

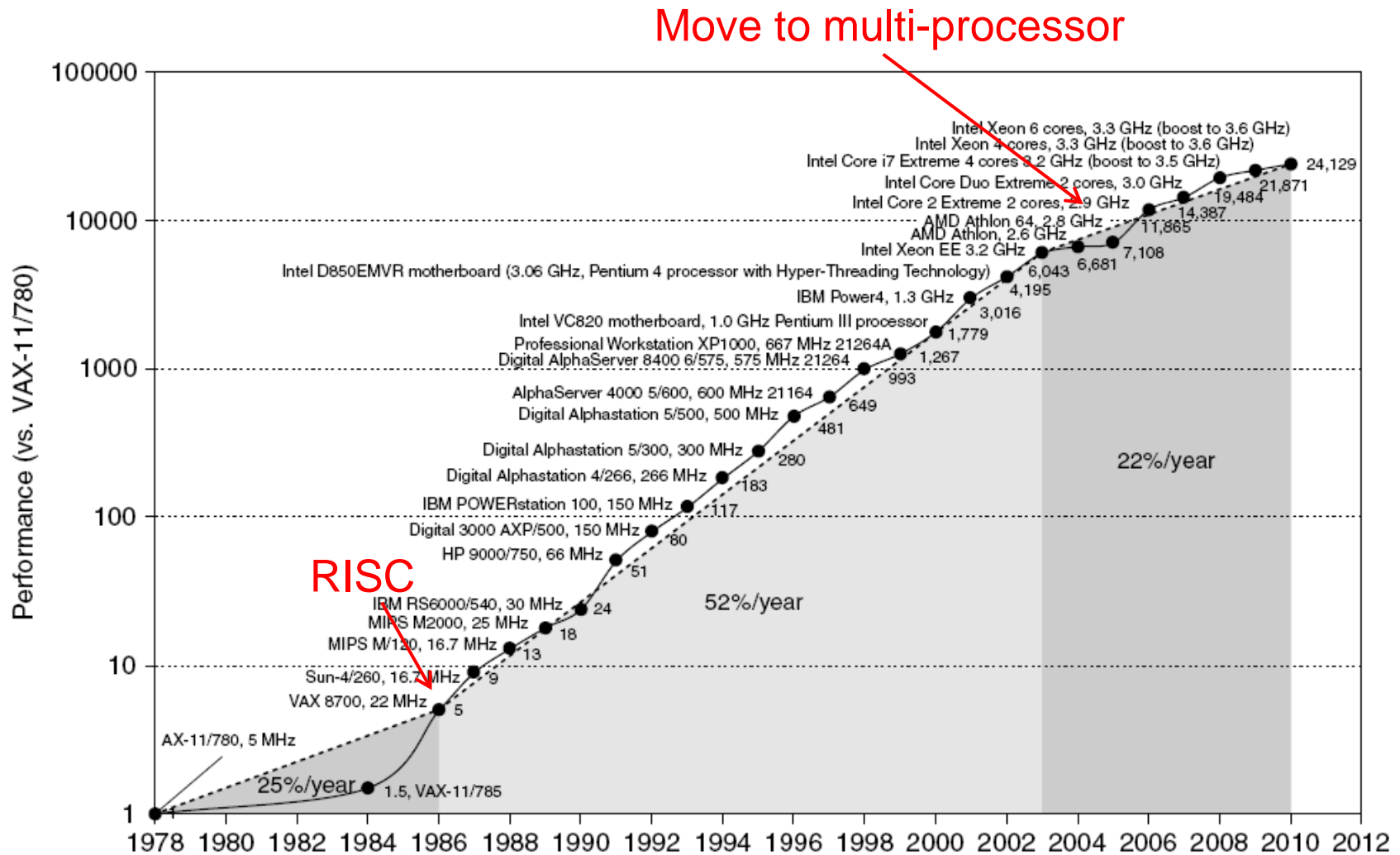
GÜÇ VE ENERJİ

# 4. HAFTA

# Bilgisayar Teknolojisi

- Performans iyileştirmeleri
  - Yarı iletken teknolojideki ilerlemeler
    - Boyut, saat hızı
  - Bilgisayar mimarisindeki iyileştirmeler
    - HLL derleyicileri, UNIX
    - RISC mimarileri
- Tüm ilerlemeler:
  - Hafif bilgisayarlar
  - Verimli yorumlanan programlama dilleri

# Tek işlemci Performansı



# Mimarideki mevcut trendler

- Komut seviyeli paralellik (ILP) ilerleyemedi
  - Tek işlemci performansının ilerlemesi 2003 yılında durdu
- Yeni performans modelleri
  - Veri seviyeli paralellik (Data-level parallelism - DLP)
  - İş parçacığı seviyeli paralellik (Thread-level parallelism -TLP)
  - İstek seviyeli paralellik (Request-level parallelism - RLP)
- Tüm uygulamaların yeniden tasarlanmasını gerektirir

# Ne tarafa Doğru Gidiyoruz?

- Modern eğilimler
- Saat hızı iyileştirmeleri yavaşlıyor
  - güç sınırlamaları
- tek çekirdekli bir işlemciyi daha iyi hale getirmek zor
- Çok-çekirdekli sistemler: her yeni işlemci üretimi daha fazla çekirdek içeriyor
- daha iyi programlama modelleri ve verimli çoklu iş parçacığı uygulamalarına ihtiyaç
- daha iyi bellek hiyerarşilerine ihtiyaç
- daha iyi enerji verimliliğine ihtiyaç
- Bazı kullanımlarda, zayıf çekirdekler çekici
- Düşük veri hareketi

# Paralellik

- Uygulamalardaki paralellik sınıfları:
  - Veri seviyeli paralellik (Data-Level Parallelism -DLP)
  - Görev seviyeli paralellik (Task-Level Parallelism -TLP)
- Mimari tabanlı paralellik sınıfları:
  - Talimat seviyeli paralellik (Instruction-Level Parallelism - ILP)
  - Vektör mimarileri / Grafik işlemciler (GPUs)
  - İş parçacığı seviyeli paralellik
  - İstek seviyeli paralellik

# Flynn Taksonomisi

- Tek talimat akışı, tek veri akışı (Single instruction stream, single data stream -SISD)
- Tek talimat akışı, çoklu veri akışı (Single instruction stream, multiple data streams - SIMD)
  - Vector mimarileri
  - Multimedya ilaveleri
  - Grafik işlemciler
- Çoklu talimat akışı, tek veri akışı (Multiple instruction streams, single data stream - MISD)
  - Ticari uygulama yok
- Çoklu talimat akışı, çoklu data akışı (Multiple instruction streams, multiple data streams - MIMD)
  - Sıkı-bağlı MIMD
  - Zayıf- bağli MIMD

# Teknolojideki Trendler

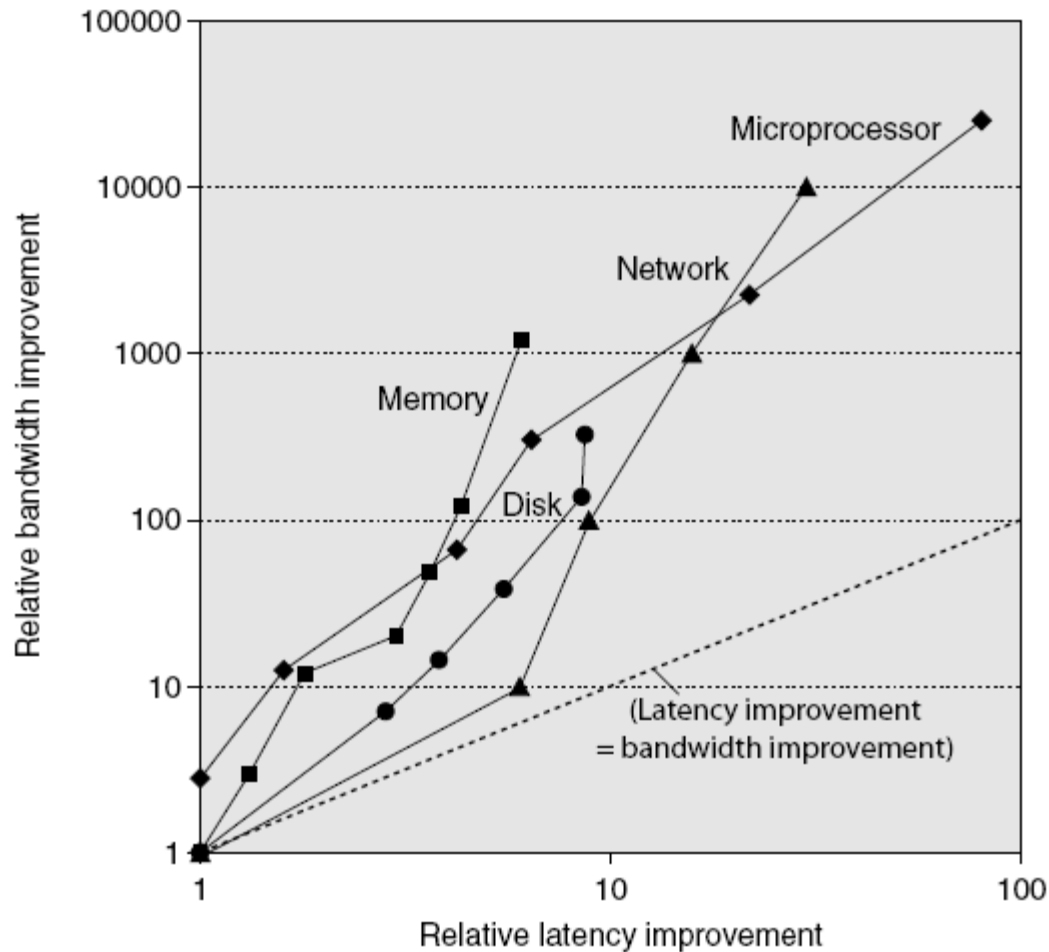
- Entegre devre teknolojisi
  - Transistor yoğunluğu: 35%/yıl
  - Kalıp boyutu (Die size): 10-20%/yıl
  - Tüm entegre: 40-55%/yıl
- DRAM kapasitesi : 25-40%/yıl (yavaşlıyor)
- Flash kapasitesi: 50-60%/yıl
  - 15-20X bir başına DRAM'dan daha ucuz
- Magnetic disk teknolojisi: 40%/yıl
  - 15-25X bit başına daha ucuz vs. Flash
  - 300-500X bit başına daha ucuz vs. DRAM



# Bandgeniřlięi ve Gecikme

- Bant geniřlięi veya verim
  - Belirli bir zamanda yapılan toplam iř
  - İřlemciler iin 10.000-25.000X iyileřtirme
  - Bellek ve diskler iin 300-1200X iyileřtirme
- Gecikme veya tepki suresi
  - Bir etkinlięin bařlaması ile tamamlanması arasındaki sure
  - İřlemciler iin 30-80X iyileřtirme
  - Bellek ve diskler iin 6-8X iyileřtirme

# Bandgenişliği ve Gecikme



Log-log plot of bandwidth and latency milestones

# Transistörler ve iletişim

- Nitelik (feature) boyutu
  - X veya Y boyutunda minimum transistör veya tel boyutu
  - 1971'de 10 mikron - 2024'de .002 mikron
  - Transistör performansı doğrusal olarak artar
    - Kablo gecikmesi nitelik boyutu ile iyileşmez!
  - Entegre yoğunluğu kuadratik olarak artar

## MOSFET scaling (process nodes)

|                   |         |
|-------------------|---------|
| 10 $\mu\text{m}$  | – 1971  |
| 6 $\mu\text{m}$   | – 1974  |
| 3 $\mu\text{m}$   | – 1977  |
| 1.5 $\mu\text{m}$ | – 1981  |
| 1 $\mu\text{m}$   | – 1984  |
| 800 nm            | – 1987  |
| 600 nm            | – 1990  |
| 350 nm            | – 1993  |
| 250 nm            | – 1996  |
| 180 nm            | – 1999  |
| 130 nm            | – 2001  |
| 90 nm             | – 2003  |
| 65 nm             | – 2005  |
| 45 nm             | – 2007  |
| 32 nm             | – 2009  |
| 22 nm             | – 2012  |
| 14 nm             | – 2014  |
| 10 nm             | – 2016  |
| 7 nm              | – 2018  |
| 5 nm              | – ~2020 |

## Future

|      |         |
|------|---------|
| 3 nm | – ~2021 |
| 2 nm | – ~2024 |

# Güç ve Enerji

- Problem: gücü al, gücü ver
- Termal Tasarım Gücü (TDP)
  - Sürdürülebilir güç tüketimini tanımlar
  - Güç kaynağı ve soğutma sistemi tasarımı için kullanılır
  - En yüksek güçten daha küçük, ortalama güç tüketiminden daha yüksek
- Saat hızı güç tüketimini kısıtlamak için çalışma anında düşürülebilir
- Görev başına enerji çoğu kez daha iyi bir ölçüttür

# Dinamik Enerji ve Güç

- Dinamik enerji:
  - Transistor 0 -> 1 veya 1 -> 0 anahtarlansın
  - $\frac{1}{2} \times \text{Kapasitif Y\ddot{u}k} \times \text{Gerilim}^2$
- Dinamik Güç:
  - $\frac{1}{2} \times \text{Kapasitif Y\ddot{u}k} \times \text{Gerilim}^2 \times \text{Anahtarlama Frekansı}$
- Saat hızını düşürmek gücü düşürür, enerjiyi değil

# Trendler

- Transistör ve voltaj başına kapasitans azalıyor, ancak transistörlerin sayısı daha hızlı bir şekilde artıyor; dolayısıyla saat frekansı sabit tutulmalıdır
- Güç kaybı da artıyor;
- Güç kaybı
  - transistör sayısı, kaçak akım ve sağlanan voltajın bir fonksiyonudur

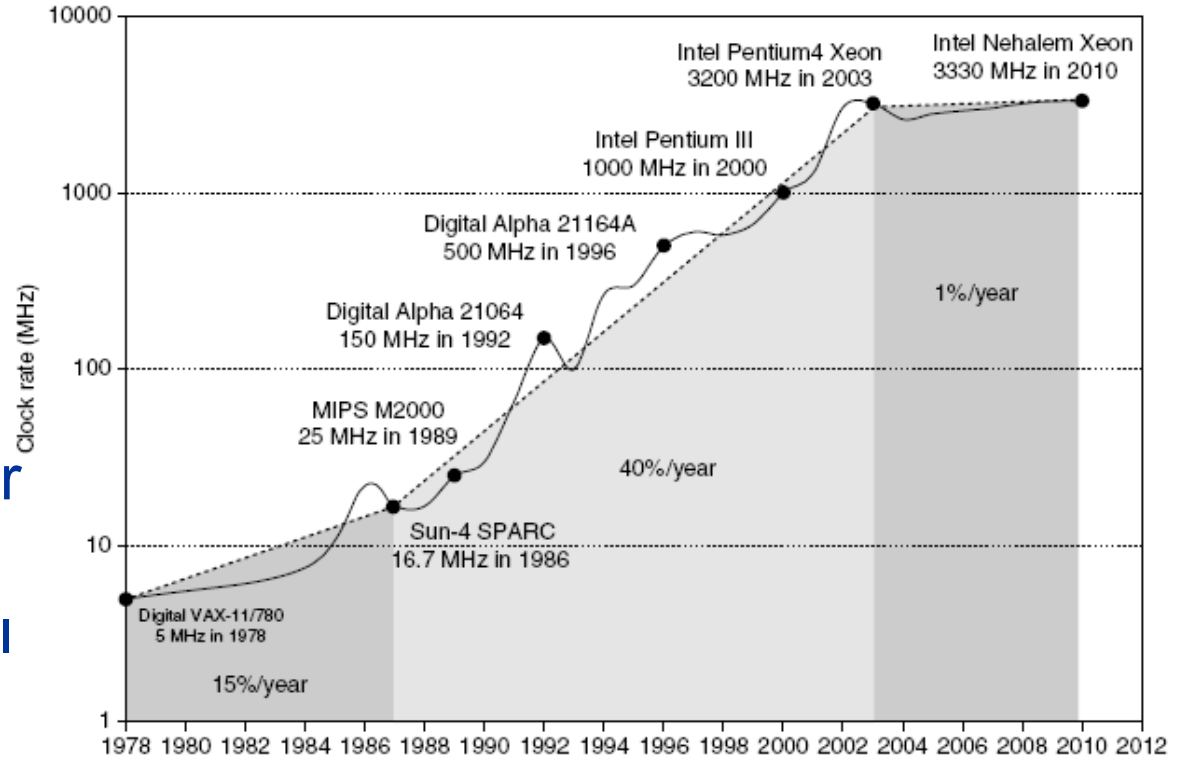
# Enerji

- Bugün yüksek performanslı işlemcilerde güç tüketimi 100-150W arasında

**Enerji = güç x zaman = (dynpower + lkgpower) x time**

# Güç

- Intel 80386 tüketim~ 2 W
- 3.3 GHz Intel Core i7 tüketimi 130 W
- Isı 1.5 x 1.5 cm boyutlu çipten uzaklaştırılmalıdır
- Hava ile soğutmanın sınırı





# Gücü Düşürme

- Gücü düşürme teknikleri:
  - Bir şeyi hakkıyla yapmama
  - Dinamik voltaj-frekans ayarı
  - DRAM ve diskler için düşük güç durumu
  - Overclocking, çekirdeklerin kapatılması

# Statik Güç

- Statik güç tüketimi
  - $Akım_{static} \times Güç$
  - Transistör sayısı ile büyür
  - Düşürmek için: power gating

# Maliyetteki Trendler

- Maliyetler, öğrenme eğrisi sayesinde düşürdü
  - verim
- DRAM: fiyat maliyeti yakından izler
- Mikroişlemciler: fiyat hacme bağlıdır
  - Hacim iki kat arttıkça 10% daha az

# Entegre devre maliyeti

## ■ Entegre devre

$$\text{Cost of integrated circuit} = \frac{\text{Cost of die} + \text{Cost of testing die} + \text{Cost of packaging and final test}}{\text{Final test yield}}$$

$$\text{Cost of die} = \frac{\text{Cost of wafer}}{\text{Dies per wafer} \times \text{Die yield}}$$

$$\text{Dies per wafer} = \frac{\pi \times (\text{Wafer diameter}/2)^2}{\text{Die area}} - \frac{\pi \times \text{Wafer diameter}}{\sqrt{2} \times \text{Die area}}$$

## ■ Bose-Einstein formülü:

$$\text{Die yield} = \text{Wafer yield} \times 1 / (1 + \text{Defects per unit area} \times \text{Die area})^N$$

- Birim alan başına kusur(Defects per unit area) = 0.016-0.057 kusur/cm<sup>2</sup> (2010)
- N = proses karmaşıklık çarpanı = 11.5-15.5 (40 nm, 2010)

# Örnek

**Example** Some microprocessors today are designed to have adjustable voltage, so a 15% reduction in voltage may result in a 15% reduction in frequency. What would be the impact on dynamic energy and on dynamic power?

**Answer** Because the capacitance is unchanged, the answer for energy is the ratio of the voltages

$$\frac{\text{Energy}_{\text{new}}}{\text{Energy}_{\text{old}}} = \frac{(\text{Voltage} \times 0.85)^2}{\text{Voltage}^2} = 0.85^2 = 0.72$$

which reduces energy to about 72% of the original. For power, we add the ratio of the frequencies

$$\frac{\text{Power}_{\text{new}}}{\text{Power}_{\text{old}}} = 0.72 \times \frac{(\text{Frequency switched} \times 0.85)}{\text{Frequency switched}} = 0.61$$

shrinking power to about 61% of the original.

---

# Die- Kalıp-entegre bloğu?

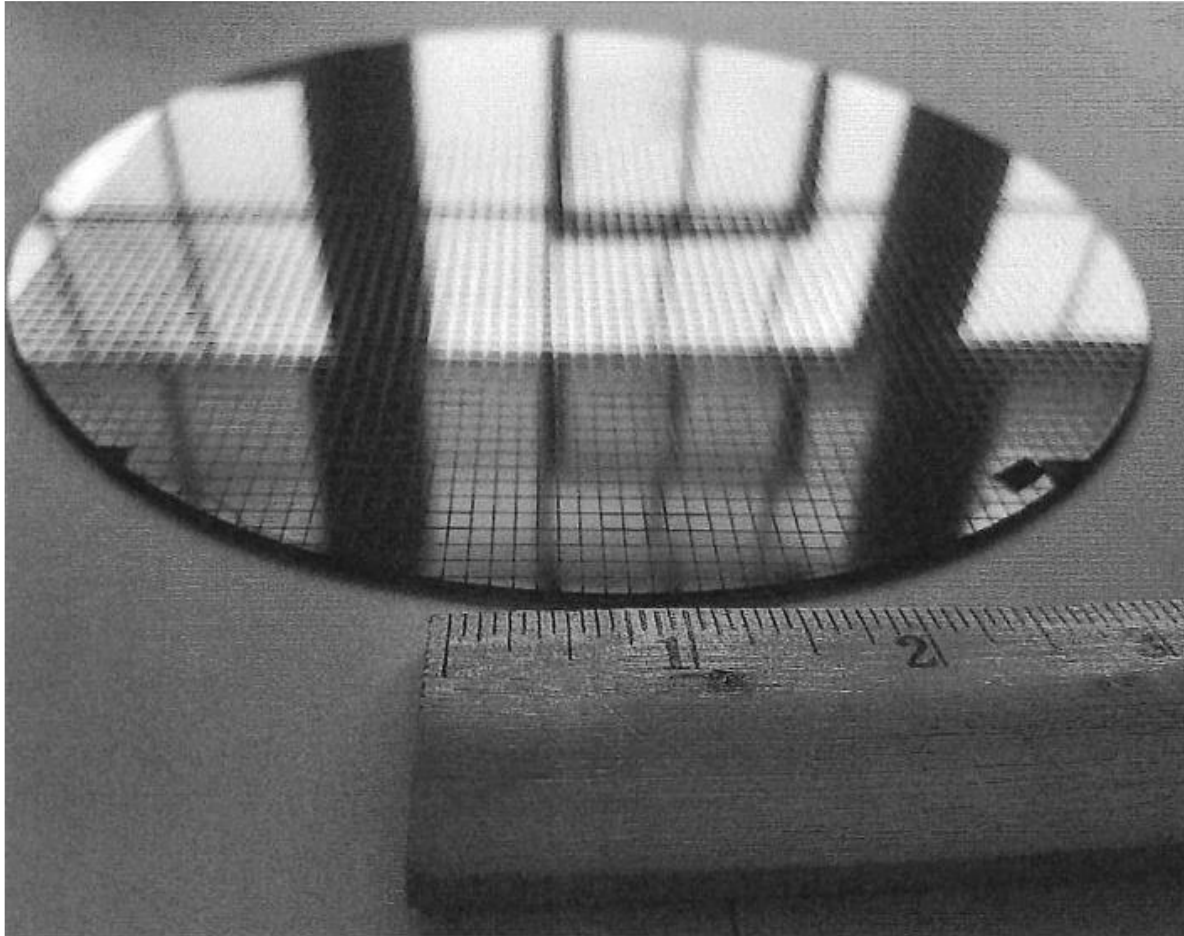


Photograph of an Intel Skylake microprocessor die,



# Wafer-silikon devre levhası

- This 200mmdiameter wafer of RISC-V dies was designed by SiFive



# Örnek

**Example** Find the number of dies per 300 mm (30 cm) wafer for a die that is 1.5 cm on a side and for a die that is 1.0 cm on a side.

**Answer** When die area is 2.25 cm<sup>2</sup>:

$$\text{Dies per wafer} = \frac{\pi \times (30/2)^2}{2.25} - \frac{\pi \times 30}{\sqrt{2 \times 2.25}} = \frac{706.9}{2.25} - \frac{94.2}{2.12} = 270$$

Because the area of the larger die is 2.25 times bigger, there are roughly 2.25 as many smaller dies per wafer:

$$\text{Dies per wafer} = \frac{\pi \times (30/2)^2}{1.00} - \frac{\pi \times 30}{\sqrt{2 \times 1.00}} = \frac{706.9}{1.00} - \frac{94.2}{1.41} = 640$$

---



# Örnek

**Example** Find the die yield for dies that are 1.5 cm on a side and 1.0 cm on a side, assuming a defect density of 0.047 per  $\text{cm}^2$  and  $N$  is 12.

**Answer** The total die areas are 2.25 and 1.00  $\text{cm}^2$ . For the larger die, the yield is

$$\text{Die yield} = 1 / (1 + 0.047 \times 2.25)^{12} \times 270 = 120$$

For the smaller die, the yield is

$$\text{Die yield} = 1 / (1 + 0.047 \times 1.00)^{12} \times 640 = 444$$

The bottom line is the number of good dies per wafer. Less than half of all the large dies are good, but nearly 70% of the small dies are good.

# Güvenilirlik

- Modül güvenirliliği
  - Ortalama hata zamanı (Mean time to failure -MTTF)
  - Ortalama tamir zamanı (Mean time to repair -MTTR)
  - Hatalar arasındaki ortalama zaman (Mean time between failures - MTBF) =  $MTTF + MTTR$
  - Erişilebilirlik =  $MTTF / MTBF$

# Güvenilirlik ve Erişilebilirlik

- Bir sistem aşağıdaki iki şey arasında değişir:
  - Servis başarısı: servis şartnameyi yerine getiriyor
  - Servis kesintisi: servisler şartnameyi yerine getiremiyor
- Bu geçiş, arızalar ve restorasyonlardan kaynaklanır
- Güvenilirlik, sürekli servis başarısını ölçer
- Genellikle ortalama arıza süresi (MTTF) olarak ifade edilir.
- Erişilebilirlik, servisin şartnamenin yerine getirdiği kısmını ölçer
- $MTTF / (MTTF + MTTR)$  olarak hesaplanır

# Örnek

**Example** Assume a disk subsystem with the following components and MTTF:

- 10 disks, each rated at 1,000,000-hour MTTF
- 1 ATA controller, 500,000-hour MTTF
- 1 power supply, 200,000-hour MTTF
- 1 fan, 200,000-hour MTTF
- 1 ATA cable, 1,000,000-hour MTTF

Using the simplifying assumptions that the lifetimes are exponentially distributed and that failures are independent, compute the MTTF of the system as a whole.

**Answer** The sum of the failure rates is

$$\begin{aligned}\text{Failure rate}_{\text{system}} &= 10 \times \frac{1}{1,000,000} + \frac{1}{500,000} + \frac{1}{200,000} + \frac{1}{200,000} + \frac{1}{1,000,000} \\ &= \frac{10 + 2 + 5 + 5 + 1}{1,000,000 \text{ hours}} = \frac{23}{1,000,000} = \frac{23,000}{1,000,000,000 \text{ hours}}\end{aligned}$$

or 23,000 FIT. The MTTF for the system is just the inverse of the failure rate

$$\text{MTTF}_{\text{system}} = \frac{1}{\text{Failure rate}_{\text{system}}} = \frac{1,000,000,000 \text{ hours}}{23,000} = 43,500 \text{ hours}$$

or just under 5 years.

# Örnek

- 100 W'da %100 kapasitede çalışan bir işlemcinin %20'si kaçak güç olarak harcanmaktadır. Bu işlemci %50 kapasitede çalıştığında toplam güç sarfiyatı ne olur?

# Örnek

- 100 W'da %100 kapasitede çalışan bir işlemcinin %20'si kaçak güç olarak harcanmaktadır. Bu işlemci %50 kapasitede çalıştığında toplam güç sarfiyatı ne olur?

$$\begin{aligned}\text{Total power} &= \text{dynamic power} + \text{leakage power} \\ &= 80\text{W} \times 50\% + 20\text{W} \\ &= 60\text{W}\end{aligned}$$

# Güç vs. Enerji

- Enerji bize sabit bir görevi yerine getirmenin gerçek “maliyetini” anlatır
- Güç (enerji / zaman) kısıtlamalar demektir; sadece güç dağıtımını veya soğutma çözümünü maksimize edecek kadar hızlı çalışabilir
- A işlemcisi B işlemcisinden 1,2 kat fazla güç tüketirse, tüketir ancak bir görevi % 30 daha kısa sürede bitirirse,

Bağıl enerjisi  $1,2 \times 0,7 = 0,84$  olur.

Proc-A daha iyidir.

# Örnek

A işlemcisi B işlemcisinden 1,4 kat fazla güç tüketir, ancak görevi % 20 daha kısa sürede bitirirse, hangi işlemciyi seçersiniz:

- (a) Eğer güç dağıtım kısıtlamaları varsa
- (b) Operasyon başına enerjiyi en aza indirmeye çalışıyorsanız?
- (c) Cevap sürelerini en aza indirmeye çalışıyorsanız?



# Örnek

A işlemcisi B işlemcisinden 1,4 kat fazla güç tüketir, ancak görevi % 20 daha kısa sürede bitirirse, hangi işlemciyi seçersiniz:

(a) Eğer güç kısıtlamaları varsa

Proc-B

(b) Operasyon başına enerjiyi en aza indirmeye çalışıyorsanız?

Proc-A  $1.4 \times 0.8 = 1.12$  kat daha çok enerji harcar

(c) Cevap sürelerini en aza indirmeye çalışıyorsanız?

Proc-A daha hızlıdır, ancak Proc-B'nin frekansını (ve gücünü) artırabilir ve Proc-A'nın tepki süresini eşleştirebiliriz (yine de güç ve enerji açısından daha iyi olur)

# Gücü ve enerjiyi düşürmek

- Aktif olmayan transistörler kapatılabilir (kayıbı azaltır)
- Tipik durumu belirle ve etkinlik belirli bir eşiği aştığında
- DFS: Dinamik frekans ölçeklendirme - yalnızca frekansı ve dinamik gücü azaltır, ancak enerjiye artırır
- DVFS: Dinamik voltaj ve frekans ölçeklendirme :
  - voltaj ve frekansı (örneğin)% 10 azalttığımızda; bir program % 8 yavaşlayabilir,
  - ancak dinamik gücü % 27 azaltabilir, toplam gücü % 23 azaltabilir, toplam enerjiyi% 17 azaltabilir
  - Not: voltaj düşmesi → transistör yavaşlaması → frekans düşmesi

# Örnek

3 GHz'de çalışan A işlemcisi 80 W dinamik güç ve 20 W statik güç tüketir. Bir programı 20 saniyede tamamlar.

- a) Frekansı% 20 oranında düşürürsem enerji tüketimi ne olur?
- b) Frekansı ve voltajı% 20 oranında düşürürsem enerji tüketimi nedir?

# Örnek

3 GHz'de çalışan A işlemcisi 80 W dinamik güç ve 20 W statik güç tüketir. Bir programı 20 saniyede tamamlar.

a) Frekansı% 20 oranında düşürürsem enerji tüketimi ne olur?

Yeni dinamik güç = 64W; Yeni statik güç = 20W

Yeni çalışma zamanı = 25 sn

Enerji = 84 W x 25 sn = 2100 Joules

b) Frekansı ve voltajı% 20 oranında düşürürsem enerji tüketimi nedir?

Yeni DP = 41W; Yeni statik güç = 16W;

Yeni çalışma zamanı = 25 sn; Enerji = 1425 Joules

# Diğer Teknoloji Trendleri

- DRAM yoğunluğu yılda % 40-60 artar, gecikme süresi 10 yılda % 33 azalır,
- Bant genişliği gecikme azaldıkça iki kat daha hızlı iyileşir
- Disk yoğunluğu her yıl % 100 artar, DRAM'a benzer gecikme süresi de artar
- DRAM ve sabit disk sürücülerini arasında bir köprü sağlayabilen NVRAM teknolojileri
- Ayrıca, güvenilirlik konusundaki artan endişeler (transistörler daha küçük olduğundan, düşük voltajlarda çalışıyor ve çok fazla var)

# Maliyet

- Maliyet birçok faktör tarafından belirlenir: hacim, verim, üretim olgunluğu, işleme adımları, vb.
- Önemli bir belirleyici: çip alanı
- Küçük alan → wafer başına daha fazla çip
- Küçük alan → herhangi bir kusur durumunda daha küçük bir alan zayi olur, yani verim artar
- Kabaca söylemek gerekirse, alanın yarısı = maliyetin üçte biridir