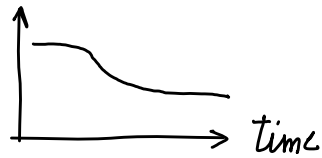


# Cost

cost per unit



Manufacturing costs drop as expertise grows, for that process

- better **methods**
- better **equipment**
- **less waste** (time, materials)

yield



Yield% = (1 - waste%)

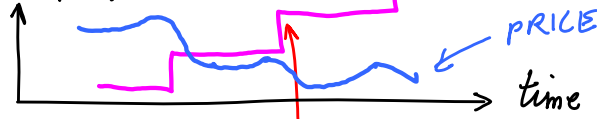
- #(devices **sellable**) **versus** #(devices **produced**)
- #(devices **sellable**) **versus** (cost to produce them)

E.g. DRAM  $\Rightarrow$  price =  $\alpha$  cost

80% contract sales to large equipment makers (hidden)

20% open market

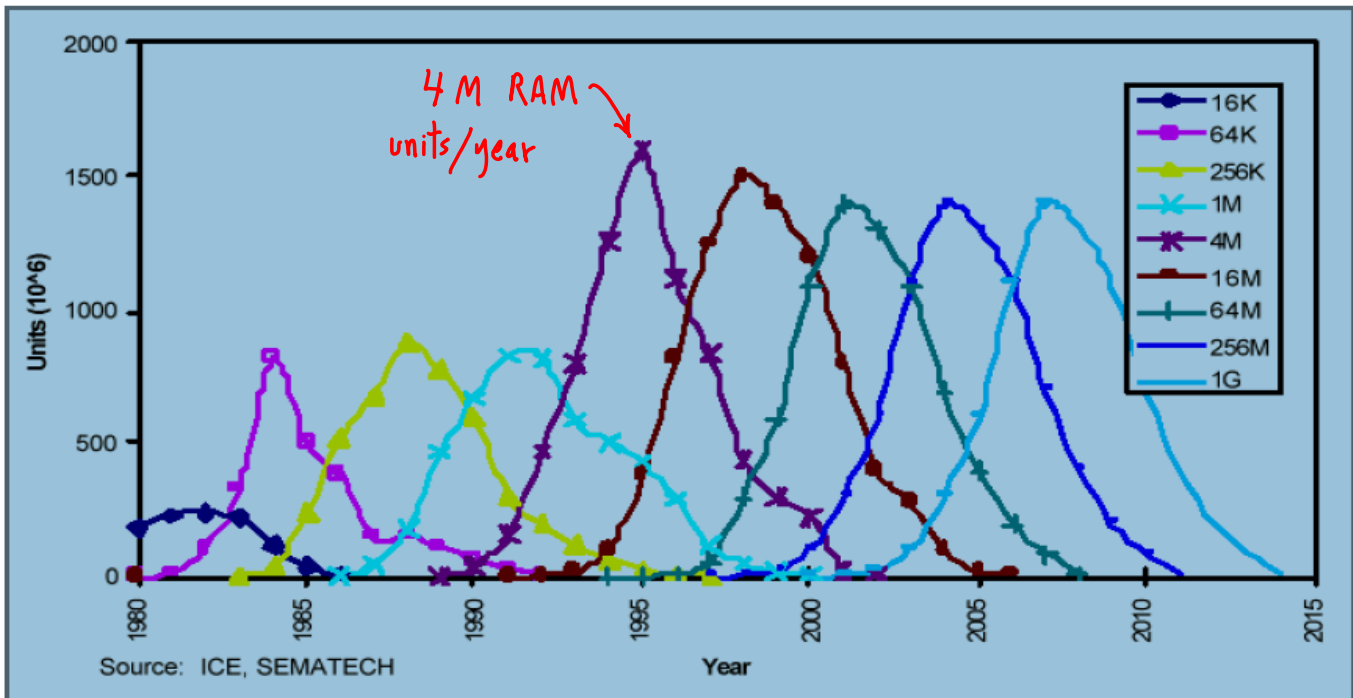
Total Capacity



Commodity market

- many vendors
- same items

new plant online : \$3B / 3 yr



DRAM Unit Volume by Generation [37]

Invest in largest **demand**  $\Rightarrow$  production **cost amortized**  $\Rightarrow$  **larger profit**  
 hot-new  $\Rightarrow$  **high price / low volume**    old-standard  $\Rightarrow$  **low price**

# Costs Drop

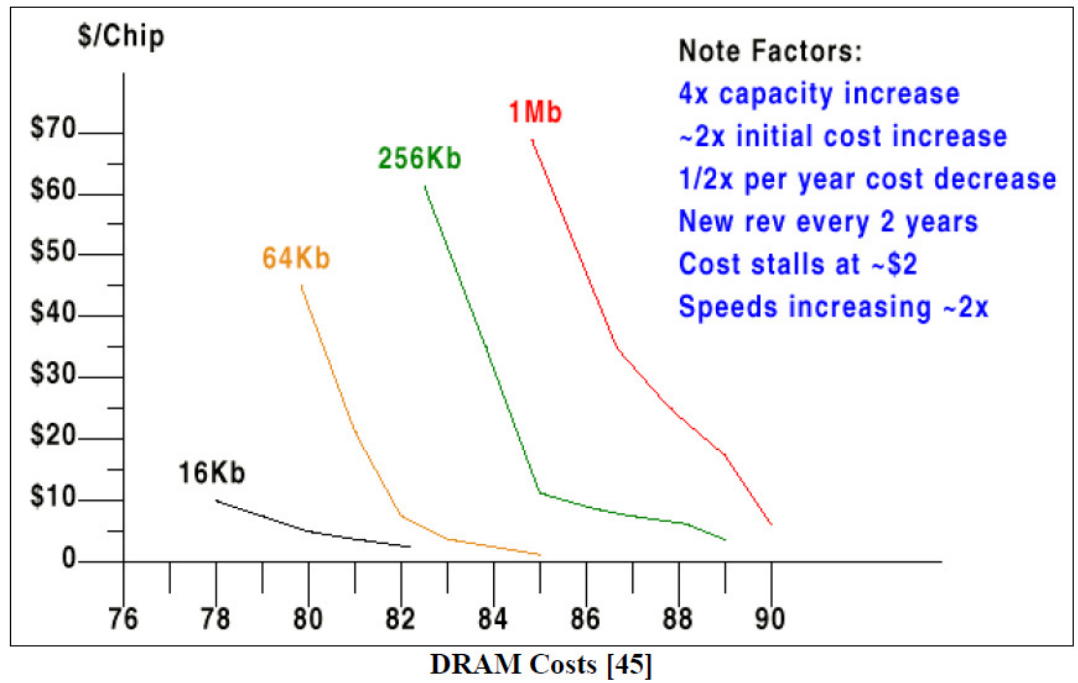
\$50/64

\$65/256

4x capacity  
2x speed

Cost per bit  
per bit/sec

Drop faster



## Changes

Telecom (routers/switches) 20% ↑ ⇒ 50% of market

latency ↓ vs bandwidth ↑  
 routers vs PCs

SRAM

12 ns

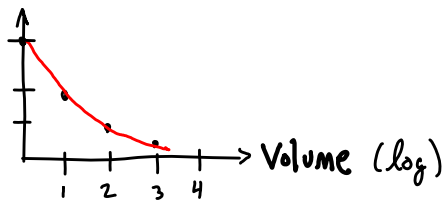
\$50/~2 MB

DRAM

40 ns

\$200/GB

Cost



$\frac{1}{2} \times \text{Cost} : 2 \times \text{Volume}$

Volume → supplier competition → lower cost

⇒ Low-end Market (Price/Performance) ↓

Standardization / Volume ⇒⇒⇒⇒ market acceptance of innovations

# Cloud Pricing

AWS

Combined efficiencies

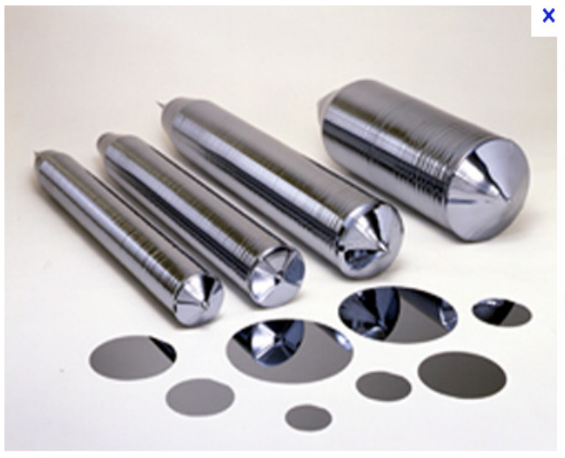
bottomed out

| Description                       | Type        | CU   | Original \$ / CU / Hour | Current \$ / CU / Hour | % Reduction | Aug 2006 | Oct 2007 | May 2008 | Oct 2009 | Feb 2010 | July 2010 | Sep 2010 | Nov 2010 | Nov 2011 |
|-----------------------------------|-------------|------|-------------------------|------------------------|-------------|----------|----------|----------|----------|----------|-----------|----------|----------|----------|
| Small - "the original"            | m1.small    | 1    | \$0.10                  | \$0.085                | 15%         | \$0.10   |          |          | \$0.09   |          |           |          |          |          |
| Large                             | m1.large    | 4    | \$0.10                  | \$0.085                | 15%         |          | \$0.40   |          | \$0.34   |          |           |          |          |          |
| Extra Large                       | m1.xlarge   | 8    | \$0.10                  | \$0.085                | 15%         |          | \$0.80   |          | \$0.68   |          |           |          |          |          |
| High-CPU Medium                   | c1.medium   | 5    | \$0.04                  | \$0.03                 | 15%         |          |          | \$0.20   | \$0.17   |          |           |          |          |          |
| High-CPU Extra Large              | c1.xlarge   | 20   | \$0.04                  | \$0.03                 | 15%         |          |          | \$0.80   | \$0.68   |          |           |          |          |          |
| High-Memory Double Extra Large    | m2.2xlarge  | 13   | \$0.09                  | 0.077                  | 17%         |          |          |          | \$1.20   |          |           | \$1.00   |          |          |
| High-Memory Quad Extra Large      | m2.4xlarge  | 26   | \$0.09                  | 0.077                  | 17%         |          |          |          | \$2.40   |          |           | \$2.00   |          |          |
| High Memory Extra Large           | m2.xlarge   | 6.5  | \$0.12                  | 0.077                  | 33%         |          |          |          | \$0.75   |          |           |          |          |          |
| Cluster Compute                   | cc1.4xlarge | 33.5 | \$0.05                  | \$0.04                 | 19%         |          |          |          |          |          | \$1.60    |          |          |          |
| Cluster Compute Eight Extra Large | cc2.8xlarge | 88   | \$0.03                  | \$0.03                 | 0%          |          |          |          |          |          |           |          |          | \$2.40   |
| Micro                             | t1.micro    | 0.9  | \$0.02                  | \$0.02                 | 0%          |          |          |          |          |          |           | \$0.02   |          |          |
| Cluster GPU Instance              | cg1.4xlarge | 33.5 | \$0.06                  | \$0.06                 | 0%          |          |          |          |          |          |           |          |          | \$2.10   |

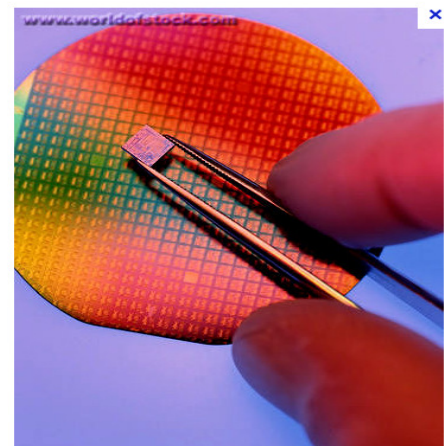
## CPUs, Chips, SOC

Si ingots slicing → wafers

masking, etching, doping

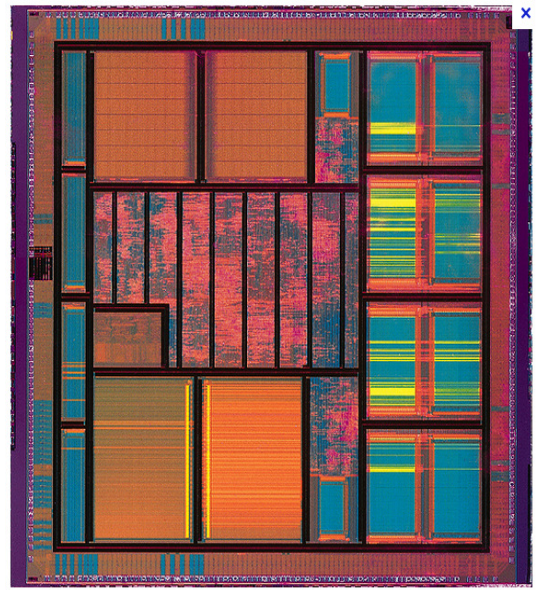
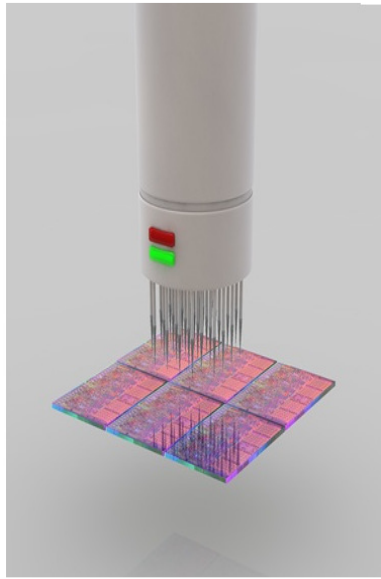


dicing →

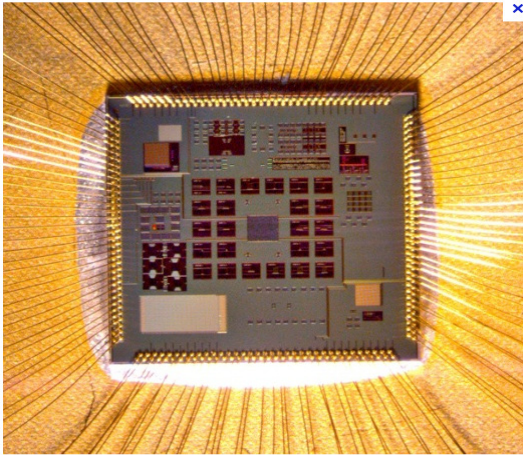




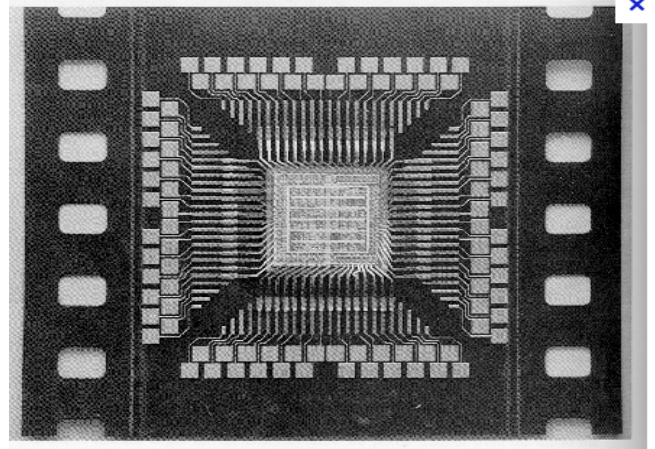
# Circuit testing



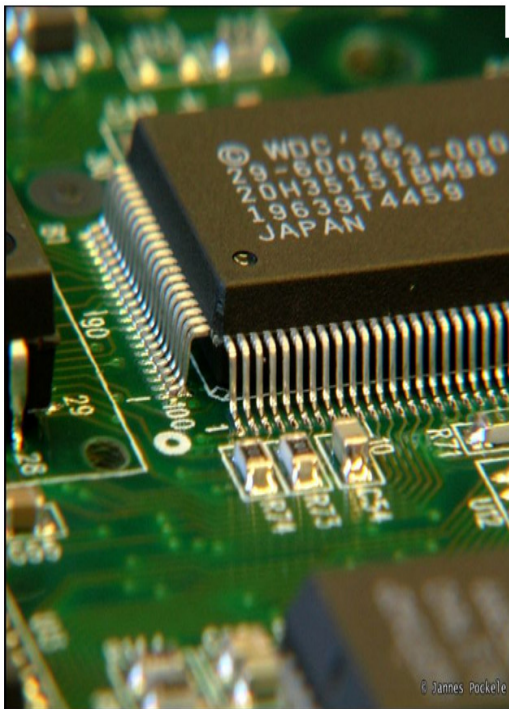
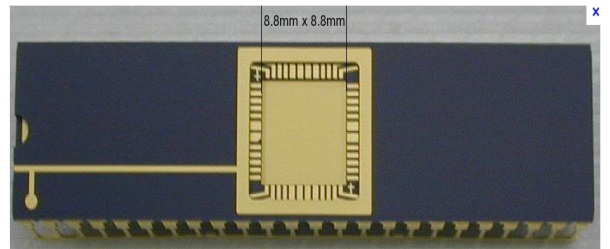
# Pad Bonding



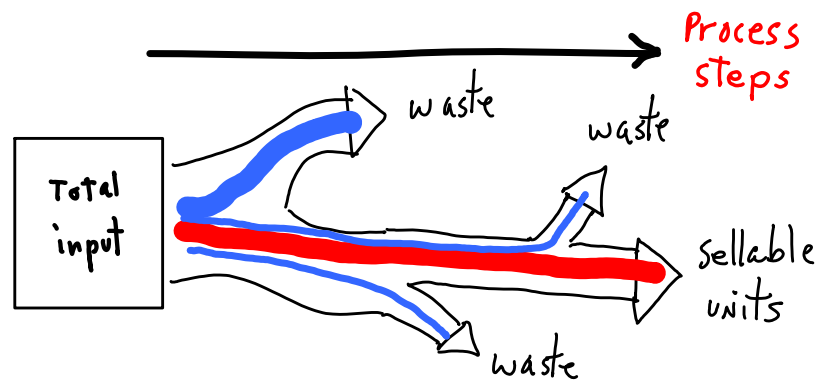
# Pin Packaging



↓ encasing

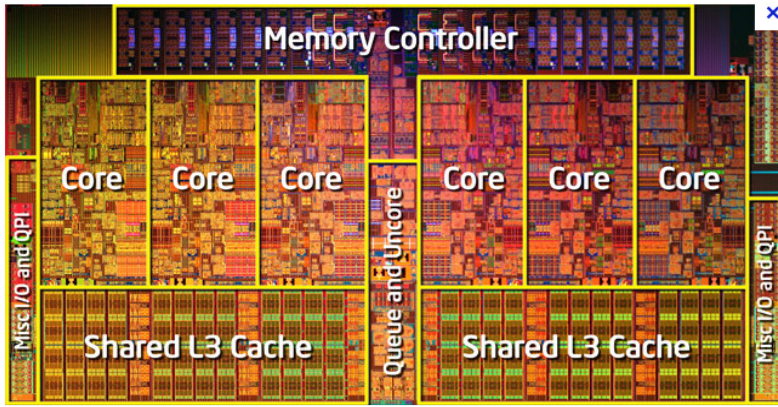


← board printing mounting

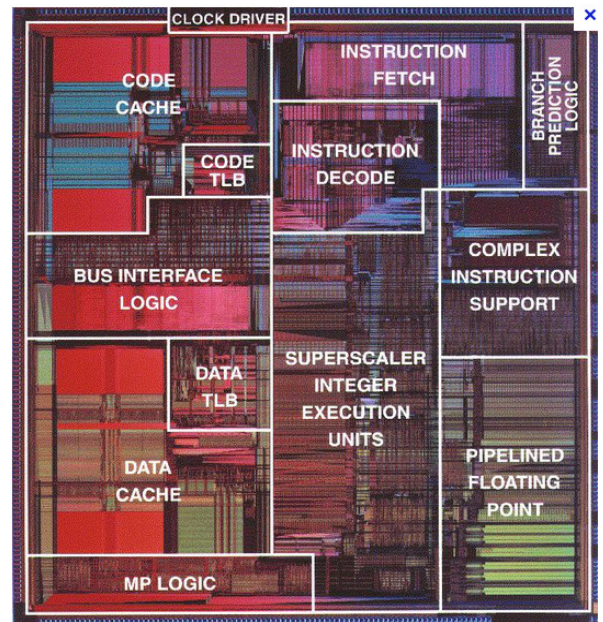




what's inside?



i7



P5

$$\text{Cost} = \frac{C_{\text{die}} + C_{\text{Test}_1} + C_{\text{package}} + C_{\text{Test}_2}}{\#(\text{sellable units})}$$

$$C_{\text{die}} = \frac{C_{\text{wafer}}}{(\# \text{dies})(\text{yield})}$$

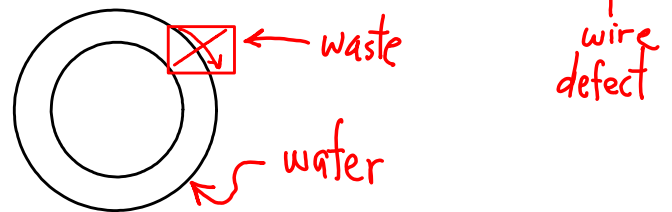
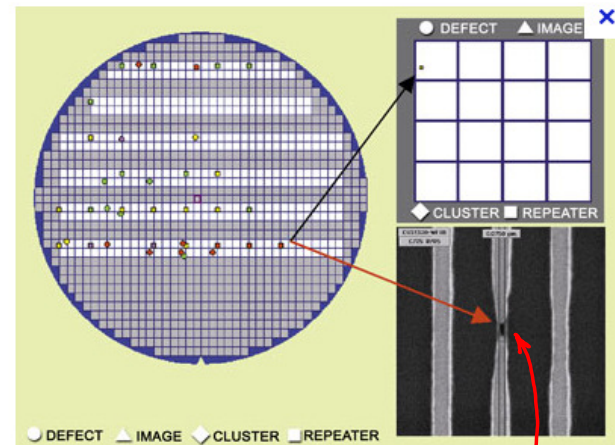
$C_{\text{wafer}} \approx \$5000$

$$\Rightarrow C_{\text{die}} \propto \frac{1}{(\# \text{dies})}$$

$$\begin{aligned} \#(\text{dies}) &= \left( \frac{\text{Area wafer}}{\text{Area die}} \right) - \left( \frac{\text{Circumference wafer}}{\text{Diagonal die}} \right) \\ &= \frac{\pi r^2}{A_{\text{die}}} - \frac{2\pi r}{\sqrt{2} \sqrt{A_{\text{die}}}} \end{aligned}$$

$$\text{yield} \approx \frac{\#(\text{good wafers}) / \#(\text{wafers})}{\left[ 1 + \frac{\#(\text{defects})}{\text{cm}^2} (A_{\text{die}} \text{ cm}^2) \right]^N}$$

Curve fitting for particular process  $\Rightarrow N \in [11.5, 15.5]$



300 mm Wafer

$$\frac{\#(\text{defects})}{\text{cm}^2} \approx 0.04$$

↖ function of time + volume

$$A_{\text{die}} = 2.25 \text{ cm}^2 \Rightarrow 109$$

$$A_{\text{die}} = 1 \text{ cm}^2 \Rightarrow 424$$

P5 Sandy Bridge

2 cm<sup>2</sup>  
\$50

embedded CPU, 32b

0.1 cm<sup>2</sup>  
\$13

printer controller

0.04 cm<sup>2</sup>  
\$0.1

Die size =  $\#(\text{Transistors}) + \#(\text{pins}) \uparrow$   
 +  
 Volume  $\downarrow$   
 +  
 Customization  $\uparrow$

}

cost

Amortized Costs

Mask = \$1M

⇒ reconfigurable gate arrays

Redundancy, e.g.



### Warehouse - Scale Costs

$$\text{Cost}_{\text{computing}} = \text{Cost}_{\text{equipment}} / \text{unit Time} + \text{Cost}_{\text{power}} + \text{Cost}_{\text{structure}} / \text{Time} + \text{Cost}_{\#} / \text{Time} + \text{Cost}_{\text{repair}}$$

(60%)
(40%)

$$\left( \frac{\text{Computers + networks}}{3 \text{ yr}} \right) + \left( \text{other} \right)$$

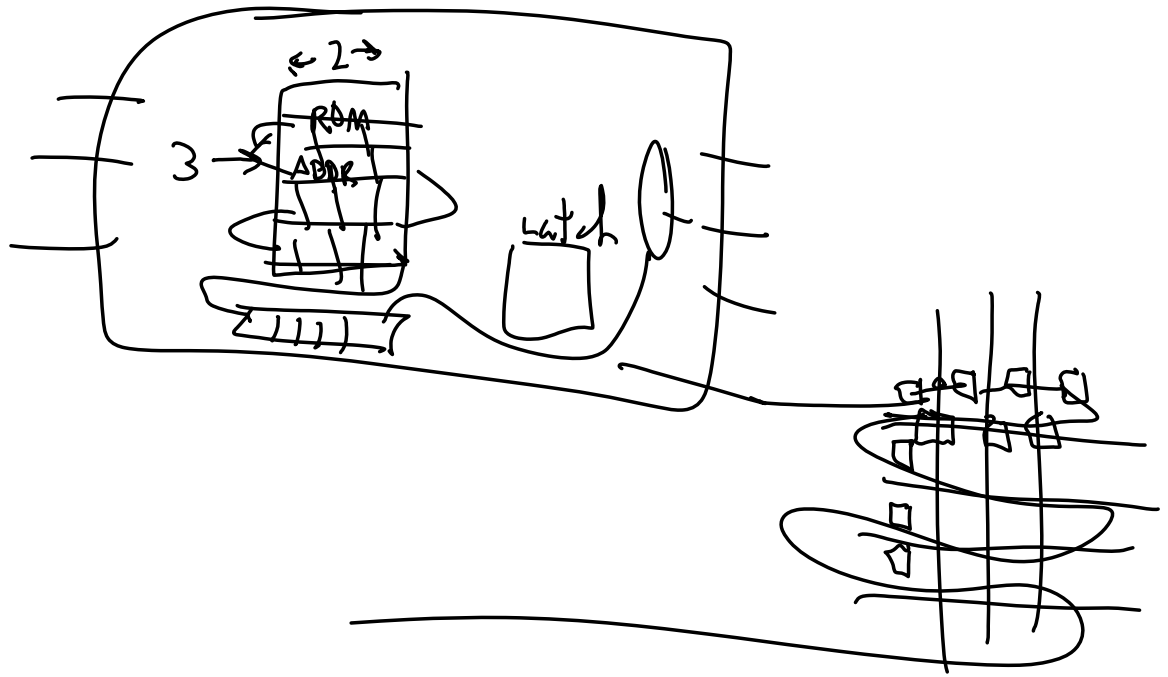
### For our purposes

$$\$/\text{die} = \frac{\$/\text{wafer}}{\#(\text{dies/wafer}) (\% \text{ good dies})}$$

↙ ↘

$$\left( \frac{A_{\text{wafer}}}{A_{\text{die}}} \right) \quad \text{yield} = \frac{1}{\left[ 1 + \left( \frac{\text{defects}}{\text{cm}^2} \right) \frac{A_{\text{die}}}{2} \right]^2}$$

# FPGA



E. G.

$$\$/\text{wafer} = \$1,500 \quad \text{wafer size} = 200 \text{ mm} \Rightarrow A = \pi r^2 \cong 3 \times 10^4 \text{ cm}^2$$

$$\#(\text{defects}/\text{cm}^2) = 0.031 \quad \text{die size} = (1 \text{ cm}) \times (1 \text{ cm}) = 1 \text{ cm}^2$$

$$\text{yield} = \frac{1}{[1 + (0.031) A_{\text{die}}/2]}^2 = \frac{1}{[1 + (0.031) 50]}^2 = \frac{1}{[2.55]}^2 = \frac{1}{2.4}$$

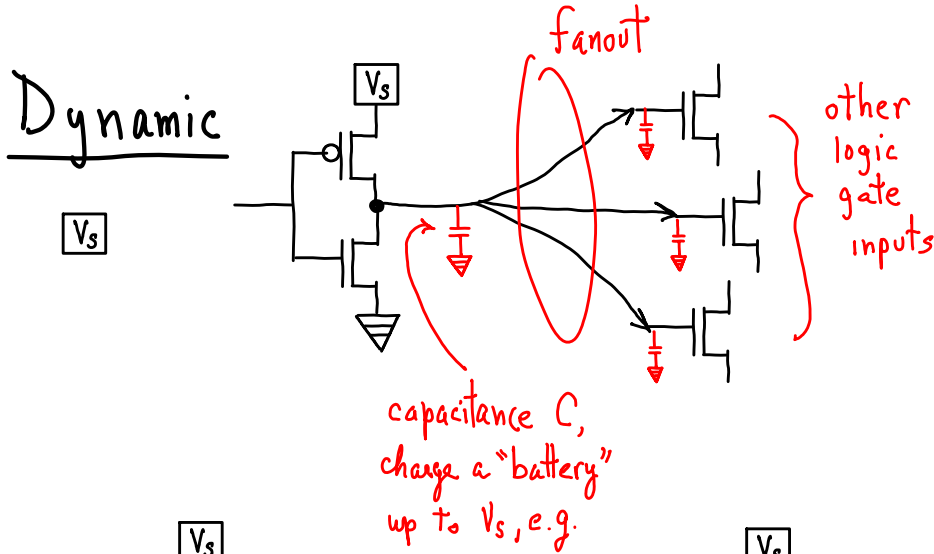
$$\# \text{dies}_{\text{good}} = \left( \frac{A_{\text{wafer}}}{A_{\text{die}}} \right) \left( \frac{1}{2.4} \right) = \left( \frac{3 \times 10^4}{10^2} \right) \left( \frac{1}{2.4} \right) = \frac{300}{2.4} = 125$$

$$\$/\text{die}_{\text{good}} = \frac{\$/\text{wafer}}{\#(\text{dies}_{\text{good}})/\text{wafer}} = \frac{\$1,500}{125} = \$12/\text{die}$$



# CMOS power and energy consumption

- Dynamic:** energy converted to heat due to switching a logic gate's output (0-1 or 1-0).
- Static:** energy converted to heat due to (steady) leakage currents.

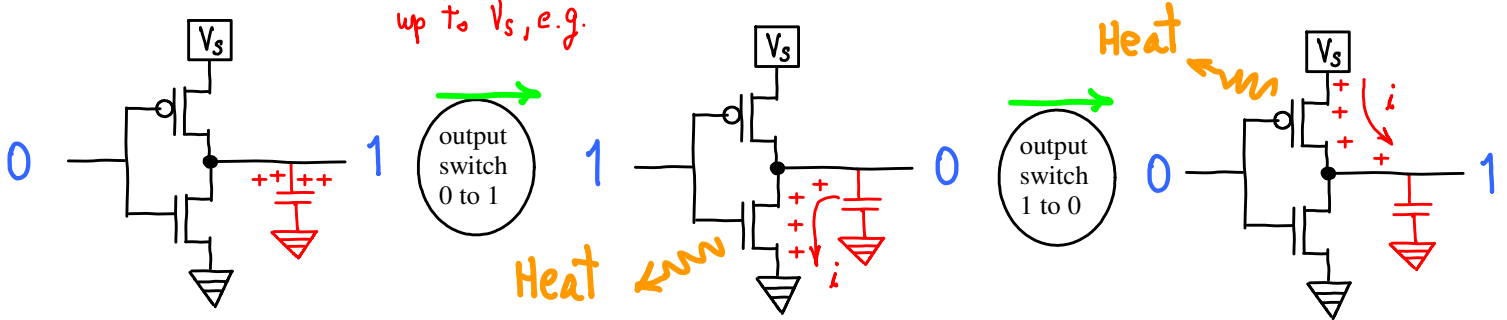


$CR_{max} \propto V$

Speed of charging C

$\longrightarrow E = V/d$

$e^- \longrightarrow$  acceleration  $\sim$  gravity



$$\frac{\text{Joules}}{\text{sec}} = \text{power} = V \left( \frac{+}{\text{sec}} \right) = Vi = (iR)i$$

$R_{\text{Transistor}}$

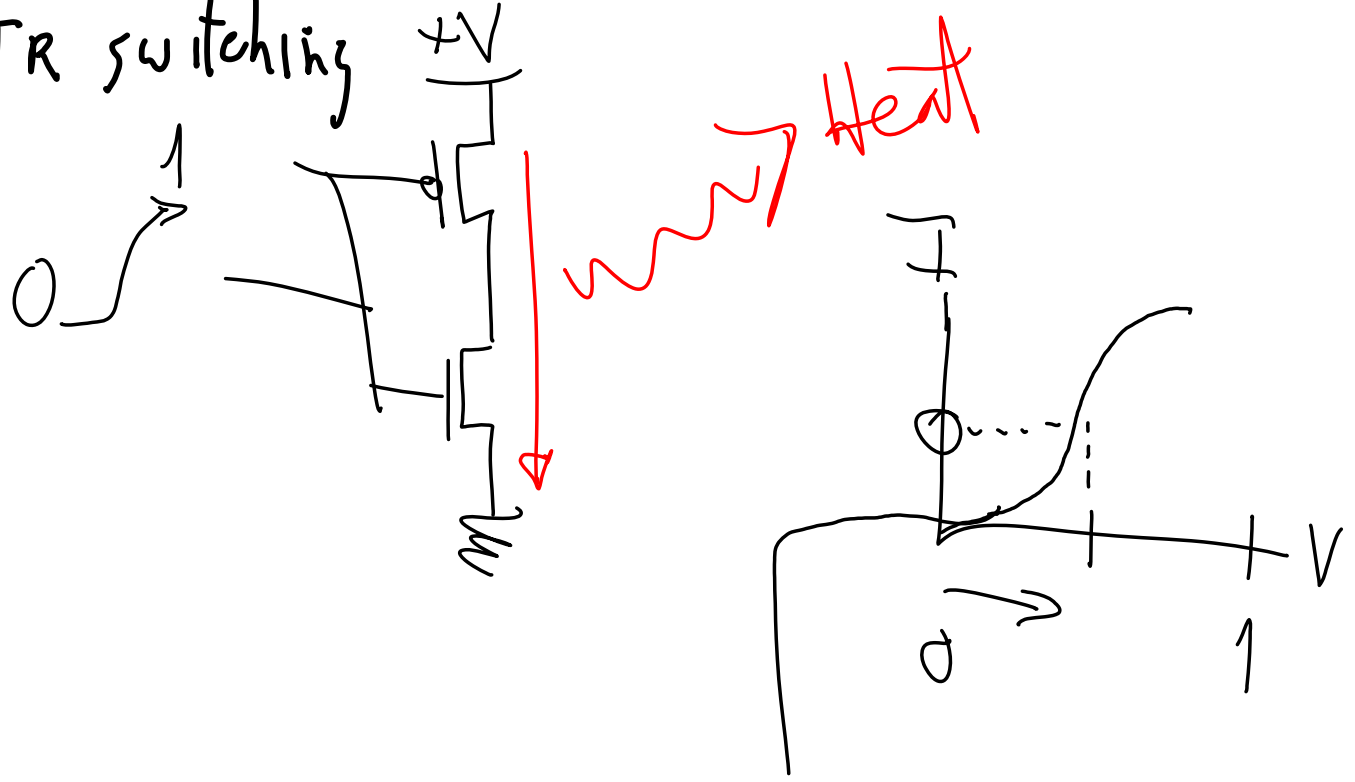
$$E = \frac{\text{Joules}}{\text{Sec}} (\Delta t \text{ sec}) = \frac{1}{2} C V^2$$

$$\Rightarrow \frac{(E / \text{Transistor})}{0-1\text{-switch}} = \frac{1}{2} C_{\text{Transistor}} V^2$$

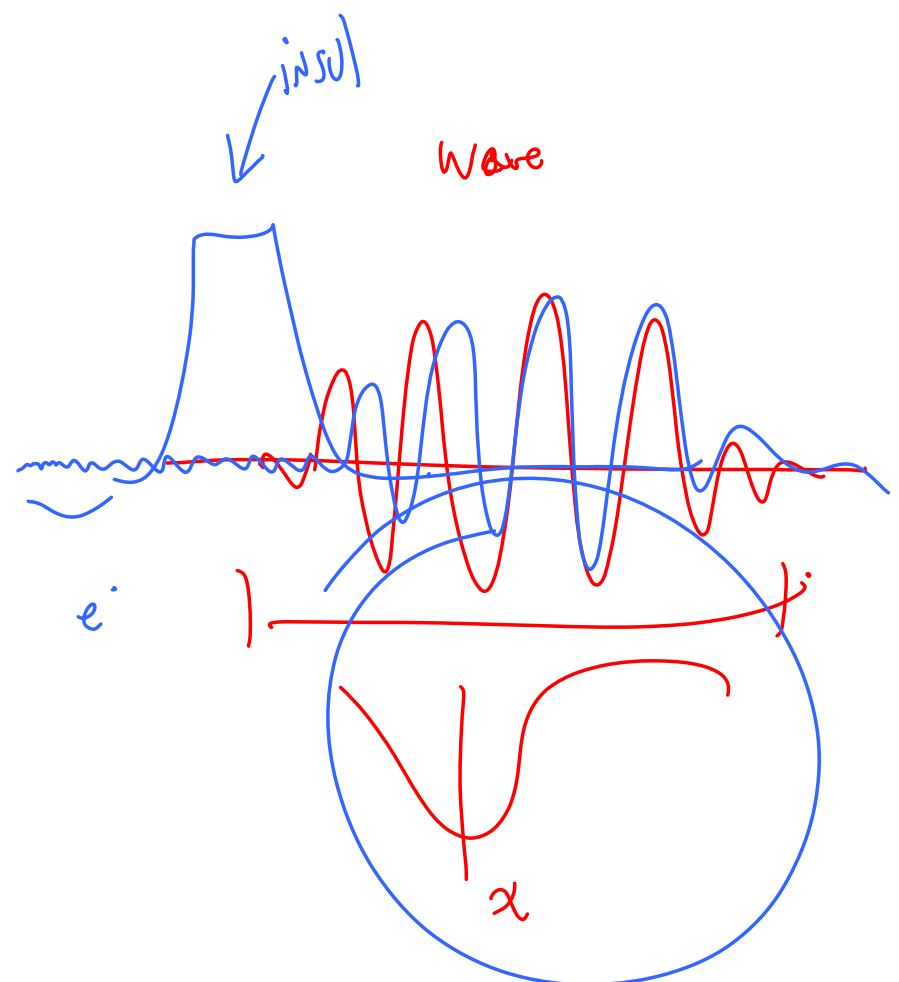
$$\Rightarrow \frac{E_{\text{Total}}}{\text{Switch}} = \sum_i^{* \text{ Transistor}} \frac{1}{2} C_i V^2 = \frac{1}{2} V^2 \sum_i C_i = \frac{1}{2} V^2 C_{\text{Total-chip}}$$

$$\Rightarrow \text{Power} = \left( \frac{E_{\text{Total}}}{\text{Switch}} \right) \left( \frac{\text{switch}}{\text{sec}} \right) = \frac{E_{\text{Total}}}{\text{Switch}} CR$$

TR switching



QM waves



E.G.  $C_{\text{Total}}^{\text{new}} = 0.85 C_{\text{Total}}^{\text{old}}$

$V^{\text{new}} = 0.85 V^{\text{old}}$

$CR \propto V$

$\Rightarrow \frac{CR_{\text{new}}}{CR_{\text{old}}} = \frac{kV_{\text{new}}}{kV_{\text{old}}} = \frac{0.85 V_{\text{old}}}{V_{\text{old}}}$

$\Rightarrow CR_{\text{new}} = 0.85 CR_{\text{old}}$

$\Rightarrow \frac{\text{Power}_{\text{new}}}{\text{Power}_{\text{old}}} = \frac{(\frac{1}{2}) C_{\text{Total}}^{\text{new}} V_{\text{new}}^2 CR_{\text{new}}}{(\frac{1}{2}) C_{\text{Total}}^{\text{old}} V_{\text{old}}^2 CR_{\text{old}}}$

$= \frac{(0.85 C_{\text{Total}}^{\text{old}}) (0.85 V_{\text{old}})^2 (0.85 CR_{\text{old}})}{C_{\text{Total}}^{\text{old}} V_{\text{old}}^2 CR_{\text{old}}}$

$= (0.85)^4 = 52\%$

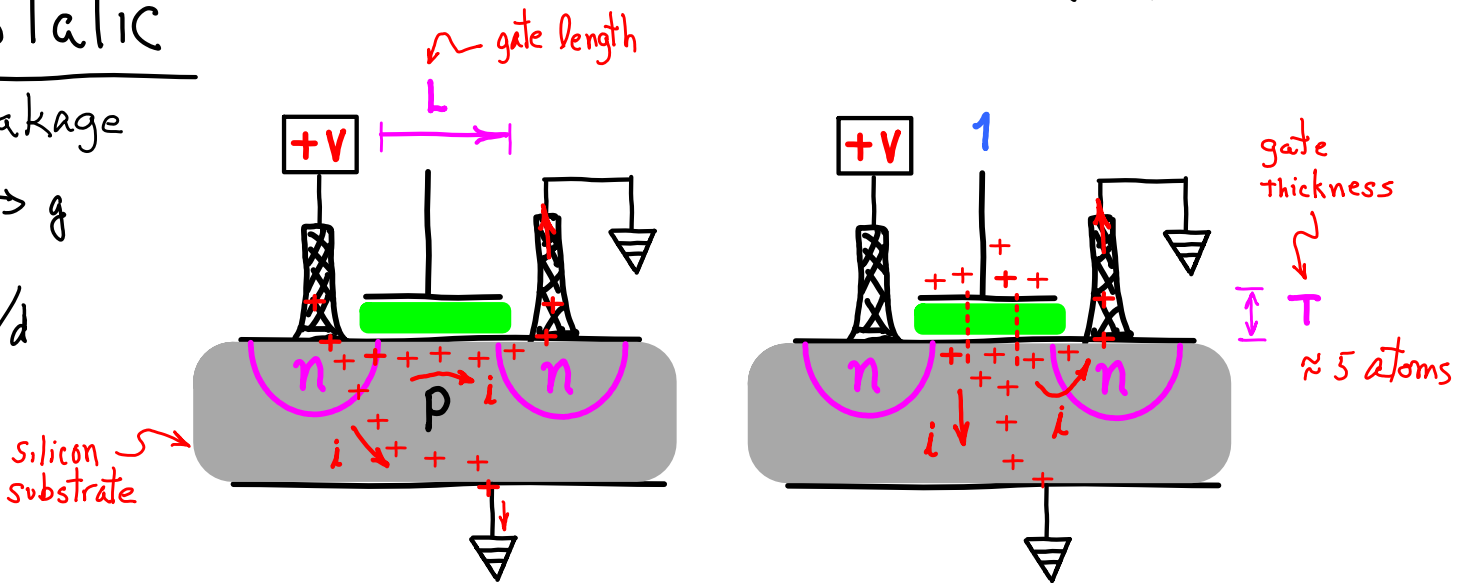
$\frac{E_{\text{new}}}{E_{\text{old}}} = \frac{k_{\text{switches}} E_{\text{switch}}^{\text{new}}}{k_{\text{switches}} E_{\text{switch}}^{\text{old}}} = \frac{k \frac{1}{2} C^{\text{new}} V_{\text{new}}^2}{k \frac{1}{2} C^{\text{old}} V_{\text{old}}^2} = \frac{(0.85 C^{\text{old}}) (0.85 V_{\text{old}})^2}{C^{\text{old}} V_{\text{old}}^2}$

$= (0.85)^3 = 61\%$

## Static

leakage

$E \rightarrow g$   
 $= V/d$



$i/\text{Area} \uparrow$  as  $L, T \downarrow \Rightarrow$

$\text{Power}_{\text{leak}} = i_{\text{leak}} V$  now greater than dynamic power

STATIC POWER  
Total Power

year

| 1985 | 1990 | 1995 | 2000 | 2005 | 2010 |
|------|------|------|------|------|------|
| 1%   | 5%   | 7%   | 20%  | 30%  | 60%  |

Geometric Mean

what's the average rate?

$r_1$   $r_2$   $r_3$   $r_4$   $r_5$   
↑ rate of change

What rate would give the same overall change over the entire time period?

⇒ Find  $\bar{r}$  s.t.  $(r_1 \cdot r_2 \cdot r_3 \cdot r_4 \cdot r_5) = (\bar{r} \bar{r} \bar{r} \bar{r} \bar{r})$

$r_1 = \frac{5}{1}$     $r_2 = \frac{7}{5}$     $r_3 = \frac{20}{7}$     $r_4 = \frac{30}{20}$     $r_5 = \frac{60}{30}$

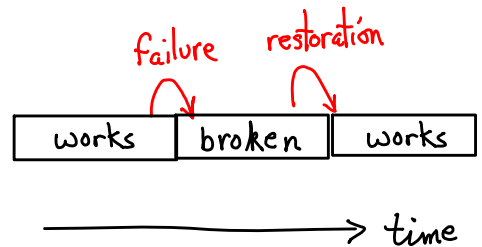
$\bar{r} = (\prod_i^n r_i)^{\frac{1}{n}} = (60/1)^{\frac{1}{5}} \cong 2.3$  per 5 years

Prediction for 2015?  $\bar{r}(60\%) = 2.3(60\%) = 1.38\%$

static power exceeds dynamic by 38%! 😡

**Failures, systemic, hierarchic**

Service Accomplishment = it works  
Service Interruption = broken



Reliability = time to failure

Mean Time To Failure (MTTF) = Avg(Reliability) = Avg(time working)

Mean Time To Repair (MTTR) = Avg(time to fix)



Failure Rate =  $1/MTTF$  (expected #failures in a year, e.g.)

Mean Time Between Failures (MTBF) =  $MTTF + MTTR$

Availability = % time working =  $\frac{MTTF}{MTBF}$

Multi-level  
Recovery

- error correction
- redo
- redundancy

Disk sub-system MTTF (Mhr)

|                  |     |
|------------------|-----|
| 10 disks         | 1   |
| 1 ATA controller | 1/2 |
| 1 power supply   | 1/5 |
| 1 fan            | 1/5 |
| 1 ATA cable      | 1   |

ASSUME: exponentially dist. failures

- independent of  $t$ , time of experiment
- No dependency between components

$$E(\# \text{ failures in } \Delta T) = \mu \Delta T$$

1 failure  $\rightarrow 1 = \mu \overline{\Delta T}_1$  ↑ avg failure rate  
↑ avg time to 1<sup>st</sup> fail

$$MTTF = \overline{\Delta T}_1 = 1/\mu$$

$$\mu = \sum_i \mu_i = \sum_i (1/MTTF_i)$$

Assume devices fail independently:  
 $P(F_1 \text{ or } F_2) = P(F_1) + P(F_2)$

$$= (10 + 2 + 5 + 5 + 1) / \text{Mhr} = 23 / \text{Mhr}$$

$$MTTF = 1/\mu = \frac{10^6 \text{ hr}}{23}$$

$$\approx \left(\frac{10^2}{25}\right) 10^4 = 40,000 \text{ hr}$$

$$\approx 4 \text{ yr.}$$

$$\text{hrs/yr} = (365 \text{ d})(24 \text{ hr/d})$$

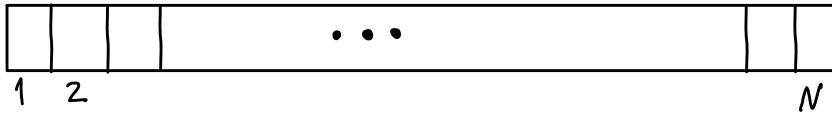
$$\approx (400)(100/4)$$

$$= 10,000$$

E.g., 2 power supplies?

Let's add a redundant power supply: if one fails we fix it while the other keeps working. How long before the system fails:

it happens that before we finish the repair, the backup also breaks?



## A Model

Suppose  $N$  time periods, each of unit length.  $N$  is large so that no devices ever survive to age  $N+1$ .

Each time period, flip a coin:

$$\begin{aligned} \text{Prob}(\text{works}) &= p \\ \text{Prob}(\text{fails}) &= (1-p) \end{aligned}$$

$$\text{Prob}(\text{fails at } i) = p^i (1-p) \quad \leftarrow \text{fails at } i+1$$

$$E(i) = \sum_{k=1}^N k \text{Prob}(\text{fails at } k)$$

$$= \sum_{k=1}^N k p^k (1-p) = (1-p) \sum_{k=1}^N k p^k \quad p < 1$$

$$= (1-p) \sum_{k=1}^N \frac{d}{dp} p^{k+1} = (1-p) p \frac{d}{dp} \sum_{k=1}^N p^k$$

$$\cong (1-p) p \frac{d}{dp} 1/(1-p) = (1-p) p / (1-p)^2$$

$$= p/(1-p) = \text{MTTF}$$

## Example

$$\text{Suppose } \text{Prob}(\text{fail}) = (1-p) = 1/2^m$$

$$\text{MTTF} = p/(1-p) = (1 - 1/2^m) / (1/2^m) = 2^m - 1$$

## 2 Power Supplies

$$\text{Prob}(\text{both work}) = p^2 \quad \implies \quad \text{Prob}(\text{not both working}) = (1-p^2)$$

$$\text{MTTF}_2 = \frac{p^2}{(1-p^2)} = \frac{p^2}{(1-p)(1+p)} = \frac{p}{(1-p)} \frac{p}{(1+p)}$$

$$= \text{MTTF} \frac{p}{(1+p)} \quad \left. \begin{array}{l} \text{for } p \text{ large} \\ (p \cong 1) \end{array} \right\} \cong \frac{1}{2} \text{MTTF}$$

What's the probability the 2<sup>nd</sup> fails in MTTR?

Time independent = begin as if only one power supply.

Looked at another way

Total time =  $T$

$N$  intervals,  $dt = T/N$

Prob(fail in  $dt$ ) =  $1/N$   
=  $1-p$

Prob(fails in interval of length  $L = n dt$ )

$$\sum_i^n \text{Prob(fail in } dt_i) = \sum_i^n (1/N) = n/N$$

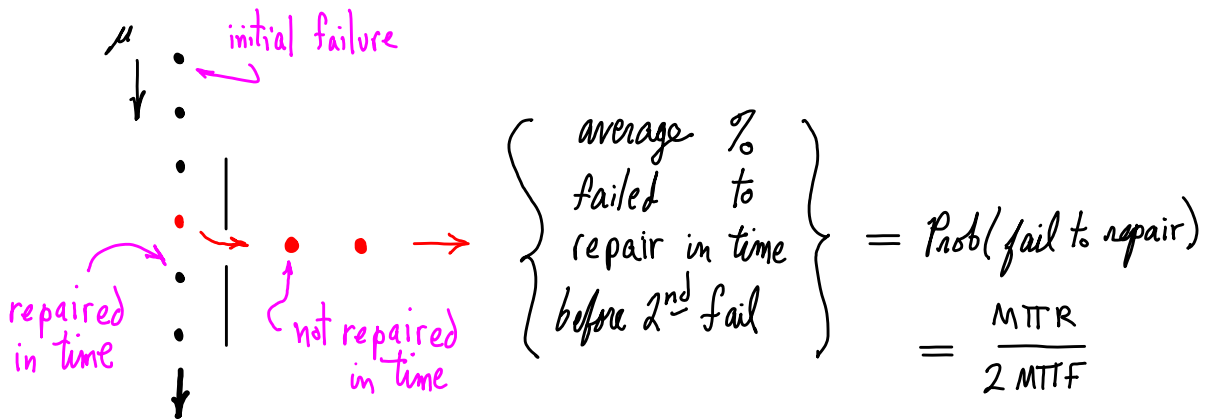
$$= \frac{L/dt}{T/dt} = \frac{L}{T} \implies \text{uniform dist.}$$

Prob(fail in  $T$ ) = 1

$MTTF = T/2$

$$\implies \text{Prob(fail in MTTR)} = \frac{MTTR}{T} = \frac{MTTR}{2 MTTF}$$

$1/(1/2 MTTF) = \mu$ , average rate of failures for 2 power supplies



rate of system failures =  $\mu$  (% not repaired)

$$= 1/(1/2 MTTF) \left( \frac{MTTR}{2 MTTF} \right) = \frac{MTTR}{MTTF^2}$$

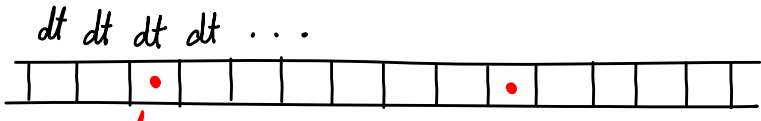
$$\implies \frac{MTTR}{MTTF^2} = \frac{24}{(1/5 M)^2} \approx \frac{25^2}{10^{12}} = \frac{(100/4)^2}{10^{12}} = \frac{10^4}{16 \cdot 10^{12}} = \frac{1}{16 \cdot 10^8}$$

$$= \frac{1 \text{ system failure}}{16 \cdot 10^4 \text{ machine yrs}}$$

$\implies$  on average, w/  $10^4$  systems, 1 fail in 16 yrs.

# Poisson

Bernoulli w/  $p = \lambda dt$



$$\text{Prob}(m, N, p) = \binom{N}{m} p^m q^{N-m}$$

approx.  $\Rightarrow \frac{e^{-pN} (pN)^m}{m!}$

Probability that a failure happens in  $[x, x+dt]$  } flip coin w/ Prob  $p$   
 inter-arrival time

let  $pN \xrightarrow[p \rightarrow 0, N \rightarrow \infty]{} \lambda$

Prob

$$\text{Prob}(m, \lambda) = \frac{e^{-\lambda t} (\lambda t)^m}{m!}, \quad \lambda = \text{rate} \quad \text{Prob}(k \text{ events in time } t)$$

$$\text{Prob}(X_n \leq x) = 1 - e^{-\lambda x}$$

$$\text{Prob}(\text{time of next arrival} \leq x)$$

$$E(X) \equiv 1/\lambda \quad \text{def'n of } \lambda$$

$$\text{Prob}(\text{arrival time } t \leq s \mid \text{one event did occur}) = s/t \quad \text{uniform dist.}$$



# Price Performance

Basic measure:  $\text{operations} / \$$

| Component       | System 1       |               | System 2       |               | System 3       |                 |
|-----------------|----------------|---------------|----------------|---------------|----------------|-----------------|
|                 | Component      | Cost (% Cost) | Component      | Cost (% Cost) | Component      | Cost (% Cost)   |
| Base server     | PowerEdge R710 | \$653 (7%)    | PowerEdge R815 | \$1437 (15%)  | PowerEdge R815 | \$1437 (11%)    |
| Power supply    | 570 W          |               | 1100 W         |               | 1100 W         |                 |
| Processor       | Xeon X5670     | \$3738 (40%)  | Opteron 6174   | \$2679 (29%)  | Opteron 6174   | \$5358 (42%)    |
| Clock rate      | 2.93 GHz       |               | 2.20 GHz       |               | 2.20 GHz       |                 |
| Total cores     | 12             |               | 24             |               | 48             |                 |
| Sockets         | 2              |               | 2              |               | 4              |                 |
| Cores/socket    | 6              |               | 12             |               | 12             |                 |
| DRAM            | 12 GB          | \$484 (5%)    | 16 GB          | \$693 (7%)    | 32 GB          | \$1386 (11%)    |
| Ethernet Inter. | Dual 1-Gbit    | \$199 (2%)    | Dual 1-Gbit    | \$199 (2%)    | Dual 1-Gbit    | \$199 (2%)      |
| Disk            | 50 GB SSD      | \$1279 (14%)  | 50 GB SSD      | \$1279 (14%)  | 50 GB SSD      | \$1279 (10%)    |
| Windows OS      |                | \$2999 (32%)  |                | \$2999 (33%)  |                | \$2999 (24%)    |
| Total           |                | \$9352 (100%) |                | \$9286 (100%) |                | \$12,658 (100%) |
| Max ssj_ops     | 910,978        |               | 926,676        |               | 1,840,450      |                 |
| Max ssj_ops/\$  | 97             |               | 100            |               | 145            |                 |

least power

server-side java ops/sec

benchmark  
ops/\$

Figure 1.18 Three Dell PowerEdge servers being measured and their prices as of August 2010. We calculated the cost of the processors by subtracting the cost of a second processor. Similarly, we calculated the overall cost of memory by seeing what the cost of extra memory was. Hence, the base cost of the server is adjusted by removing the estimated cost of the default processor and memory. Chapter 5 describes how these multi-socket systems are connected together.

## ops/watt vs load

avg utilization  
≈ 10-50%

What's the message?

cpu as % performance?  
% cost?

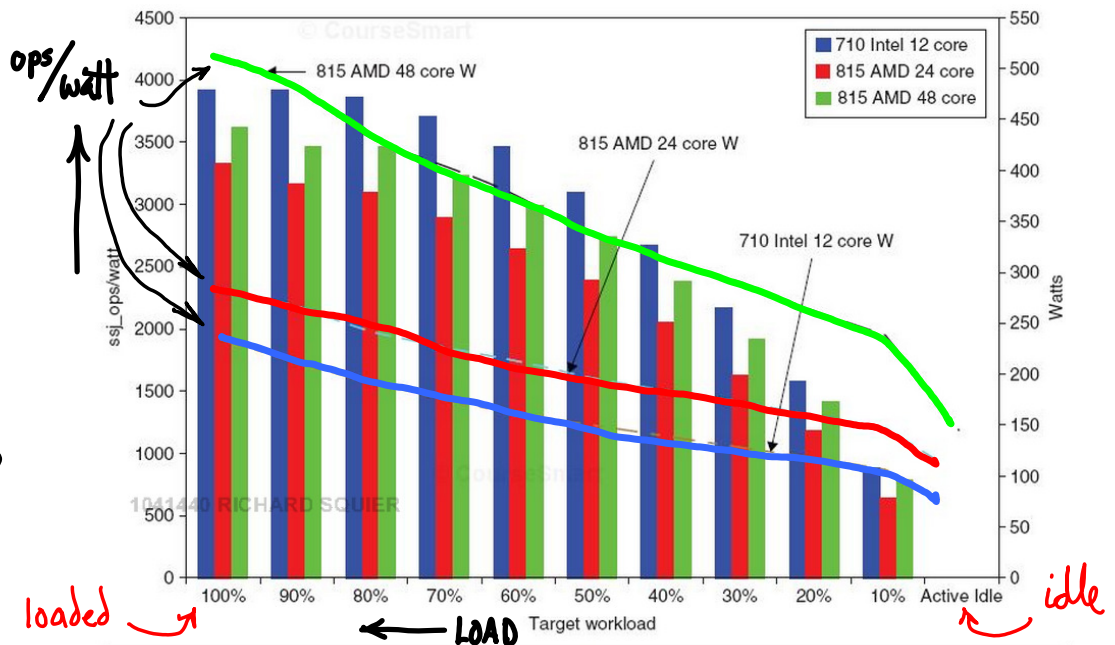


Figure 1.19 Power-performance of the three servers in Figure 1.18. Ssj\_ops/watt values are on the left axis, with the three columns associated with it, and watts are on the right axis, with the three lines associated with it. The horizontal axis shows the target workload, as it varies from 100% to Active Idle. The Intel-based R715 has the best ssj\_ops/watt at each workload level, and it also consumes the lowest power at each level.

