

# Hot Chips: Atoms to Heat Sinks

## ECE 598EP

**Prof. Eric Pop**  
Dept. of Electrical and Computer Engineering  
Univ. Illinois Urbana-Champaign

<http://poplab.ece.uiuc.edu>



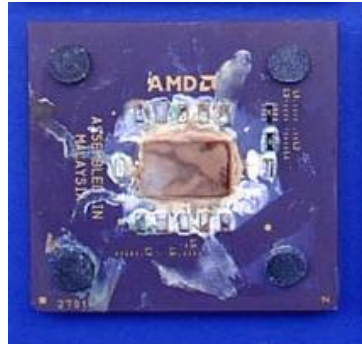
## The Big Picture



**AMD**  
XP1500+ CPU

[http://phys.ncku.edu.tw/~htsu/humor/fry\\_egg.html](http://phys.ncku.edu.tw/~htsu/humor/fry_egg.html)

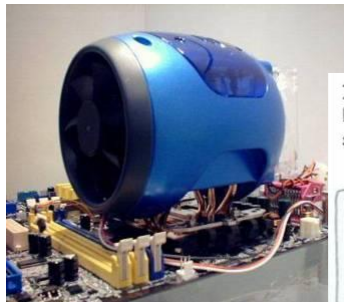
## Another CPU without a Heat Sink



Source: Tom's Hardware Guide

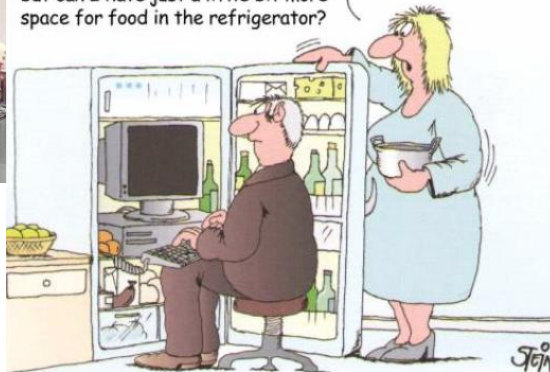
<http://www6.tomshardware.com/cpu/01q3/010917/heatvideo-01.html>

## Thermal Management Methods



ASUSTeK cooling solution (!)

I believe that your CPU needs extra cooling  
but can I have just a little bit more  
space for food in the refrigerator?



# Impact on People & Environment

- Fast computers run HOT
- COOL computers are slow...
- Huge data centers need significant power generation and cooling investment
- Impact on environment?!

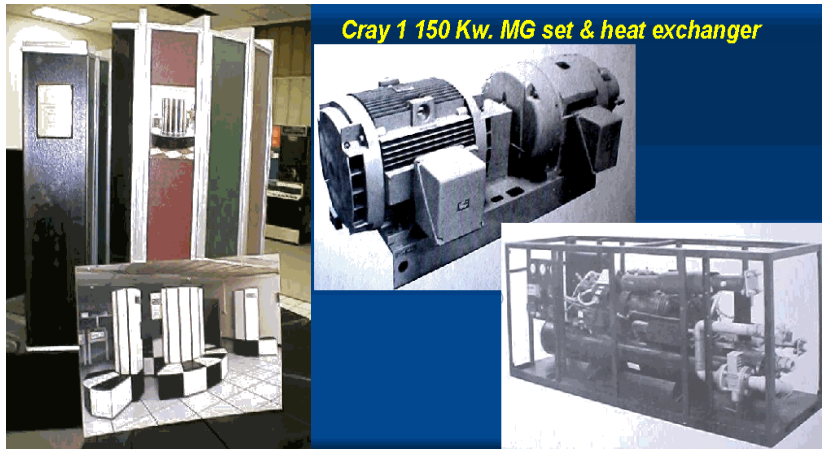
The industry now calls them "portables" or "notebooks" not "laptops"

The screenshot shows a BBC News page with the following content:

- Navigation: NEWS SPORT WEATHER WORLD SERVICE A-Z INDEX
- Header: BBC NEWS WORLD EDITION
- Location: You are in: Health
- Date/Time: Friday, 22 November, 2002, 12:55 GMT
- Section: News Front Page
- Article Title: Burned groin blamed on laptop
- Image: A man in a suit holding a laptop.
- Text: Hot stuff: Could laptopping be a painful business? A Swedish scientist who rested his laptop computer on his lap for just an hour needed medical treatment for extensive blistering.
- Section: Talking Point
- Text: A concerned doctor wrote to The Lancet medical journal after encountering the distressed patient.
- Section: Country Profiles
- Section: In Depth
- Text: He is warning the public of the potential dangers of using a laptop "in the literal sense".
- Section: Programmes
- Text: The 50-year-old father-of-two used the laptop machine, of unknown origin, to write a report while sitting in an armchair.
- Section: BBC SPORT
- Section: BBC WEATHER
- Section: SERVICES
- Text: Dr Claes-Göran Ostenson, from Sweden's Karolinska Institute, told the journal: "He had placed his laptop computer on his lap while writing for about one hour. "The next day he noticed irritation."

# Packaging cost

From Cray (local power generator and refrigeration)...



<http://www.research.microsoft.com/users/gbell/craytalk/>

# Packaging cost

To today...

- Grid computing: power plants co-located near computer farms
- IBM S/390: refrigeration

350-V bulk power subassembly  
(under cover)

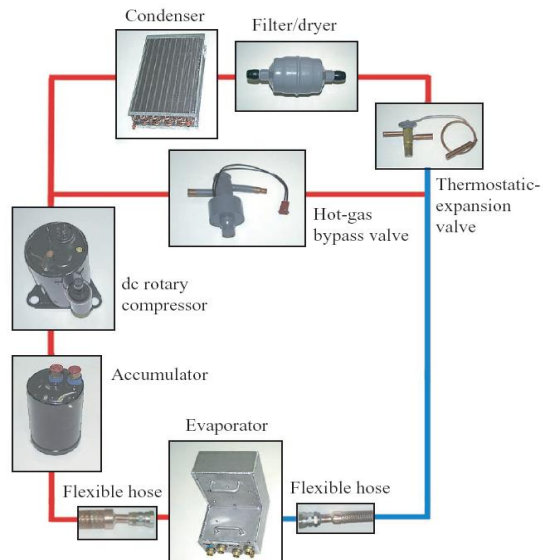


- Processor cage:
  - Contains processors, memory and I/O
  - Dual redundant three-phase line cord
  - Distributes power to system
  - Five I/O slots
- MCM/evaporator
- 350-V integrated battery
- Modular cooling unit
- Expansion cage:
  - Powered from processor cage
  - Twenty-two I/O slots

Source: R. R. Schmidt, B. D. Notohardjono "High-end server low temperature cooling"  
IBM Journal of R&D

# IBM S/390 refrigeration

- Complex and expensive

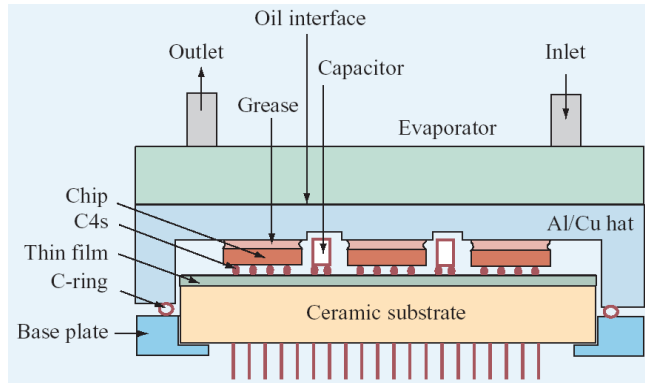


Source: R. R. Schmidt, B. D. Notohardjono "High-end server low temperature cooling" *IBM Journal of R&D*

# IBM S/390 processor packaging

Processor subassembly: complex!

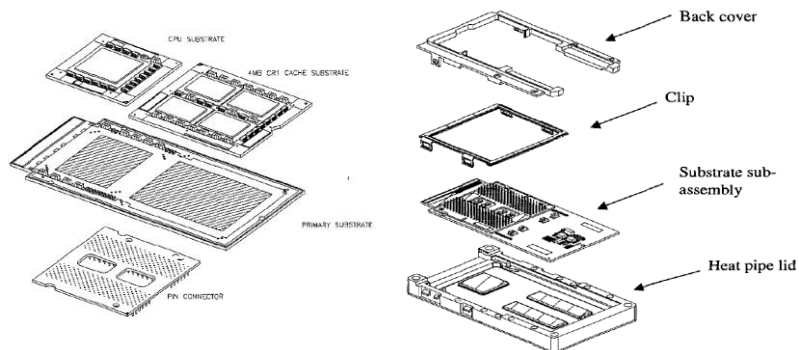
C4: Controlled Collapse Chip Connection (flip-chip)



Source: R. R. Schmidt, B. D. Notohardjono "High-end server low temperature cooling"  
IBM Journal of R&D

# Intel Itanium packaging

Complex and expensive (note heatpipe)

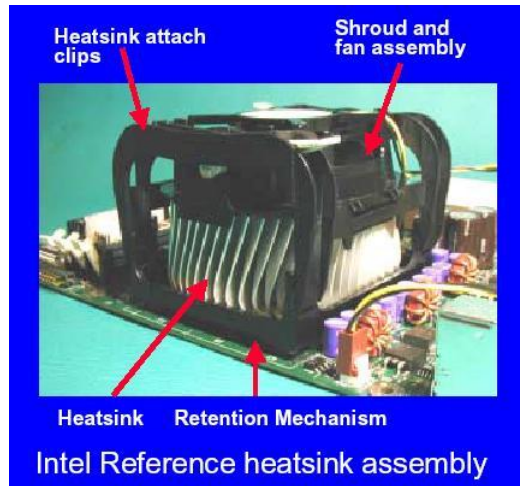


Source: H. Xie et al. "Packaging the Itanium Microprocessor"  
Electronic Components and Technology Conference 2002



## Intel Pentium 4 packaging

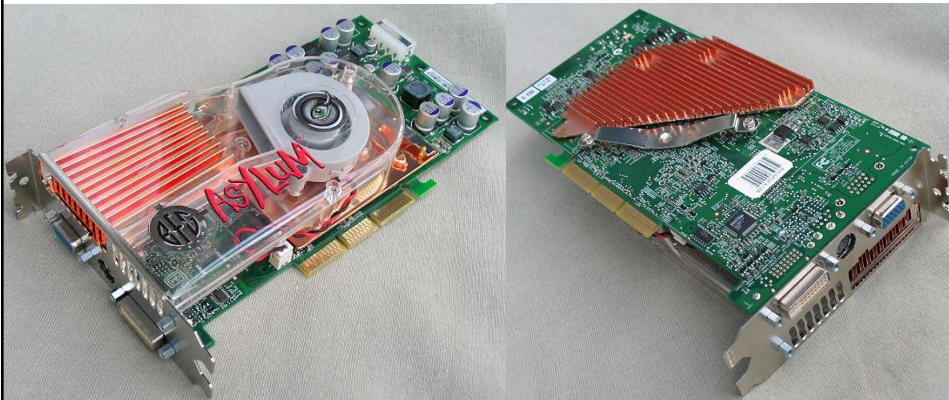
- Simpler, but still...



Source: Intel web site

## Graphics Cards

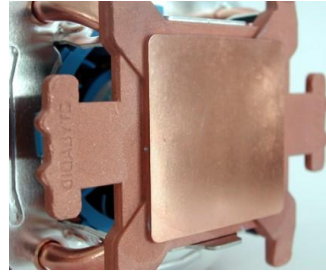
- Nvidia GeForce 5900 card



Source: Tech-Report.com

# Under/Overclocking

- Some chips need to be underclocked
  - Especially true in constrained form factors

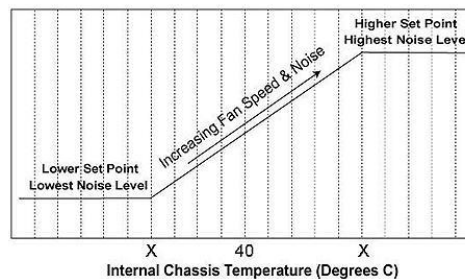


- Try fitting this in a laptop or Gameboy!

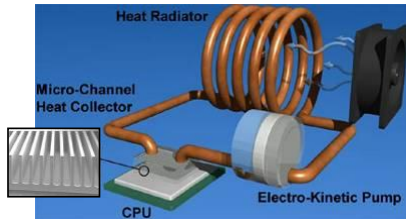
Ultra model of Gigabyte's 3D Cooler Series  
Source: Tom's Hardware Guide

# Environment

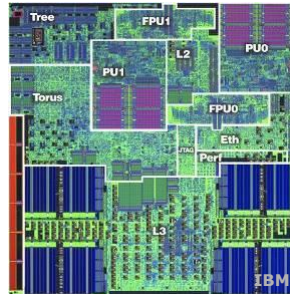
- Environment Protection Agency (EPA): computers consume 10% of commercial electricity consumption
  - This incl. peripherals, possibly also manufacturing
  - A DOE report suggested this percentage is much lower
  - No consensus, but it's probably significant
- Equivalent power (with only 30% efficiency) for AC
- CFCs used for refrigeration
- Lap burn
- Fan noise



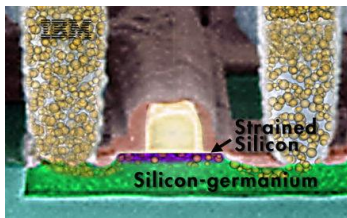
# Thermal Management Methods



**System Level**  
→ Active Microchannel Cooling (Cooligy)

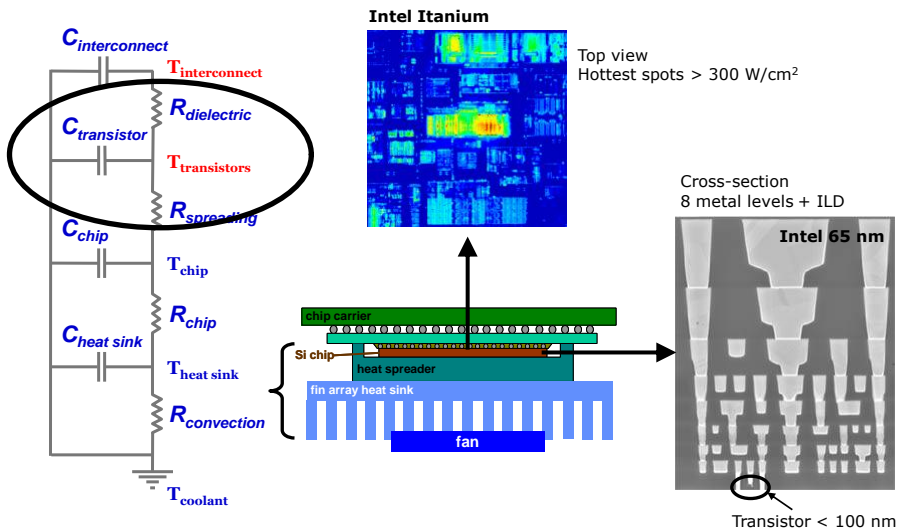


**Circuit + Software Level**  
→ active power management  
(turn parts of circuit on/off)



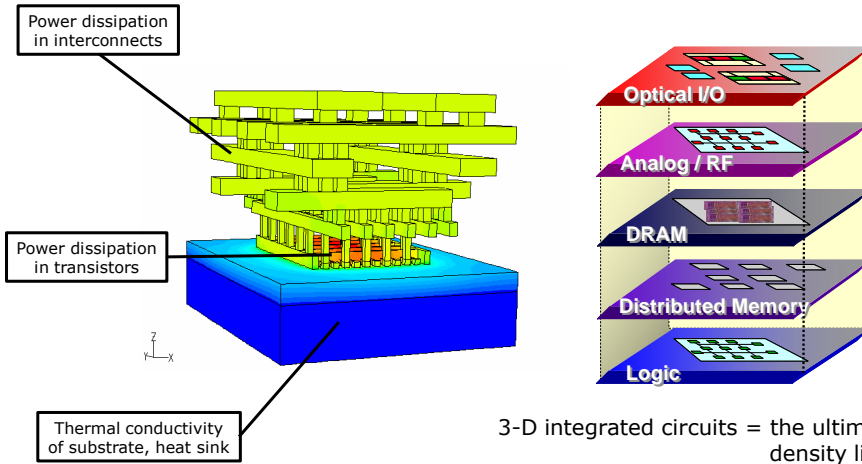
**Transistor Level**  
→ electro-thermal device design

# Where Does the Heat Come From?



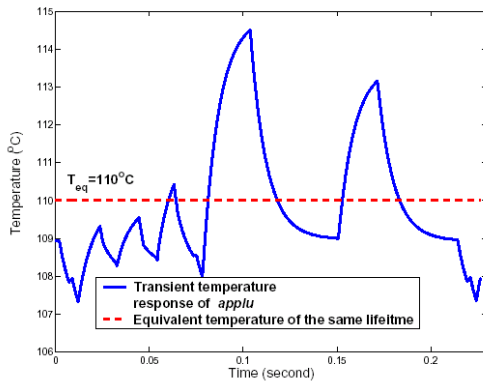


# More on Chip-Level Complexity

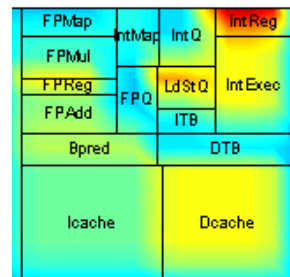


3-D integrated circuits = the ultimate density limit  
 How do we get the power in?  
 How do we take the heat out?

# Temporal, Spatial Variations



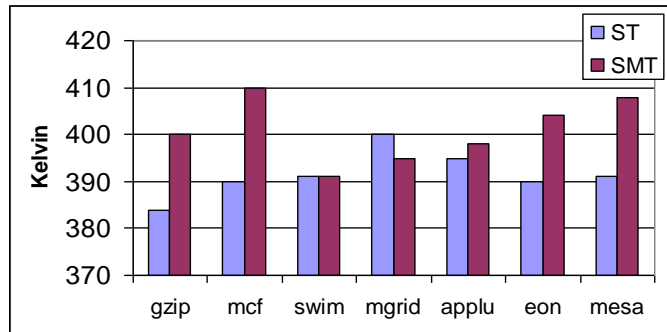
Temperature variation of SPEC applu over time



Hot spots increase cooling costs  
 ⇒ must cool for hot spot

## Variations Depending on Application

- Wide variation across applications
- Architectural and technology trends are making it worse, e.g. *simultaneous multithreading* (SMT)
  - Leakage is an especially severe problem: *exponentially dependent on temperature!*



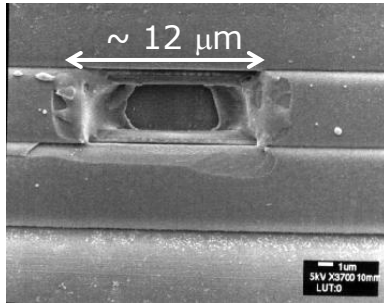
## Temperature Affects:

- Circuit performance
- Circuit power (leakage exponential)
- IC reliability (exponential)
- IC and system packaging cost
- Environment

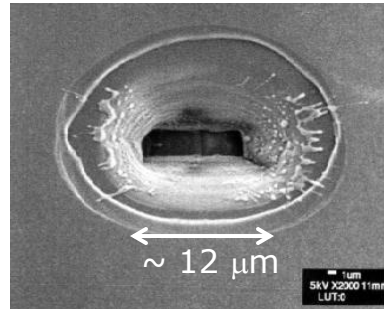
# Thermal Interconnect Failure

## Open Circuit Interconnect Failure

Banerjee, Kim, Amerasekera, Hu, Wong, and Goodson, IRPS 2000



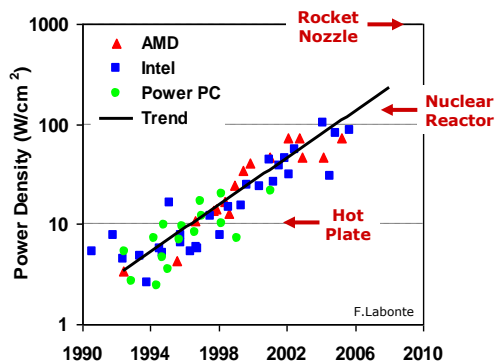
Metal 4



Metal 1

- Passivation fracture due to the expansion of critical volume of molten AlCu. (@ 1000 °C)

# Chip-Level Thermal Challenges



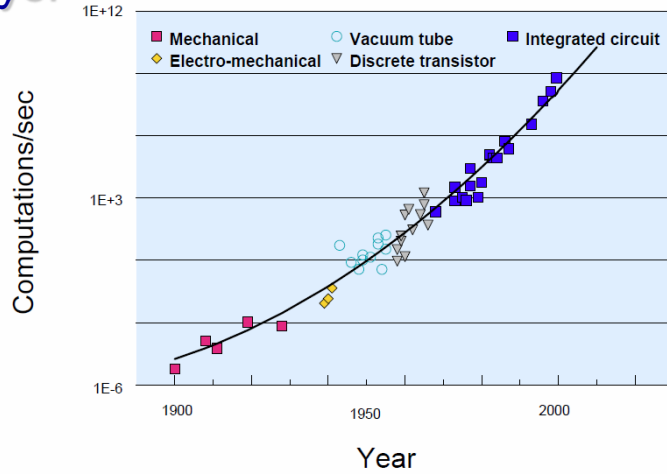
Device Level:  
Confined Geometries, Novel Materials



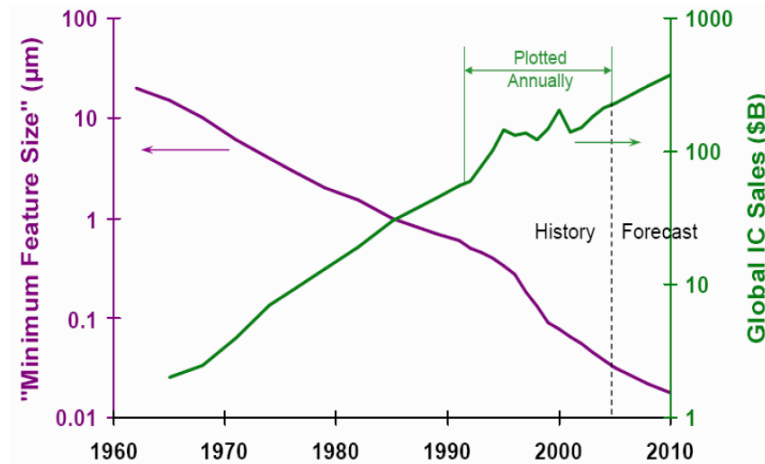
Material	$k$ (W/m/K)
Si	148
Ge	60
Silicides	40
Si (10 nm)	13
SiO <sub>2</sub>	1.4

Source: E. Pop (Proc IEEE 2006)

# Why (down)Scaling? To increase speed & complexity! \$1000 buys:



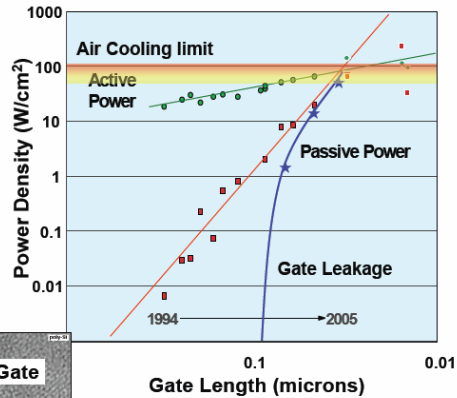
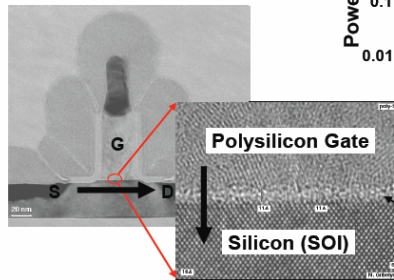
# Scaling = Progress in Electronics



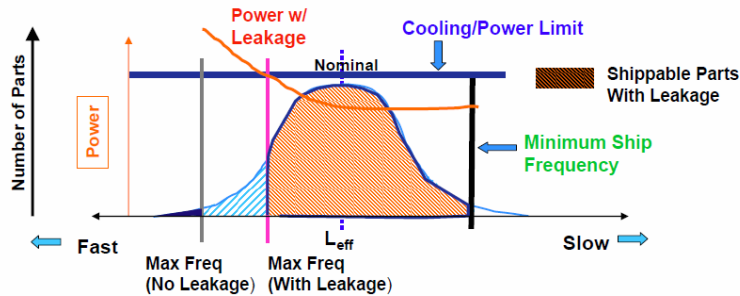
Smaller features → Better performance & cost/function  
→ More apps → Larger market

# CMOS Power Issue: Active vs. Passive

- ◆ Power components:
  - ❖ Active power
  - ❖ Passive power
    - Gate leakage
    - Source - Drain leakage



# Power & Heat Limit Frequency Scaling

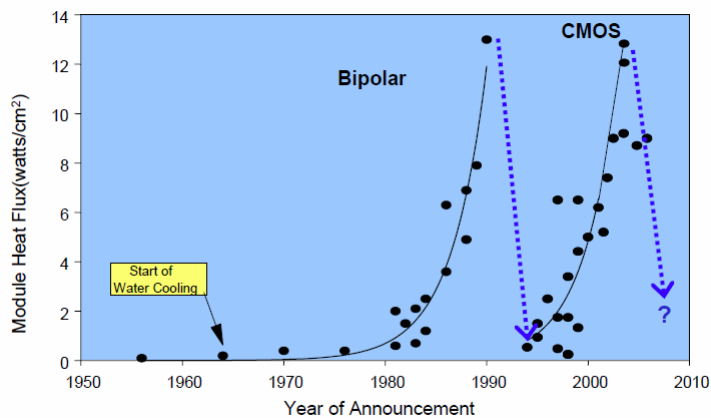


- ◆ Server microprocessors cannot simultaneously utilize all their transistors at full frequency due to power
  - ❖ Workload demands are highly variable – must exploit power management
- ◆ Moore's law is continuing for transistor density (maybe at a reduced pace)
  - *New methods to utilize silicon density scaling will have to be developed to accommodate diverse workloads while managing power constraints*

# Industry Developed ITRS Guide (Intl. Technology Roadmap for Semic.) <http://www.itrs.net>

Production Year:	2001	2004	2007	2010	2013	2016 ...
DRAM Half-Pitch [nm]:	130	90	65	45	32	22
Overlay Control [nm]:	45	32	16	11	8	5.5
Gate Length [nm]:	65	37	25	18	13	9
CD Control [nm]:	6.3	3.3	2.6	1.9	1.3	0.9
T <sub>OX</sub> (equivalent) [nm]:	1.3-1.6	1.2	1.1	0.65	0.5 (UTB)	0.5 (MUG)
I <sub>ON</sub> (NMOS) [ $\mu\text{A}/\mu\text{m}$ ]:	900	1110	1200	2050	2198	2713
I <sub>OFF</sub> (NMOS) [ $\mu\text{A}/\mu\text{m}$ ]:	0.01	0.05	0.2	0.28	0.29	0.11
Interconnect K <sub>EFF</sub> :	-	3.1-3.6	2.7-3.0	2.5-2.8	2.1-2.4	1.9-2.2

## Has This Ever Happened Before?





## Implications for Nanoscale Circuits

- ◆ **Circuit heat generation is the main limiting factor for scaling of device speed and switch circuit density**
- ◆ Scaling to molecular dimensions may not yield performance increases
  - ❖ We will be forced to trade-off between speed and density
- ◆ Optimal dimensions for electronic switches should range between 5 and 50 nm
  - ❖ Likely achievable with Si – easily within the scope of ITRS projections
- ◆ Going to other materials for FETs will likely achieve only “one-time” percentage gains
  - ❖ Need a new device mechanism for new scaling path

## Transistor-Level Thermal Challenges

- Small geometry
  - High power density (device-level hot spot)
  - Higher surface-to-volume area, i.e. higher role of thermal interfaces between materials
- Lower thermal conductivity
- Lowering power (but can it ever be low enough?)
- Device-level thermal design (phonon engineering)

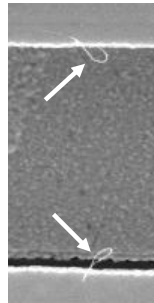
*Device Level:  
Confined Geometries, Novel Materials*



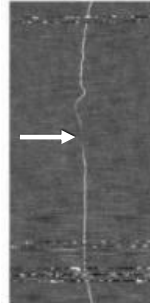
Material	$k$ (W/m/K)
Si	148
Ge	60
Silicides	40
Si (10 nm)	13
SiO <sub>2</sub>	1.4

Source: E. Pop (Proc IEEE 2006)

# The Tiny Picture



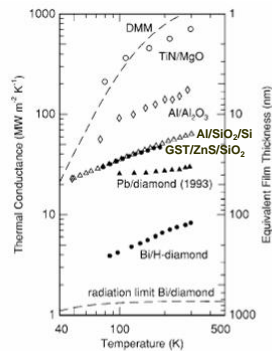
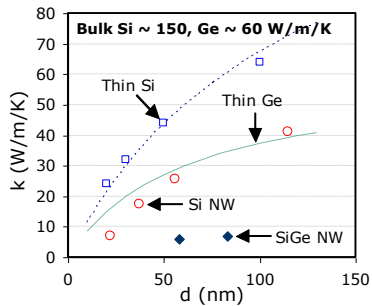
Suspended



On substrate

Carbon nanotubes burn at high enough applied voltage

# K of Nano{wires;layers}, $R_B$ of Interfaces



Thermal conductivity (K) of thin films and nanowires:

- Decrease due to phonon confinement and boundary scattering
- Up to an order of magnitude decrease from bulk values

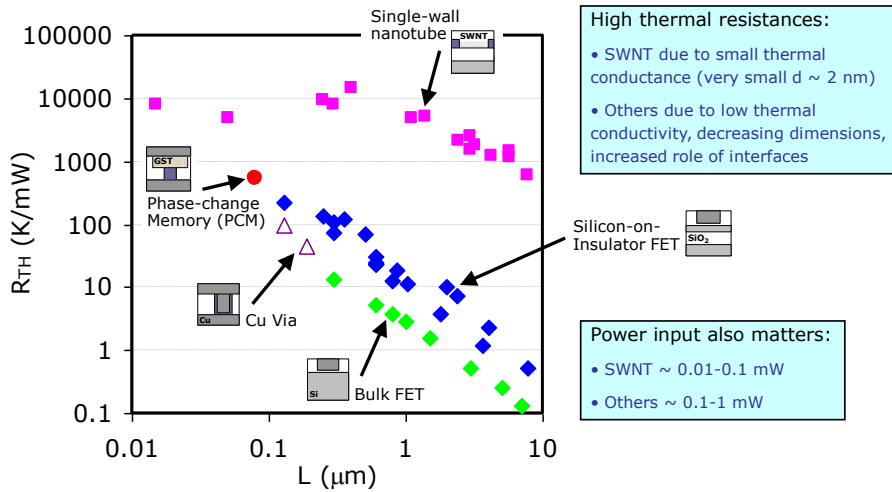


Thermal interface resistance  $\sim 10$  nm SiO<sub>2</sub>

Data: Li (2003), Liu (2005); Model: Pop (2004)

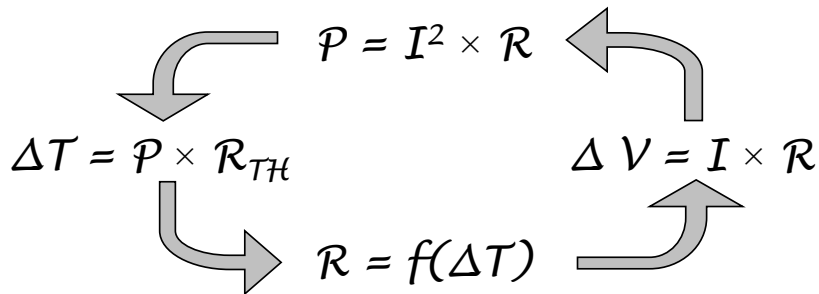
Lyeo (2006)

# Thermal Resistance at Device Level



Data: Mautry (1990), Bunyan (1992), Su (1994), Lee (1995), Jenkins (1995), Tenbroek (1996), Jin (2001), Reyboz (2004), Javey (2004), Seidel (2004), Pop (2004-6), Maune (2006).

# Thermal Resistance, Electrical Resistance



Fourier's Law (1822)



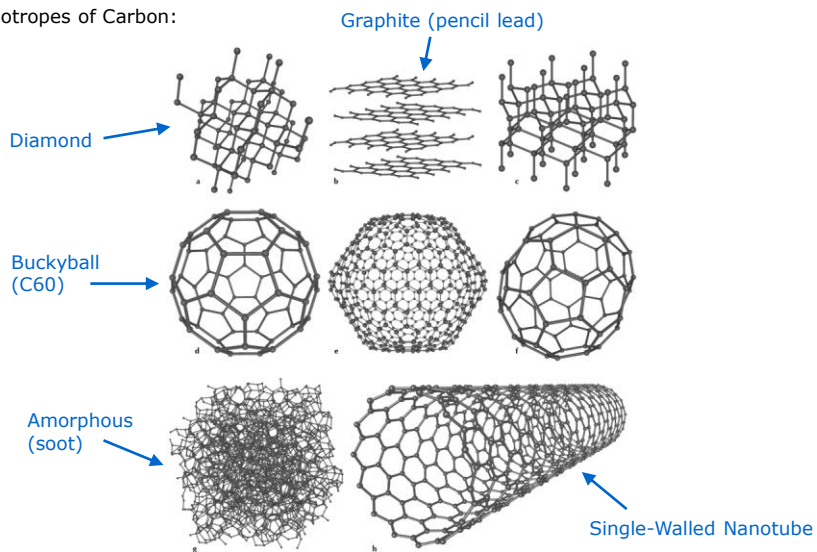
Ohm's Law (1827)

# This Heating Business is Not All Bad...

*IF we can control it!*

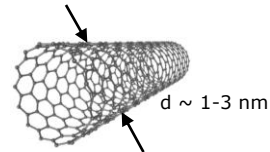
## Nanotubes in the Carbon World

Allotropes of Carbon:

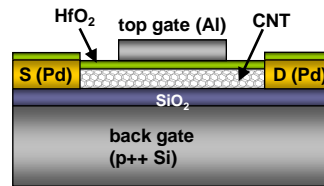


# Why Carbon Nanotubes & Graphene?

- Carbon nanotube = rolled up graphene sheet
- Great electrical properties
  - Semiconducting → Transistors
  - Metallic → Interconnects
  - Electrical Conductivity  $\sigma \approx 100 \times \sigma_{Cu}$
  - Thermal Conductivity  $k \approx k_{diamond} \approx 5 \times k_{Cu}$



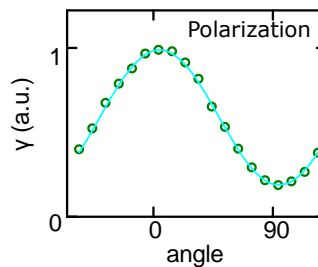
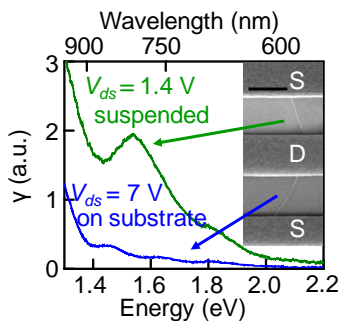
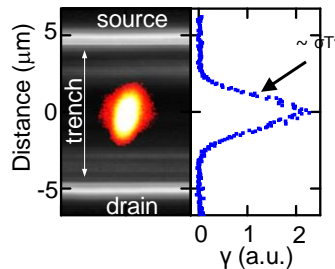
- Nanotube challenges:
  - Reproducible growth
  - Control of electrical and thermal properties
  - Going “from one to a billion”



# Light Emission from Metallic SWNTs

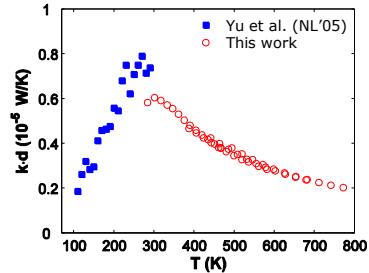
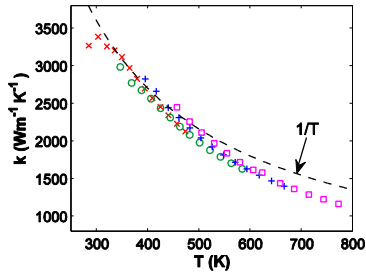
D. Mann *et al.*, *Nature Nano* 2, 33 (2007)

- Joule-heated tubes emit light:
  - Comes from center, highly polarized
  - Emitted photons at higher energy than applied bias (high energy tail)
  - World’s smallest light bulb?



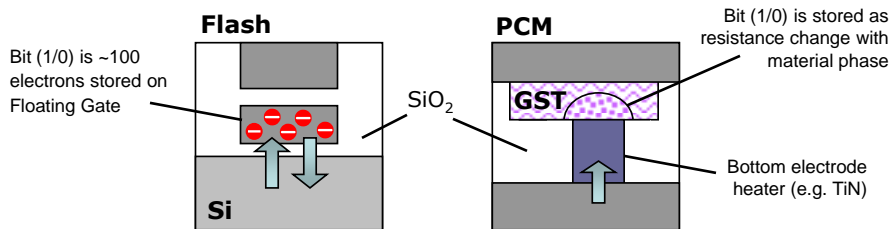
# Extracting SWNT Thermal Conductivity

E. Pop et al., *Nano Letters* 6, 96 (2006)



- Numerical extraction of  $k$  from the high bias ( $V > 0.3$  V) tail
- Comparison to data from 100-300 K of UT Austin group (C. Yu, *NL Sep'05*)
- Result: first “complete” picture of SWNT thermal conductivity from 100 – 800 K

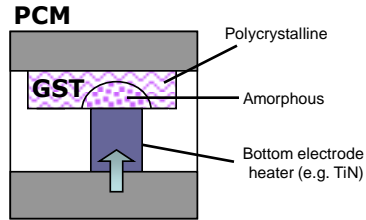
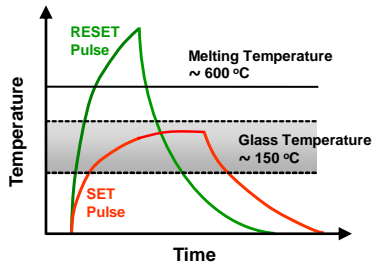
# What Is Phase-Change Memory?



- PCM: Like Flash memory (non-volatile)
- PCM: Unlike Flash memory (resistance change, not charge storage)
- Faster than Flash (100 ns vs. 0.1–1 ms), smaller than Flash (which is limited by ~100 electrons stored/bit)
- For: iPod nano, mobile phones, PDAs, solid-state hard drives...



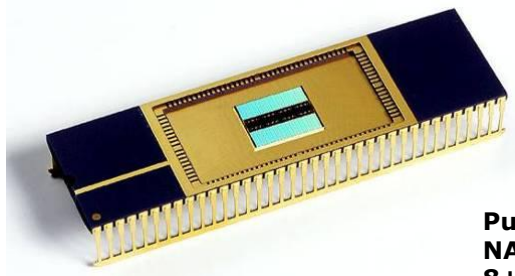
# How Phase-Change Memory Works



- Based on  $\text{Ge}_2\text{Sb}_2\text{Te}_5$  reversible phase change:  $R_{\text{amorph}} / R_{\text{xtal}} > 100$
- Short (10 ns), high pulse (0.5 mA) melts, amorphizes GST
- Longer (100 ns), lower pulse (0.1 mA) crystallizes GST
- Small cell area (sits on top of heater), challenge is reliability and lowering programming current (BUT, helped by scaling!)

# Samsung 512 Mb PCM Prototype

Sep 11, 2006



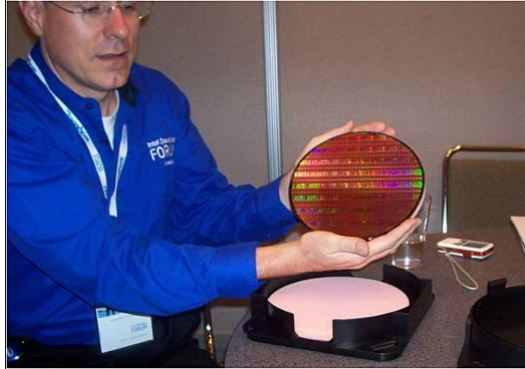
**Put in perspective:  
NAND Flash chips of  
8+ Gb in production**

"Samsung completed the first working prototype of what is expected to be the main memory device to replace high density Flash in the next decade – a Phase-change Random Access Memory (PRAM). The company unveiled the 512 Mb device at its sixth annual press conference in Seoul today." Source:

[http://samsung.com/PressCenter/PressRelease/PressRelease.asp?seq=20060911\\_0000286481](http://samsung.com/PressCenter/PressRelease/PressRelease.asp?seq=20060911_0000286481)

# Intel/ST Phase-Change Memory Wafer

Sep 28, 2006



“Intel CTO of Flash Memory Ed Doller holds the first wafer of 128 Mbit phase change memory (PCM) chips, which has just been overnigheted to him from semiconductor maker STMicroelectronics in Agrate, Italy. Intel believes that PCM will be the next phase in the non-volatile memory market.” Source: <http://www.eweek.com/article2/0,1895,2021841,00.asp>