

VLSI microarchitecture

Scaling

Toni Juan

Short index

- **Administrative**
 - glass schedule
 - projects?

- **Topic of the day: **Scaling****

Administrative

Class Schedule

Day	Topic	Lecturer
06 Nov 00	Introduction Tech. & μ Arch.	T+E
13 Nov 00	Introduction To MOS VLSI	E
20 Nov 00	Scaling of wires & Trans.	T
27 Nov 00	Implem./Techn. MOS	E
04 Dec 00	Delay μ Architecture	T
11 Dec 00	Delay VLSI	E
18 Dec 00	Power μ Architecture	T
08 Dec 00	Power VLSI	E
15 Jan 01	Area μ Architecture	T
18 Jan 01	Case studies & alternatives	T+E

Projects

- In groups of two people

- Take your favorite structure and:
 - architect several implementations (block level)
 - find critical path
 - delay estimation
 - delay measurement (SPICE)
 - power estimation
 - area estimation
 - write report

Introduction

Why Scaling?

- **We care about technology improvements**

- **Scaling (reduction of minimum feature size)**
 - faster devices
 - more devices per area unit

- **Other important things (that we won't consider...)**
 - die size
 - metal levels
 - new conductors and dielectrics
 - ...

SCALING for Architects

- **What is scaling**
 - summary of impact of technology improvements
- **How is it measured, what does it mean**
- **Scaling for**
 - gates
 - wires
- **Limits for scaling**
- **Impact on**
 - delay
 - power
 - area

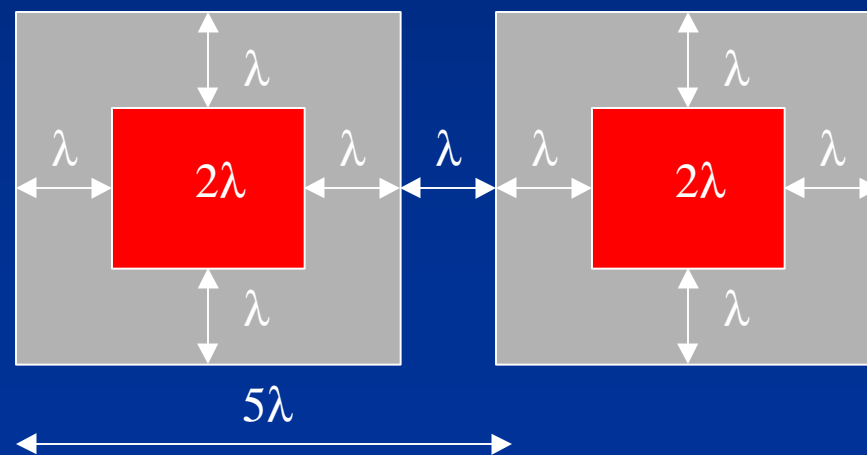
Scaling

Disclaimer

- We refer to CMOS VLSI scaling
- Very simple models for μ architects
 - Gates
 - Wires

Scaling factor and technology

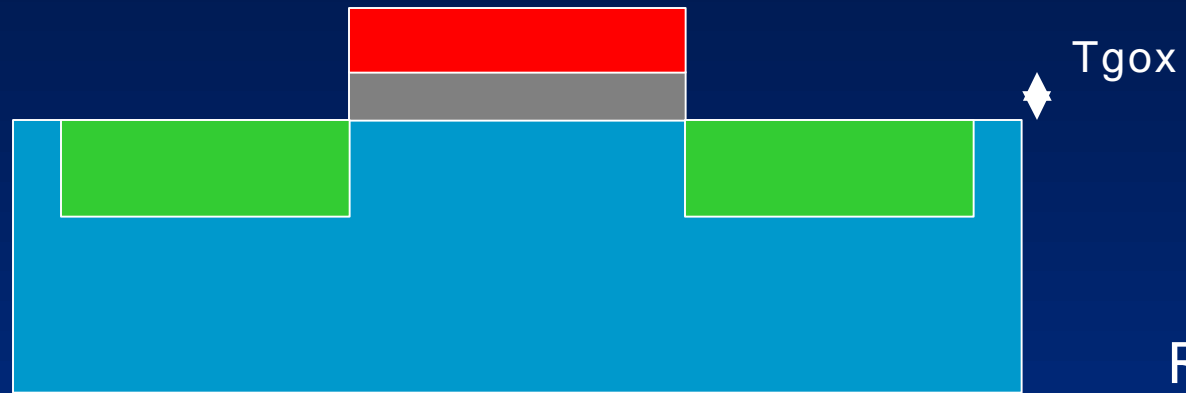
- Depends on the lithography and components tech.
- **f**: Minimum feature size
 - $f = 2\lambda$
 - λ is the mask registration variation
 - smallest device is $5\lambda \times 5\lambda$
 - S: scale = $f_{\text{old}}/f_{\text{new}}$



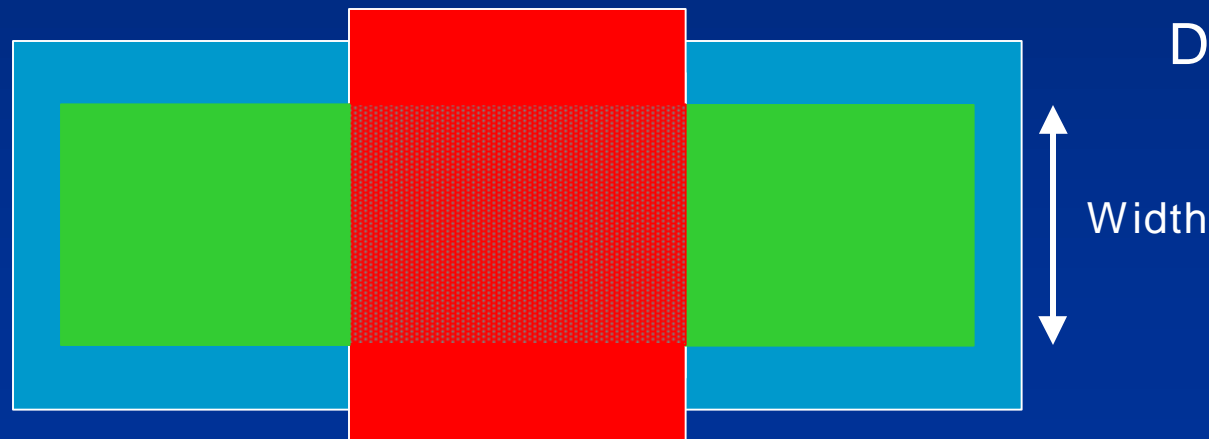
Gate scaling

How does the gate speed depend on technology?

Shape of transistors



Length



$$R1 = L/W$$

$$C1 = W \times L/T_{gox}$$

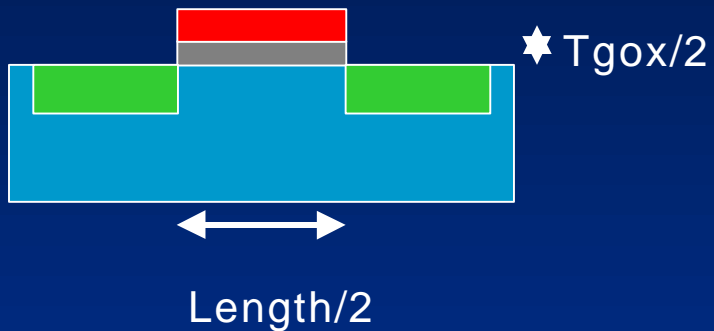
$$\text{Delay} = R1 \times C1$$

Ideal scaling of MOS transistors

Dimensions (W, L, T _{gox})	1/S
Voltages (V _{dd} , V _{tn} , V _{tp})	1/S
Current per device (I _{ds})	1/S
Gate capacitance	1/S
Transistor on-resistance	1
Intrinsic gate delay	1/S
Area per device (A= WxL)	1/S ²
Power dissipation per gate (IxV)	1/S ²
Power-delay product per gate	1/S ³
Power-dissipation density (IxV/A)	1

Scaling transistors

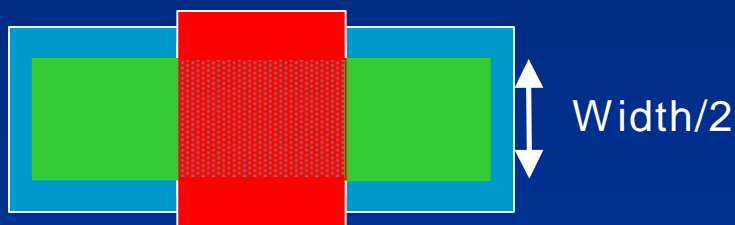
$$\alpha = 0.5 = \frac{1}{2} \quad (S = 2)$$



$$R_2 = (L/S) / (W/S) = R_1$$

$$C_2 = (L/S \times W/S) / (T_{gox}/S) = C_1/S$$

$$\text{Delay} = R_2 \times C_2 = R_1 \times C_1/S$$

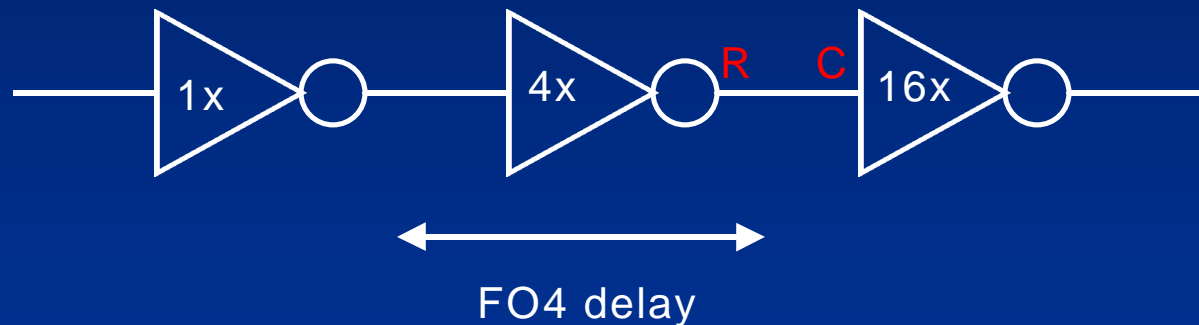


Gate delay (ideal scaling)

- R remains constant
- C is divided by the scaling factor
- **Delay = RC** is divided by S when scaling by S

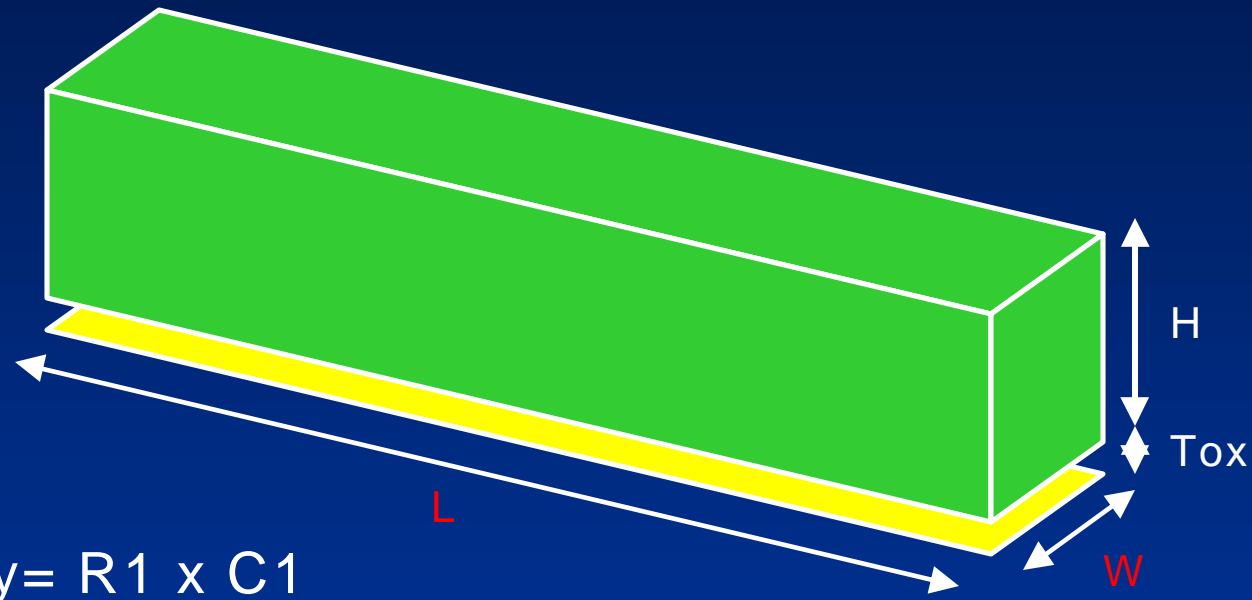
Fanout of 4 inverter metric

- Delay of an inverter with $C_{out}/C_{in} = 4$



Wire scaling

Shape of a wire



$$\text{Delay} = R1 \times C1$$

$$R1 = L / (W \times H)$$

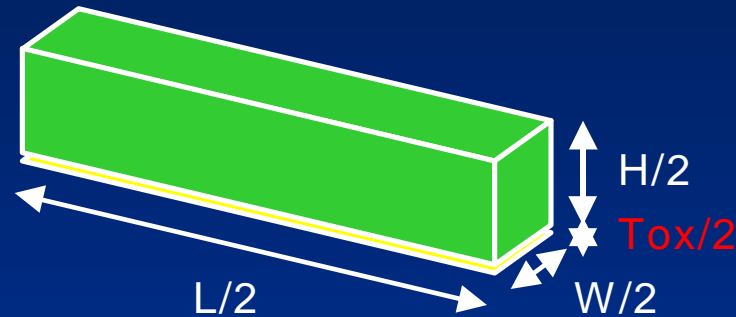
$$R1 / \mu\text{m} = \rho / (W \times H)$$

$$C1 = L \times W$$

$$C1 / \mu\text{m} = \epsilon \times W$$

Scaling of a wire (uniform)

$$\alpha = 0.5 = \frac{1}{2} \quad (S = 2)$$



$$\text{Delay} = R_2 \times C_2 = R_1 \times C_1$$

$$R_2 = R_1 \times S$$

$$C_2 = C_1 / S$$

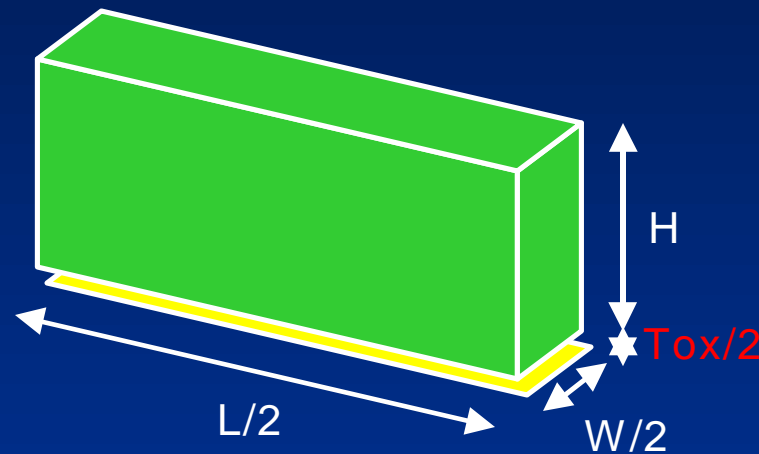
$$R_2 / \mu\text{m} = R_1 \times S^2$$

$$C_2 / \mu\text{m} = C_1$$

Scaling of a wire (non uniform)

Aspect ratio= 2

$\alpha = 0.5 = \frac{1}{2}$ (S= 2)



$$\text{Delay} = R_3 \times C_3 = R_1 \times C_1 / S$$

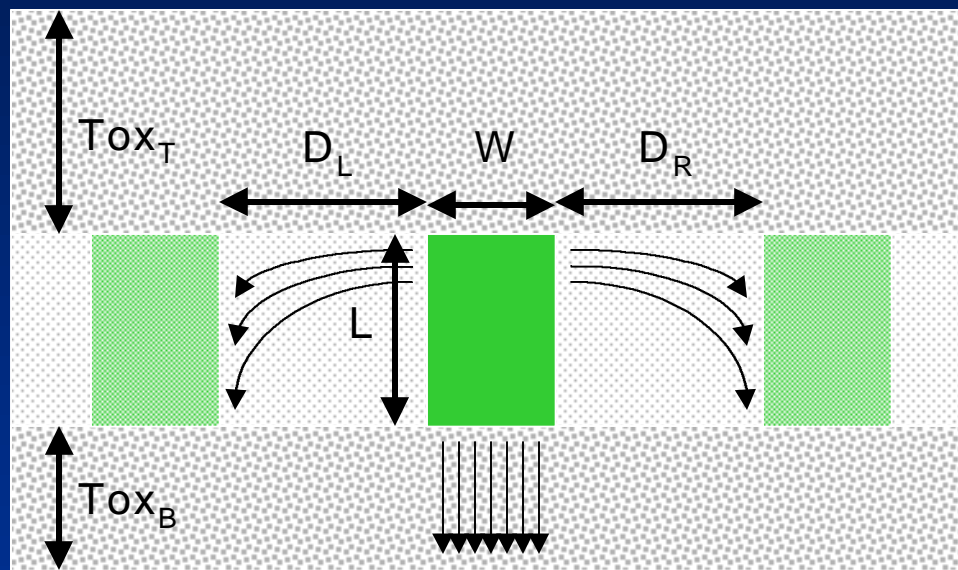
$$R_3 = R_1$$

$$R_3 / \mu\text{m} = R_1 \times S$$

$$C_3 = C_1 / S$$

$$C_3 / \mu\text{m} = C_1 / S^2$$

Not so easy ...



Even worse

□ Three types of levels

- M1
 - ◆ Local interconnect
 - ◆ Highest resistance, finest pitch
- M2, M3, M4 (...)
 - ◆ Intra block
- M5, M6, ...
 - ◆ Global wires
 - ◆ Thick coarse metal

□ When scaling introduces too thinner metal

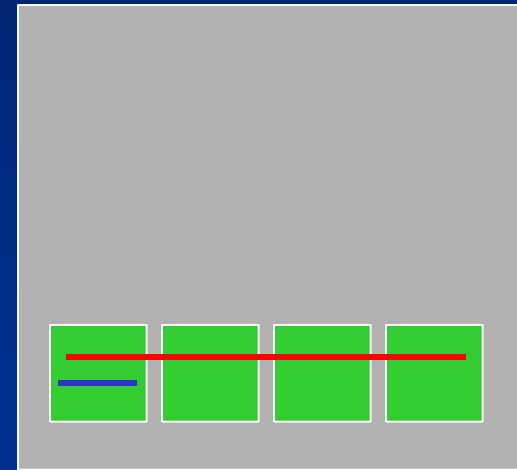
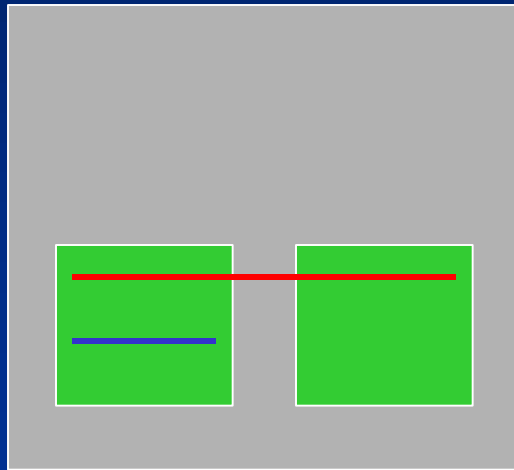
- Create new top layer

Performance of wires?

- Measured relative to logic gates
- Divide speed of circuit by speed of FO4 inverter
- Many options and complex trade-offs
- Local wiring vs Global wiring

Scaling of global wires

- R gets worse with scaling
- C remains constant (+ or -)
- RC grows

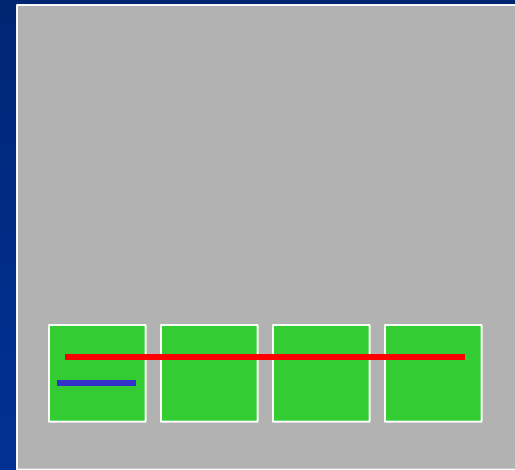
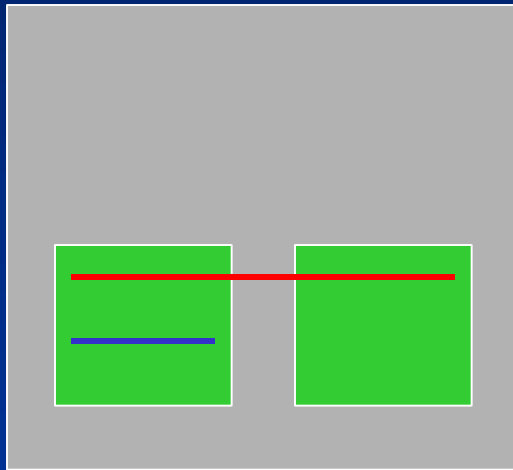


Solution to global wires

- **Use wider or thicker wires**
 - Higher metal levels are faster
- **Use repeaters**
 - Break wire into segments
 - Delay becomes linear with length
- **Repeaters introduce delay and complexity**

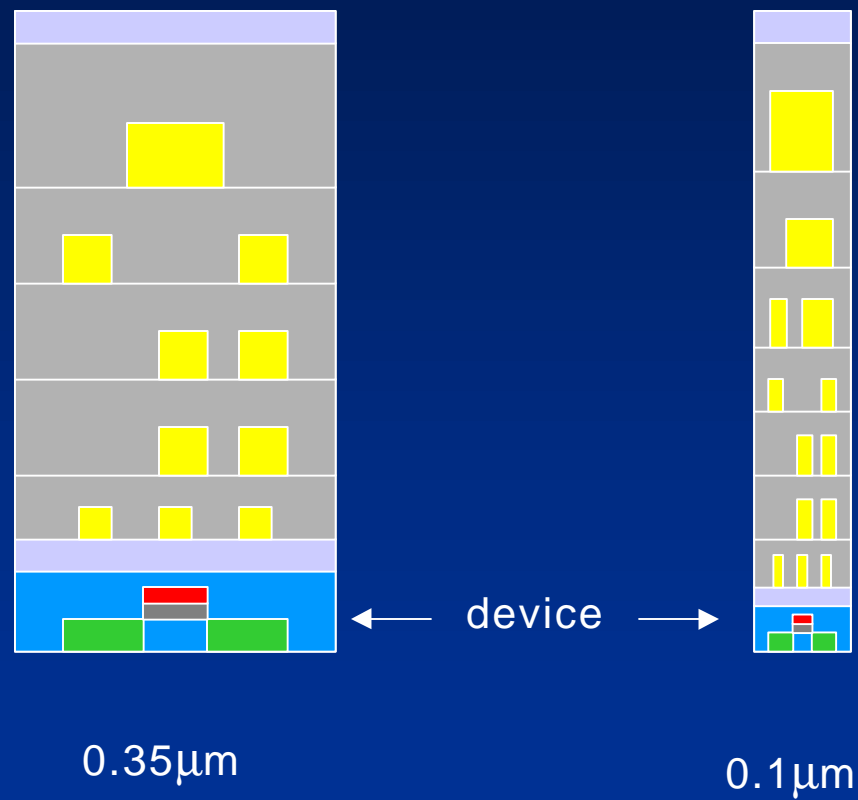
Scaling of local wires

- R is constant (+ or -)
- C falls linearly
- RC is reduced like gate delay (+ or -)



All together

Gates & Wires



Impact of scaling

From the point of view of an architect

Main topics

- **Delay**
- **Power**
- **Area**
- **μArchitecture**
- **Performance**
- **Complexity**

The equation

□ $T = I * CPI * T_c$

- When we scale we improve T_c (reduce it)
- We also add more stages to further reduce T_c
- But more stages increase CPI
 - ◆ load misses, branches, ...
- So we put better predictors and caches and ...
 - ◆ to compensate for the lost CPI due to more stages

□ Future

- some more T_c reduction due to technology
 - ◆ where is the limit?
- no more reduction on number of FO4 per stage?
 - ◆ only new stages for wire delays???? where is the limit?
- so ... we better improve the μ architecture...

The equation (II)

□ Now

- $T_1 = I \times \text{CPI}_1 \times T_{c1}$

□ Usual improvements

- $T_2 = I \times \text{CPI}_1 * \alpha_1 * \beta_1 \times T_{c1} * \gamma * \delta$

$\alpha_1 > 1$ (more stages because of $T_{c1} * \gamma * \delta$)

$\beta_1 < 1$ (better μ architecture)

$\gamma < 1$ (better technology)

$\delta < 1$ (reduced FO4 per stage)

□ Future

- $T_3 = I \times \text{CPI}_1 * \alpha_2 * \beta_2 \times T_{c1} * \gamma$

◆ $\alpha_2 > 1$ (more stages because of $T_{c1} * \delta$)

The equation (III)

- **Another improvement in T_c comes from improvements in circuit design**
 - dynamic vs static logic
 - latches with embedded logic
 - ...

Impact on μ architecture

- **When using global resources we run into wire limitations....**
- **Maintain number of stages**
 - Multilevel (to maintain number of
 - ◆ Caches
 - ◆ Predictors
 - ◆ Register files
 - Reduce structure size
- **Maintain structure size**
 - Increase pipe stages
- **Multibank**
- **Clustering**

More impact on μ architecture

- **As geometries get smaller and stored charge gets smaller...**
 - soft errors (alpha particles, gamma radiation, ...) may happen frequently (we need fault tolerance)

- **Memory will be at 1000 cycles from processor**
 - and caches can not grow as we would like
 - and for some codes caches don't help at all
 - ◆ low temporal locality codes (like data bases)

More architecture stuff

□ More ILP

- bigger Reg file
- bigger IQ
- more functional units
- more ports

□ Longer distance, global wires, slow ...

- amount of state that can be reached in 1 clock

It's a 2D world

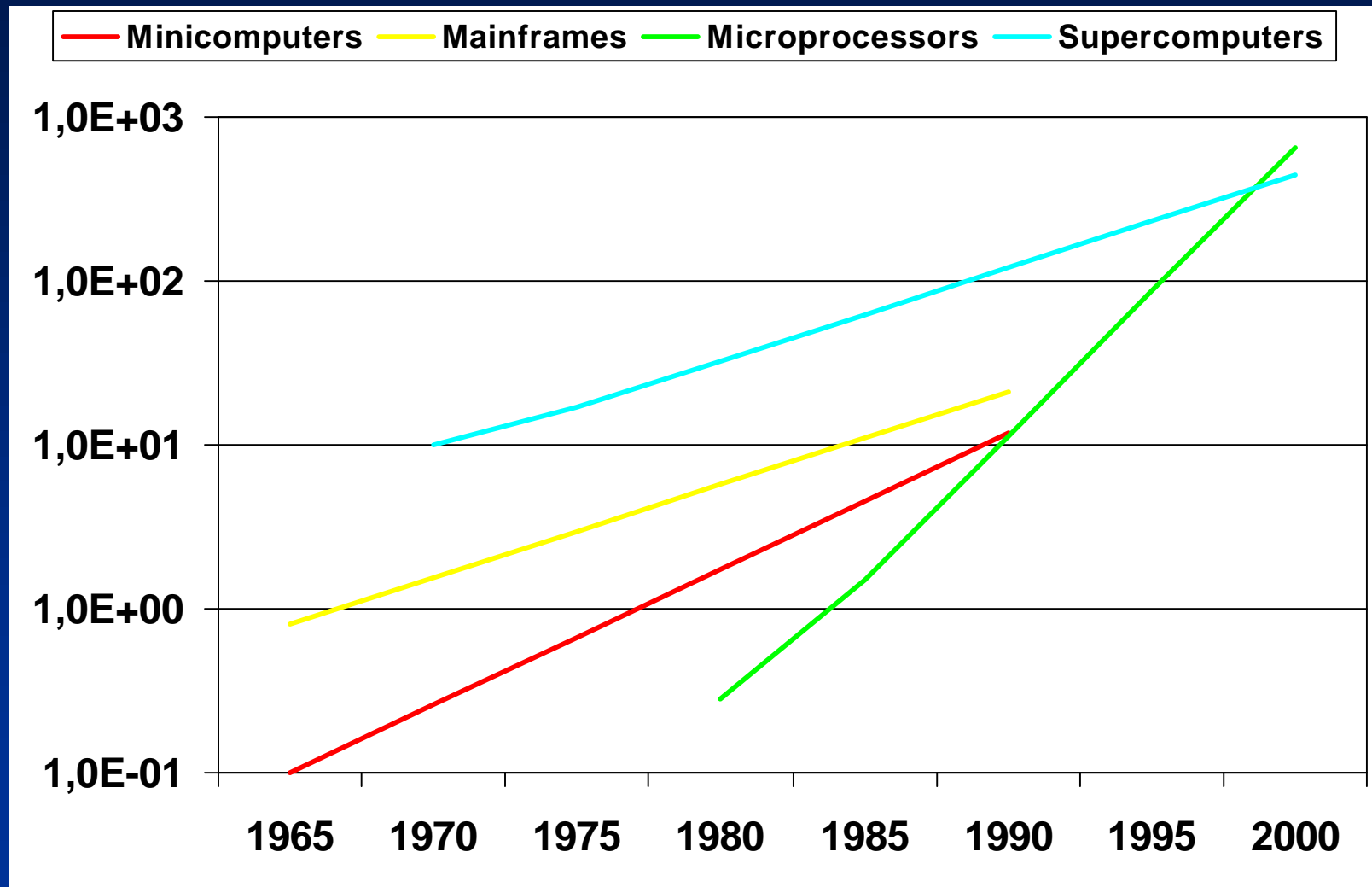
□ Perimeter vs Area

- as time goes by there are **more transistors per pin**
 - ◆ how do you communicate them?
 - ◆ how do you debug them?
- avg. distance between communicating transistors
 - ◆ try to get as many communications as possible local
- similar to living things and 3D (area vs volume)

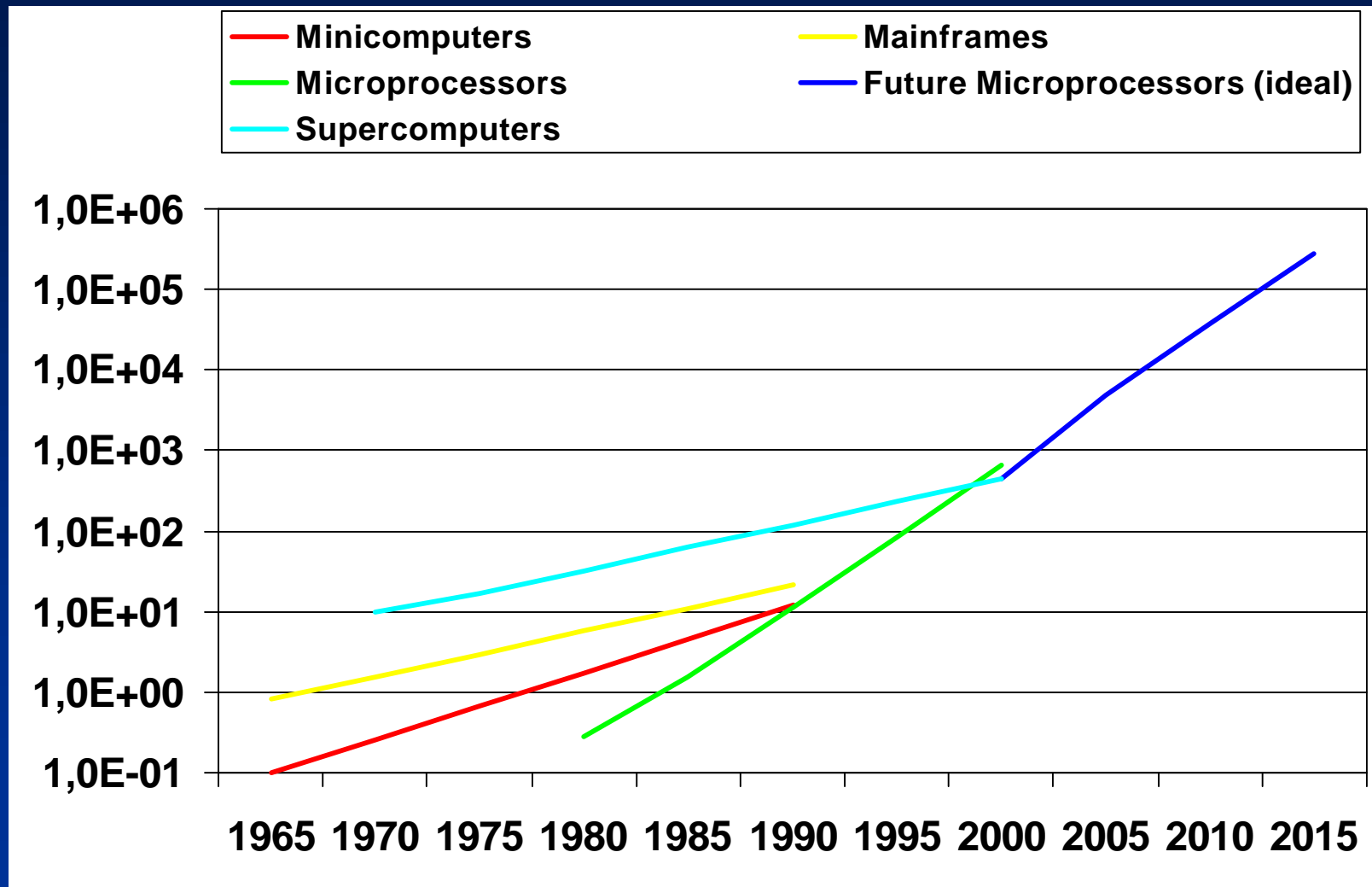
□ At some point bigger is not better

- the extra transistors do not compensate
 - ◆ because of the communication latency

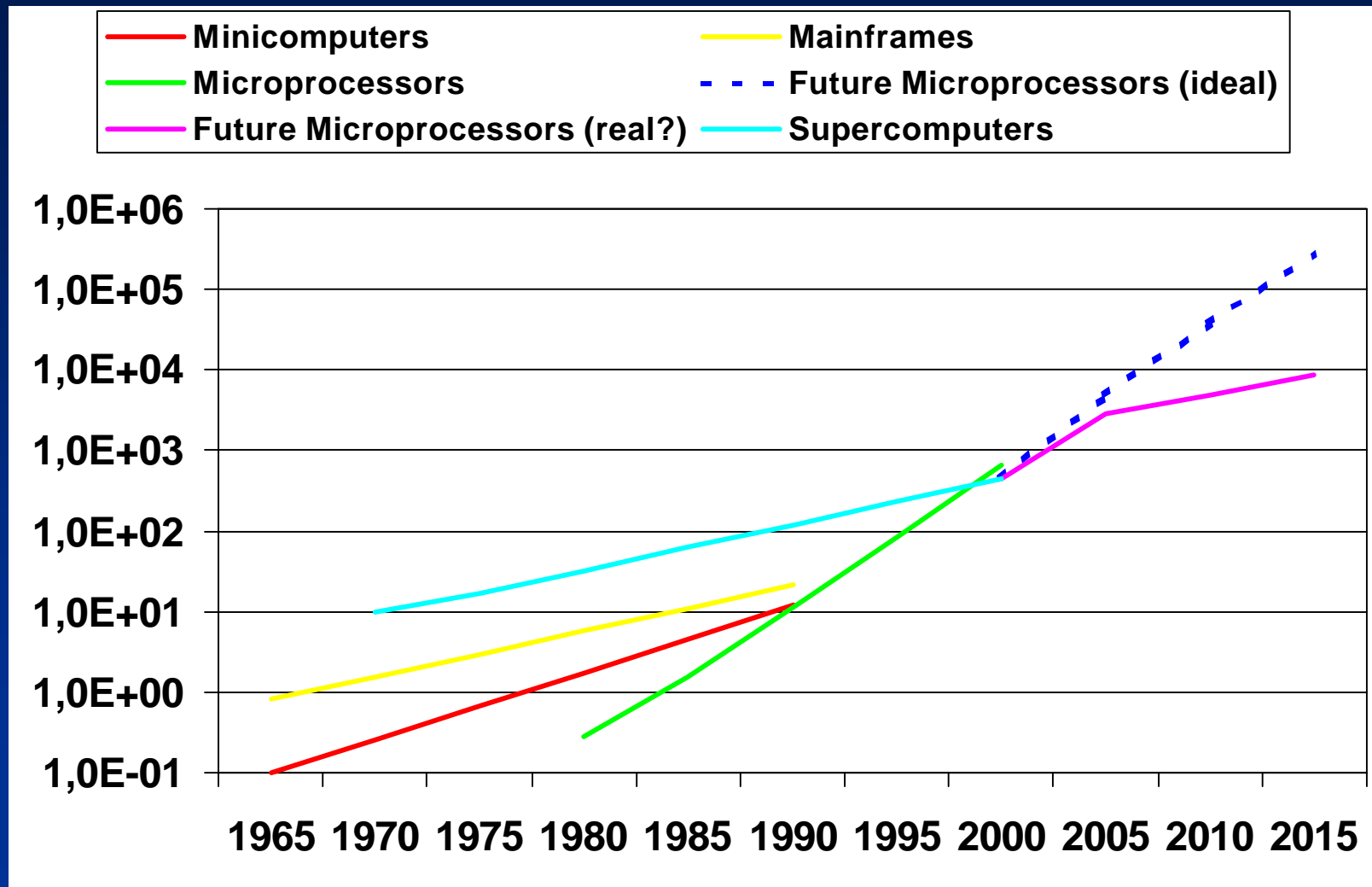
Performance trends



Ideal future



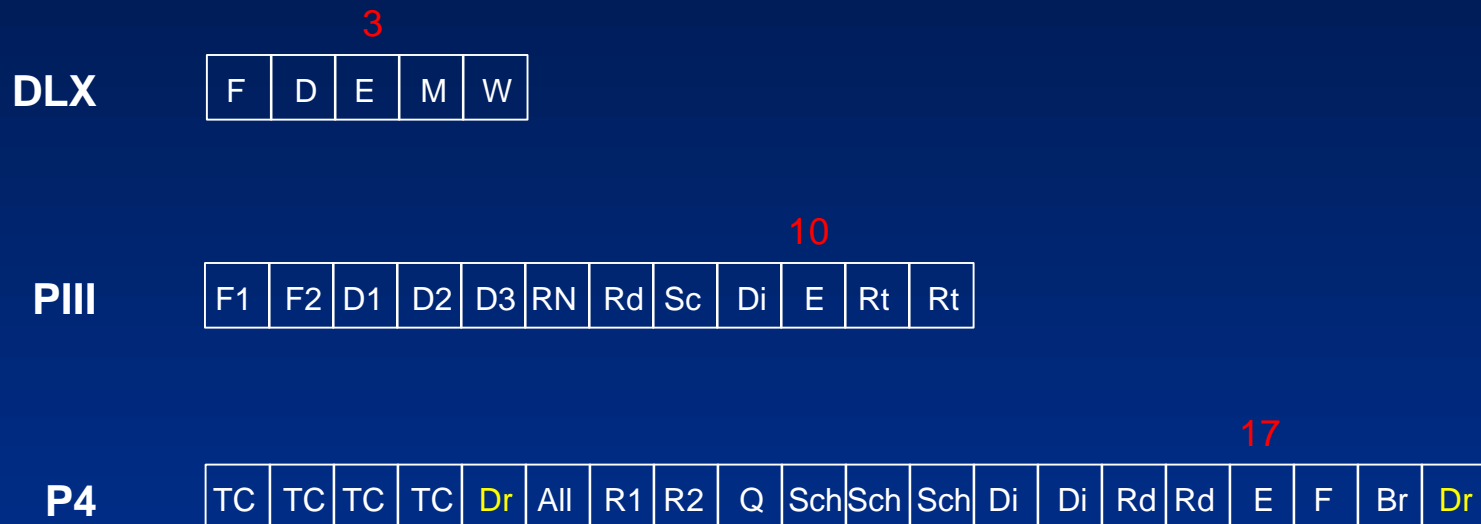
Real future?



Cicle time in FO4

- **Lots of performance from reduced FO4s per stage**
 - Used to be 100 FO4
 - Now close to 16 FO4
 - Difficult below (8 FO4)
 - ◆ Fastest 64 bit add 5.5 FO4
 - ◆ Add latch overhead bypass delay and so on ...
 - SIA aggressive predictions assume 6 FO4

X86 Pipelines



Ppro, PII, PIII share the basic same core

Impact on Complexity

- ❑ **Size of design**
- ❑ **Size of teams**
- ❑ **Worse yield**
- ❑ **Cost is growing exponentially**
 - design+development cycle from 1 to 6+ years
 - time to market -> miss the market
 - fab costs
- ❑ **Intellectual complexity**
 - number of transistors per person?
 - Who understands the whole design?
 - Number (and complexity) of bugs
 - lack of tools

Bibliography

- **Mark Horowitz (ISCA2K panel)**
- **Jim Smith (ISCA2K panel)**
- **Clock Rate vs IPC: The End of the Road for Conventional Microarchitectures, Agarwal et Al., ISCA2k**
- **Circuits, Interconnections and Packagin for VLSI, H.B. Bakoglu, 1990**
- **Computer Technology and Architecture: An Evolving Interaction, John L. Hennessy and Norman P. Jouppi, Sept. 1991, Computer**