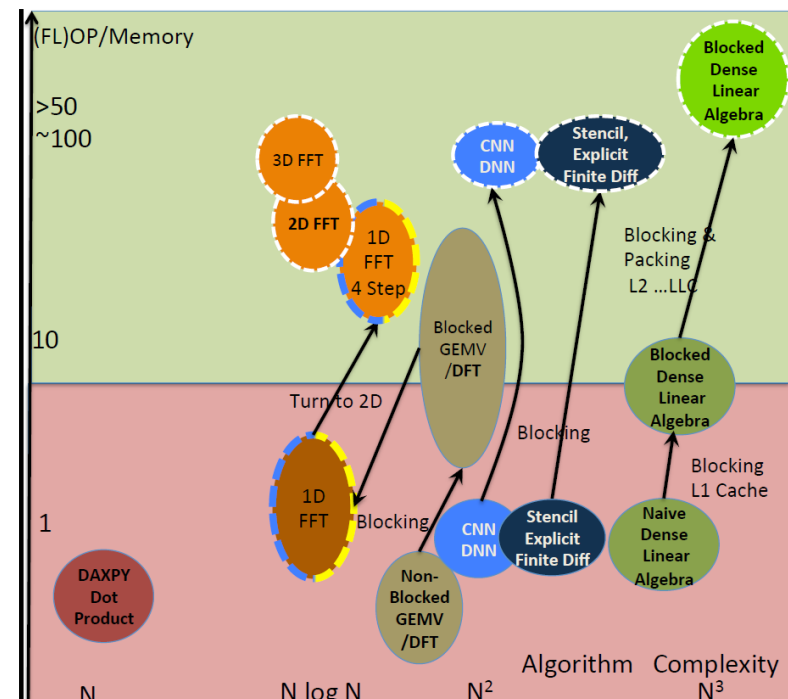
# Will Memristors Light the (Dark) Memory?

- Nonvolatility – low static energy

- Dense memory – short wires

- Still large memory -> relatively long wires, not a fundamental change in energy

# Fundamental Solution – SW-HW

- Minimizing memory accesses – algorithm execution

- High chip-level locality

- Memristive accelerators can help

A. Pedram, S. Galal, S. Richardson, S. Kvatinsky, and M. Horowitz, "Dark Memory and Accelerator-Rich System Optimization in the Dark Silicon Era," IEEE Design & Test (submitted)

# **Memristive Accelerators**

- Resistive Associative Processor

  (ReAP, Yavits et al. CAL 2015)

- Resistive GP-SIMD (Morad et al., TACO 2016)

- Neuromorphic (Soudry et al. TNNLS 2015)

- Memory Processing Unit (MPU, Kvatinsky et al.

  TVLSI 2014, TCAS II 2014, Levy et al. MEJ 2014)
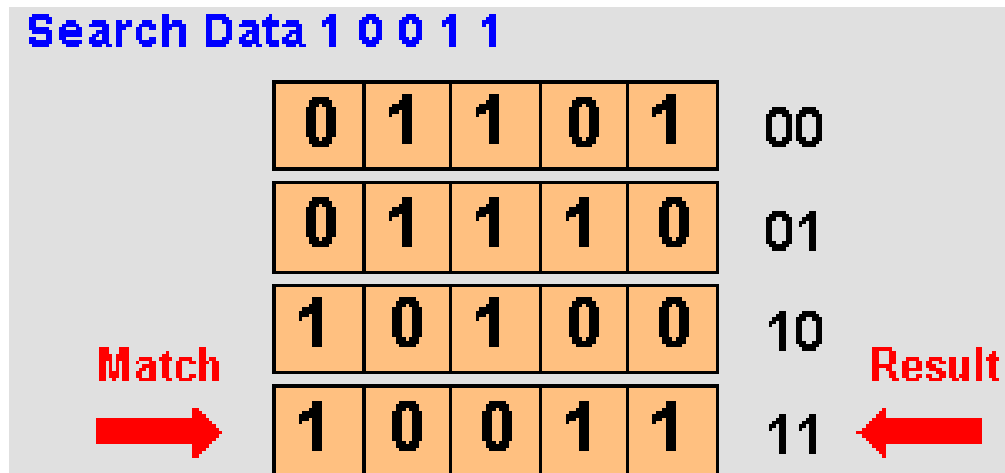
# Memristive Accelerators

- **Resistive Associative Processor**
  (ReAP, Yavits et al. CAL 2015)

- Resistive GP-SIMD (Morad et al., TACO 2016)

- Neuromorphic (Soudry et al. TNNLS 2015)

- Memory Processing Unit (MPU, Kvatinsky et al.
  TVLSI 2014, TCAS II 2014, Levy et al. MEJ 2014)

# Associative Processor

- Processing in-memory (PiM), using CAM

- AP is similar to a look-up table

- Computation is a series of "compare" and "write" operation



Search Data 1 0 0 1 1

| 0 | 1 | 1 | 0 | 1 | 00 |
| 0 | 1 | 1 | 1 | 0 | 01 |
| 1 | 0 | 1 | 0 | 0 | 10 |
| 1 | 0 | 0 | 1 | 1 | 11 |

Match →

Result ←

# Example: Associative Vector Addition

## ASSOCIATIVE PROCESSOR: MEMORY MAP

| | 255 | 12 | 11 | 8 | 7 | | | 4 | 3 | | | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| | | | | | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| | | | | | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| | | | | | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| | | | | | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| | | | | | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 |
| | | | | | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 |
| | | | | | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| | | | | | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| | | | | | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 |

**C   S   =   B   +   A**

25

# Example: Associative Vector Addition

## ASSOCIATIVE PROCESSOR: MEMORY MAP

| 255 | 12 11 | 8 7 | 4 3 | 0 |
|---|---|---|---|---|

| | | 0 0 0 1 | 0 1 0 0 | |
| | | 0 1 0 1 | 0 1 0 1 | |
| | | 0 0 0 0 | 0 0 1 0 | |
| | | 0 1 0 0 | 0 1 1 0 | |
| | | 1 0 0 1 | 0 0 0 0 | |
| | | 0 1 0 0 | 0 1 0 1 | |
| | | 0 0 1 1 | 1 1 0 1 | |
| | | 0 0 0 0 | 0 1 0 1 | |
| | | 1 0 0 1 | 0 0 1 0 | |
| | | 0 0 1 1 | 1 0 1 0 | |

**C    S    =    B    +    A**

| cout | s | $c_{in}$ | a | b |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

26

# Example: Associative Vector Addition

SELECTING BIT COLUMN 0

| 255 | 12 | 11 | | | 8 | 7 | | | 4 | 3 | | | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | | | | | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| | 0 | | | | | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| | 0 | | | | | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| | 0 | | | | | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| | 0 | | | | | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| | 0 | | | | | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 |
| | 0 | | | | | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 |
| | 0 | | | | | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| | 0 | | | | | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| | 0 | | | | | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 |

**C   S   =   B   +   A**

| cout | s | $c_{in}$ | a | b |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

27

# Example: Associative Vector Addition

COMPARE



| cout | s | $c_{in}$ | a | b |
|------|---|----------|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

C   S   =   B   +   A

# Example: Associative Vector Addition

WRITE



|255|12|11|8|7| | | |4|3| | | |0|
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| cout | s | $c_{in}$ | a | b |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

C   S   =   B   +   A

# Example: Associative Vector Addition

COMPARE



| cout | s | $c_{in}$ | a | b |
|------|---|----------|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

C   S   =   B   +   A

# Example: Associative Vector Addition

WRITE



| cout | s | $c_{in}$ | a | b |
|------|---|------|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

C   S   =   B   +   A

# Example: Associative Vector Addition

**SELECTING BIT COLUMN 1**



| cout | s | $c_{in}$ | a | b |
|------|---|----------|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

C   S   =   B   +   A

36

# Example: Associative Vector Addition

END OF COMPUTATION



C  S  =  B  +  A

| cout | s | $c_{in}$ | a | b |
|------|---|----------|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

# AP Complexity

- Arithmetic:
  - Fixed point
    - $m$ bit add / sub:   $O(m)$ cycles
    - $m$ bit mult/div:    $O(m^2)$ cycles
- Pattern match:          $O(1)$ cycles
- Finding max/min:        $O(1)$ cycles
- Independent of the dataset size:
  <span style="color:red">The larger the problem, the better the performance of the Associative Processor!</span>

# Resistive Associative Processor



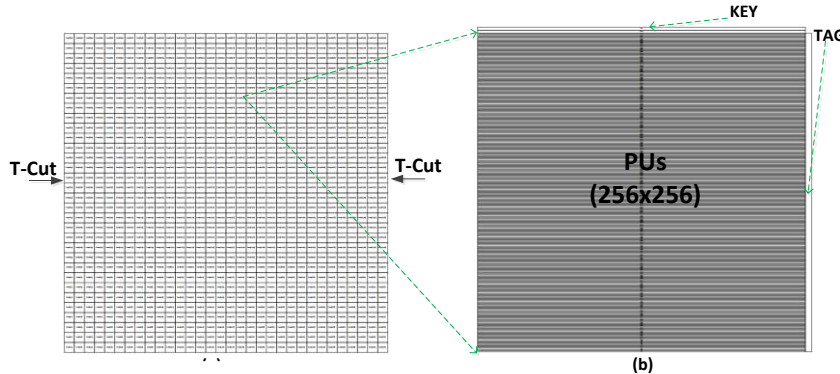- Converting a memory crossbar into a massively parallel SIMD processor
- ❖ Enabling a 100M PU-AP

# What AP is Good for

- Dense and sparse linear algebra
- K-means clustering
- Linear SVM classification
- FFT, convolution, feature extraction
- Sequence alignment (Smith-Waterman)
- Graph processing (Dijkstra's shortest path finding)
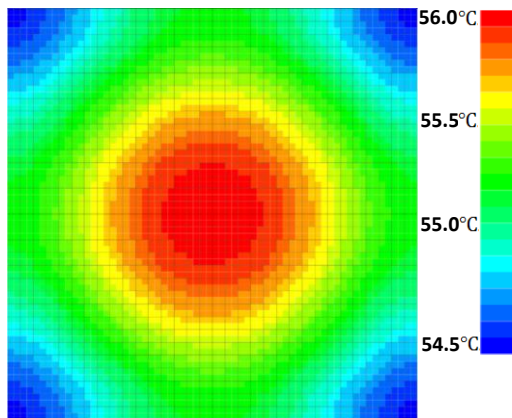
# **Performance and Power Consumption**



- ReAP size (and consequently performance) are constrained by memristor write energy

- Max Dense Matrix Multiplication performance is 5TFLOPS under this constraint
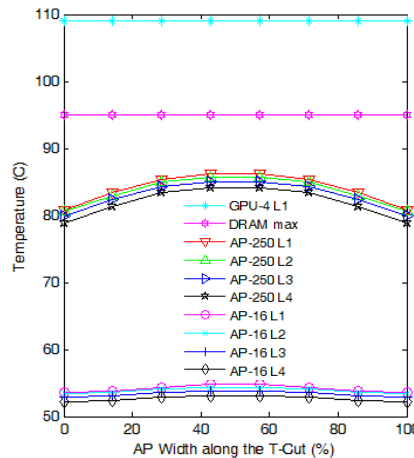
# Thermal View



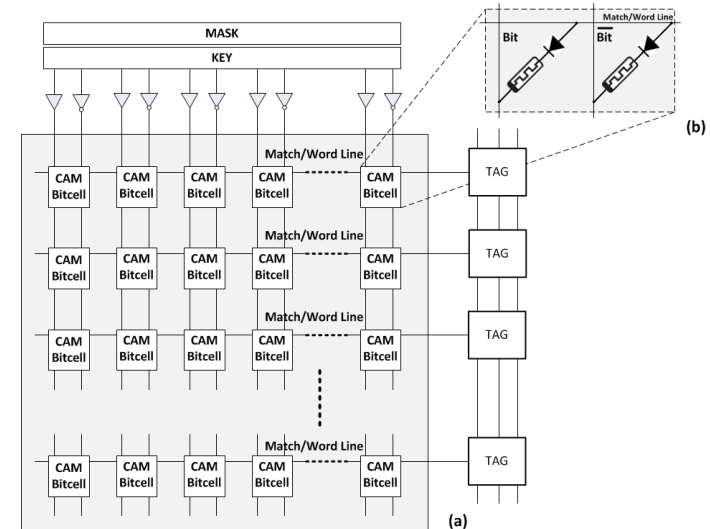ReAP Floorplan



ReAP Thermal Map



ReAP Temperature (vs. DRAM)

- Temperature and hot spots are the reason 3D integration of CPUs and DRAM is stalling

- AP does not have this problem due to its (almost) uniform thermal distribution

# Summary

- The dark (silicon and memory) age

  – Main source of inefficiency is data movement

- The solution: accelerators and HW-SW

  awareness

- Memristive accelerators!



43

# Thanks!

[shahar@ee.technion.ac.il](mailto:shahar@ee.technion.ac.il)