

# Küme Bilgisayarlar

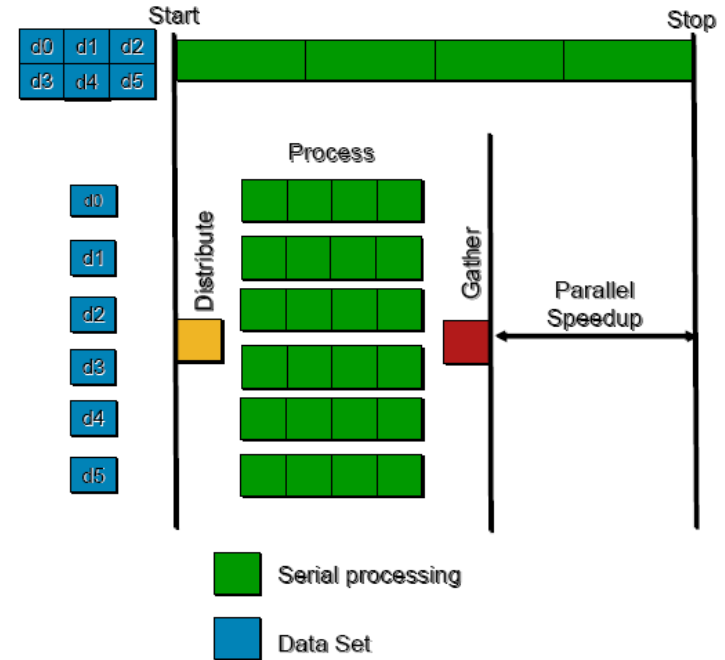
*Onur Temizsoylu*

*Bilgisayar Mühendisliği, ODTÜ, Ankara*

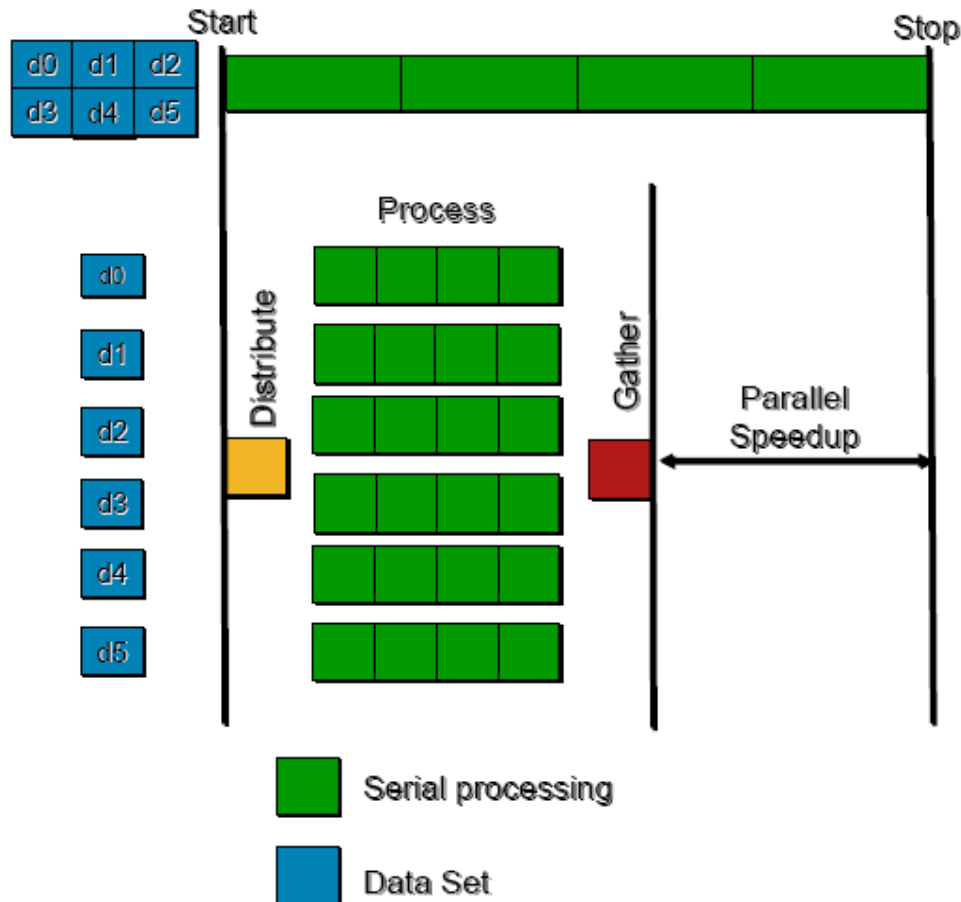
*04/02/2008*

- Terminoloji
- Paralleleştirme, Paralel Donanım Mimarileri
- Neden hesaplamada kümeleme?
- Kümeleme nedir?
  - Yüksek kullanılabilirlik kümeleri
  - Yük dengeleme kümeleri
  - Veritabanı kümeleri
  - Yüksek başarımlı hesaplama kümeleri
- YBH küme mimarileri
  - Sunucular
  - Ağ bağlantısı
  - Depolama
  - İşletim sistemi
  - Orta katman yazılımları
  - Uygulama programlama arayüzleri

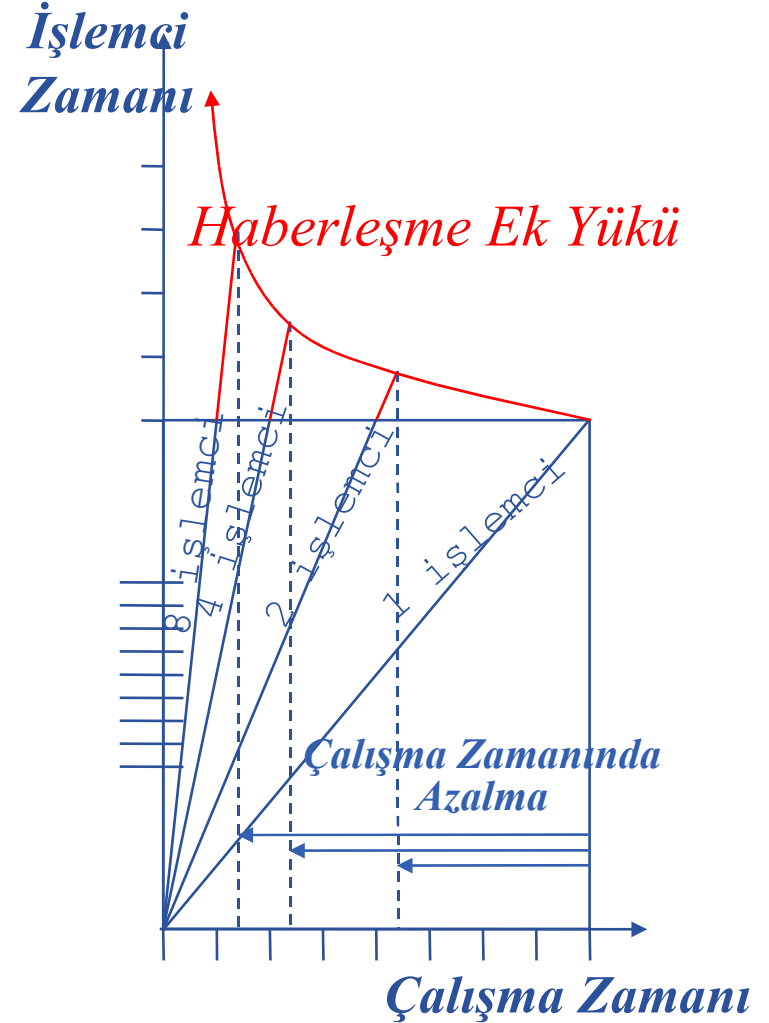
- Süreç (“Process”)
- İş Parçacığı (“Thread”)
- Görev (“Task”)
- Hızlanma (“Speedup”)
- Ölçeklenebilirlik (“Scalability”)
- Verimlilik
- Senkronizasyon (“Synchronization”)
- Paralel Ek Yükü (“Parallel Overhead”)
- Süperbilgisayar



- Bir işin paralelleştirilmesinde programın toplam çalışma zamanını azaltmak amaçlanır.



- Ek yük:
  - İşlemcilerde fazladan geçen süre
  - İletişim ek yükü
  - Senkronizasyon ek yükü
  - Programın paralel olmayan/olamayan parçaları
- Paralel programlamada ek yük ve çalışma zamanı hızlanma ve verimlilik ile ifade edilir.



- $i$  sayıda işlemcide programın toplam işlemci zamanını  $Z(i)$  olarak ifade edelim.

$$\text{Hızlanma } (i) = Z(1) / Z(i)$$

$$\text{Verimlilik } (i) = \text{Hızlanma } (i) / i$$

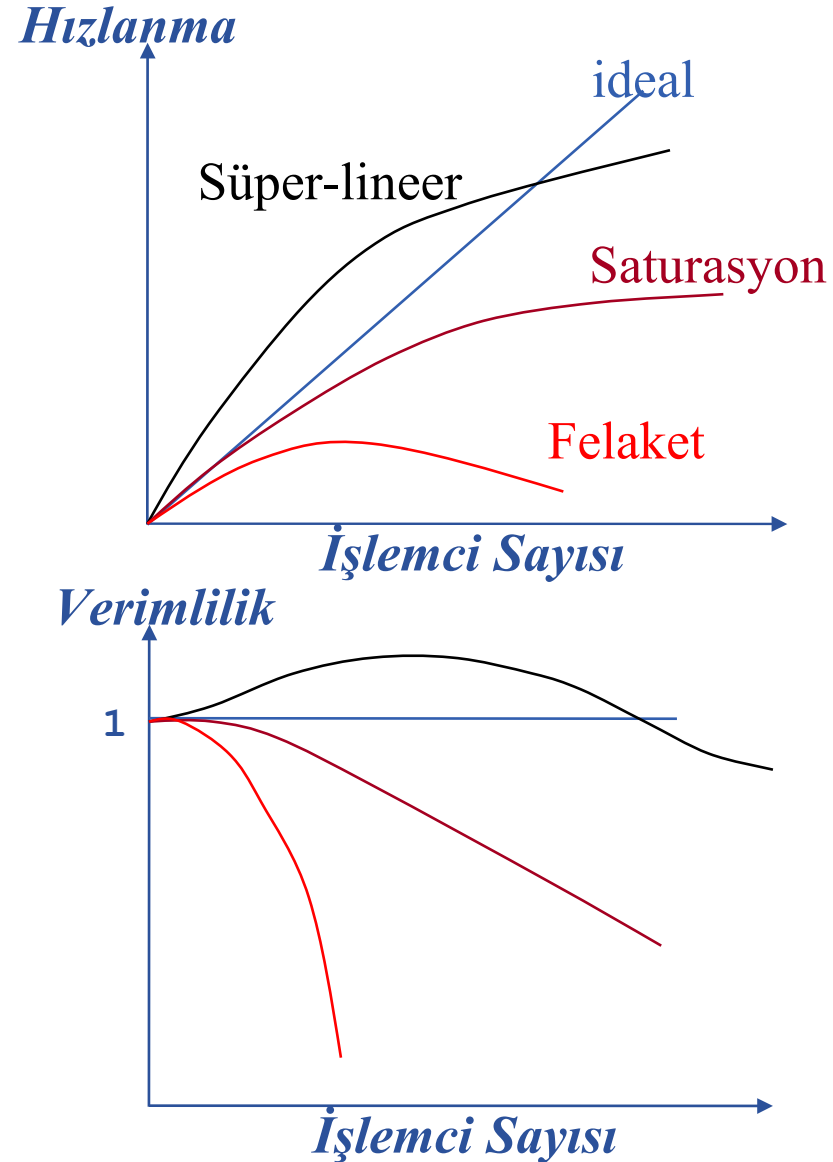
- İdeal durumda:

$$Z(i) = Z(1) / i$$

$$\text{Hızlanma } (i) = i$$

$$\text{Verimlilik } (i) = 1$$

- Ölçeklenebilir programlar büyük işlemci sayılarında bile verimli kalırlar.



- Amdahl yasası:
  - “Kodun paralel olmayan kısmı (ek yük), kodun ölçeklenebilirliği konusunda üst limiti oluşturur.”
- Kodun seri kısmını  $s$ , paralel kısmını  $p$  olarak ifade edersek:

$$1 = s + p$$

$$Z(1) = Z(s) + Z(p)$$

$$= Z(1) * (s + p)$$

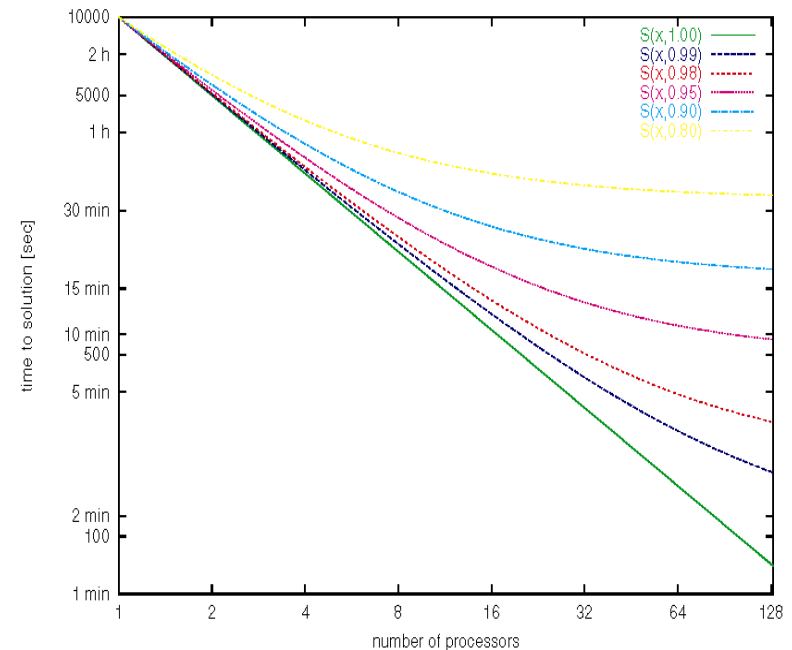
$$= Z(1) * (p + (1-p))$$

$$Z(i) = Z(1) * (p/i + (1-p))$$

$$\text{Hızlanma}(i) = Z(1) / Z(i)$$

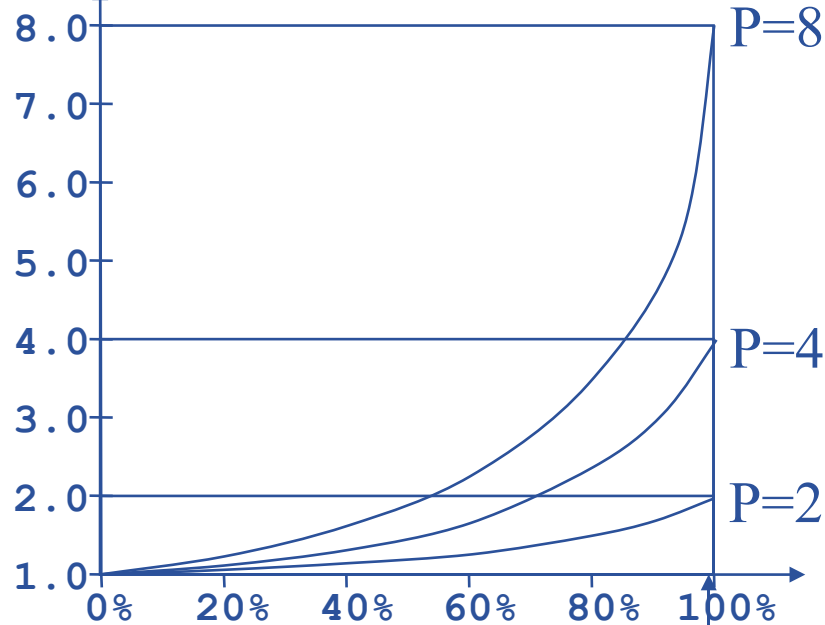
$$= 1 / (p/i + 1 - p)$$

$$\text{Hızlanma}(i) < 1 / (1 - p)$$



- Pratikte programları paralelleştirmek Amdahl yasasında görüldüğü kadar zor değildir.
- Ancak programın çok büyük bir kısmını paralel işlem için harcaması gereklidir.

*Hızlanma*



*Kodda Paralel Kısım*

1970s →

1980s →

1990s →

*En iyi paralel kodlar  
~99% diliminde*

David J. Kuck,  
Hugh Performance Computing,  
Oxford Univ.. Press 1996

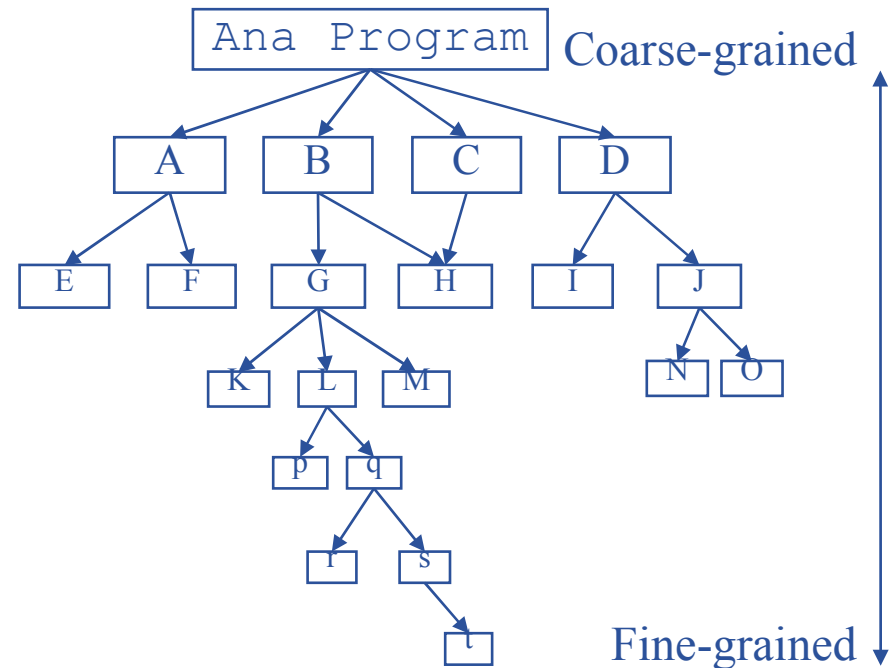


## — Fine-Grained:

- Genelde her döngüde paralelleştirme vardır.
- Çok sayıda döngü paralelleştirilir.
- Kodun çok iyi bilinmesine gerek yoktur.
- Çok fazla senkronizasyon noktası vardır.

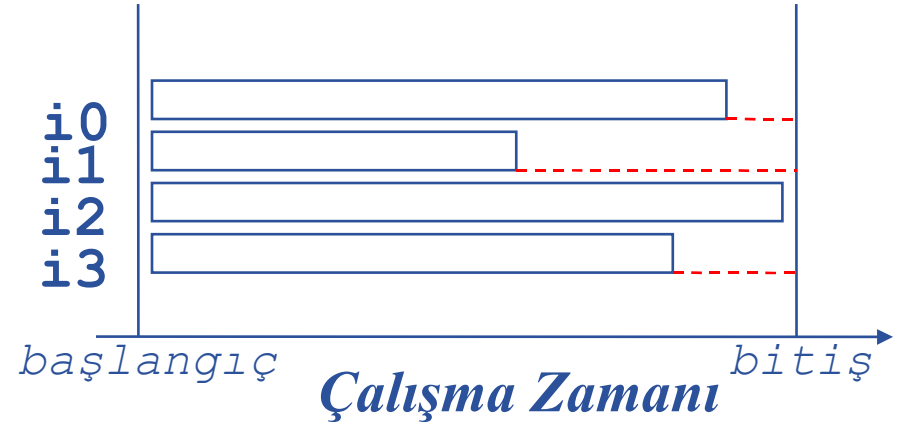
## — Coarse-Grained:

- Geniş döngülerle paralelleştirme yapılır.
- Daha az senkronizasyon noktası vardır.
- Kodun iyi anlaşılması gerekir.



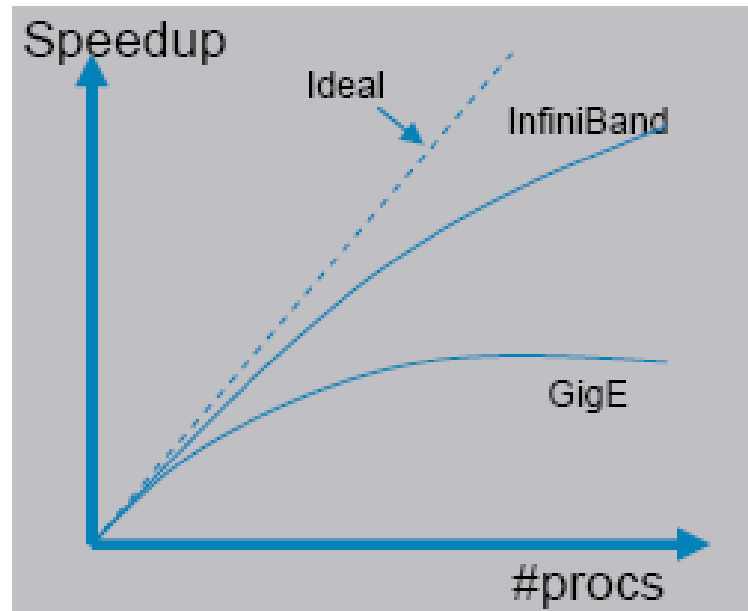
— Ölçeklenebilirliği etkileyen diğer faktörler:

- **İş parçacıkları arası yük dengesizliği** : Bir kodun herhangi bir paralel kısmının çalışma zamanı en uzun süren iş parçacığının çalışma zamanıdır. Coarse-Grained programlamada ortaya çıkması daha olasıdır.
- **Çok fazla senkronizasyon**: Kodda küçük döngüler sırasında her seferinde senkronizasyon yapılırsa bu ek yük getirir. Fine-Grained programlamada ortaya çıkması daha olasıdır.

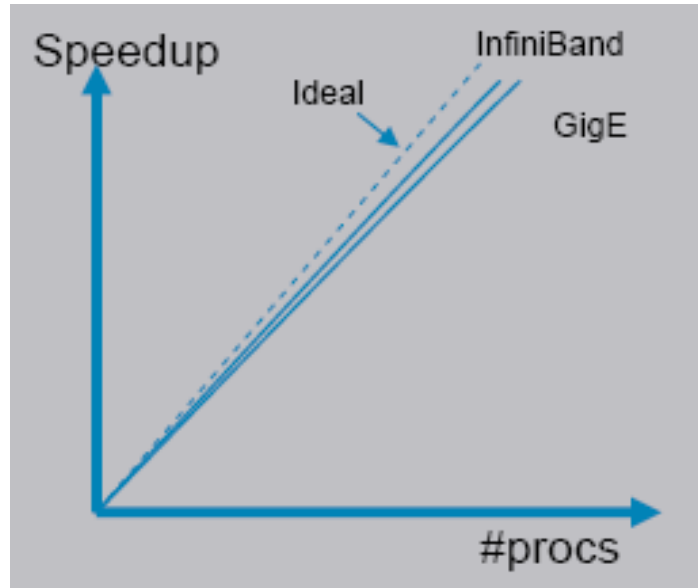


— Sıkı bağlı sistemler:

- Süreçler arasında yoğun haberleşme
- Gecikme süresine hassas
- Ortak Bellek Paralel
- Dağıtık Bellek Paralel

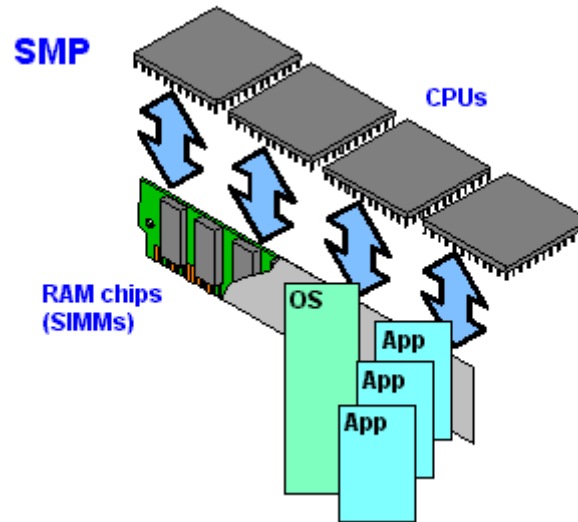


- Gevşek bağlı sistemler:
  - Süreçler arasında haberleşme azdır veya hiç yoktur.
  - Gecikme süresine hassas değildir. Ancak bant genişliği veri transferi için etkili olabilir.
- Parametrik çalışan uygulamalar
  - Süreçler arasında haberleşme yoktur.
  - Kümelerde, grid altyapılarında çalışan uygulamaların çoğunluğunu oluştururlar.

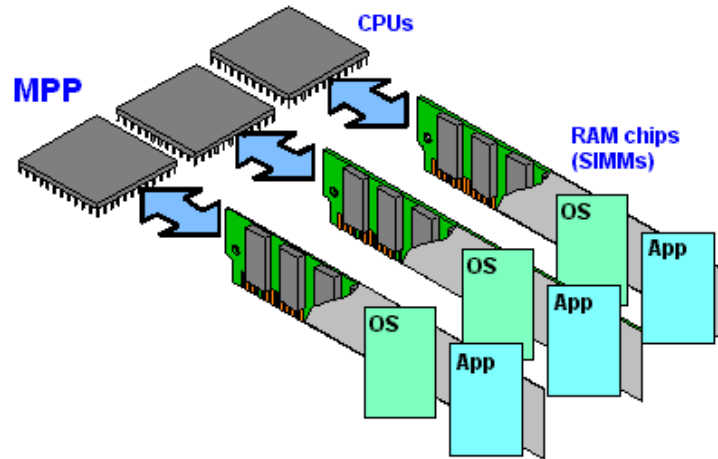


- SMP makineler
- MPP makineler
- NUMA makineler
- Superscalar işlemciler
- Vektör makineler
- Küme bilgisayarlar

- SMP, birden fazla eş işlemcinin ortak bir belleğe bağlandığı çok işlemcili bir bilgisayar mimarisidir.
- SMP sistemler, görevleri işlemciler arasında paylaşabilirler.
- SMP sistemler, paralel hesaplama için kullanılan en eski sistemlerdir ve hesaplamalı bilimlerde yoğun bir şekilde kullanılırlar.

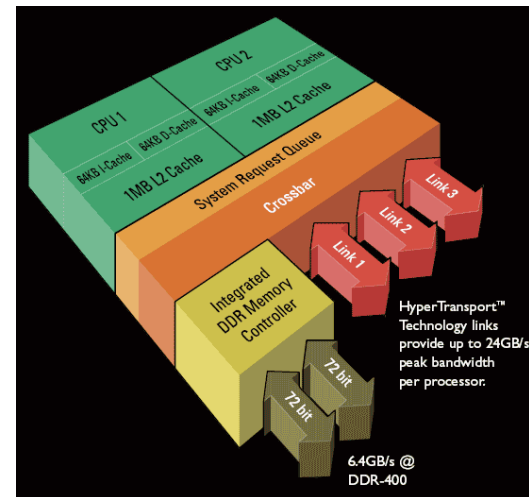
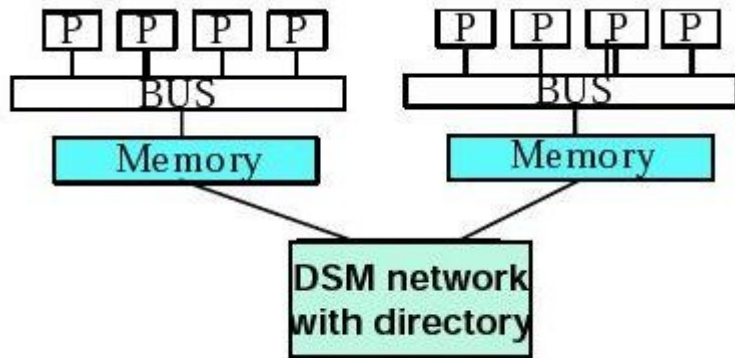


- MPP, binlerce işlemci kullanılabilen çok işlemcili bir mimaridir.
- Bir MPP sisteminde her işlemci kendi belleğine ve işletim sistemi kopyasına sahiptir.
- MPP sistemler üzerinde çalışacak uygulamalar eş zamanda çalışacak eş parçalara bölünebilmelidirler.
- MPP sistemlere yeni işlemci ekledikten sonra uygulamalar yeni paralel kısımlara bölünmelidirler. SMP sistemler ise bundan çok iş parçacığı çalıştırabilir yapıları sayesinde hemen faydalanırlar.



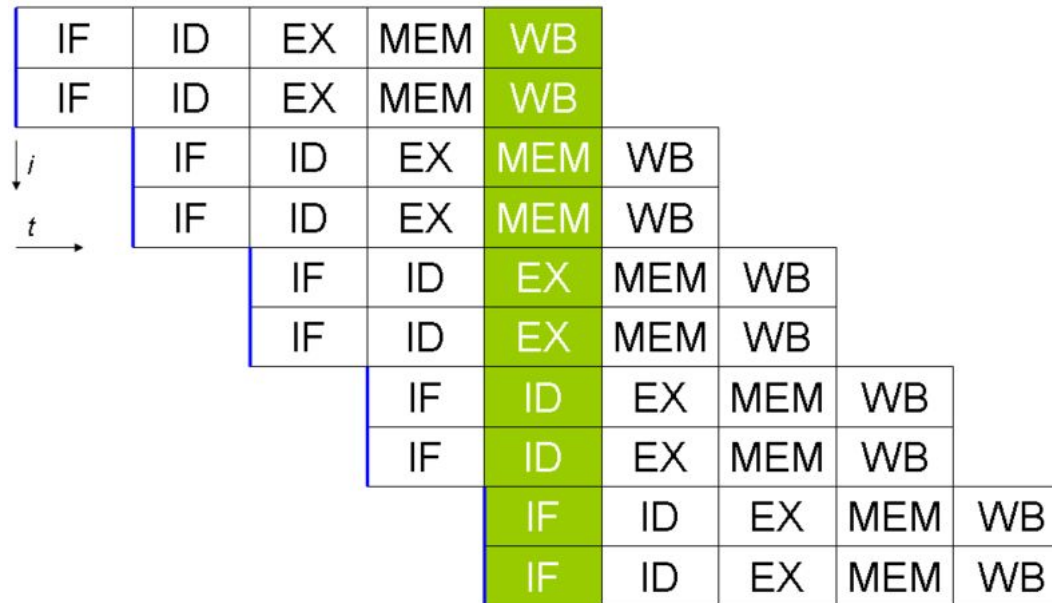
From Computer Desktop Encyclopedia  
© 1998 The Computer Language Co., Inc.

- NUMA, çok işlemcili makinelerde bellek erişim zamanının bellek yerine göre değiştiği bir bellek tasarımıdır.
- İlk defa 1990'larda ortaya çıkmıştır.
- Modern işlemciler, belleklere hızlı bir şekilde erişmeye ihtiyaç duyarlar. NUMA, istenen verinin “cache” bellekte bulunamaması, belleğin başka işlemci tarafından kullanılması gibi performans sorunlarını her işlemciye bellek vererek aşar.
- Intel Itanium ve AMD Opteron işlemciler ccNUMA tabanlıdır.



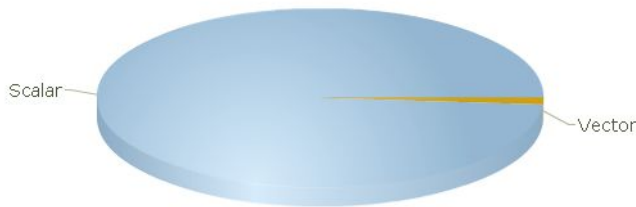


- 1998 senesinden beri üretilen bütün genel amaçlı işlemciler “superscalar” işlemcilerdir.
- “Superscalar” işlemci mimarisi, tek bir işlemcide makine kodu seviyesinde paralellik sağlar.
- “Superscalar” bir işlemci tek bir basamakta birden fazla işlem yapar.

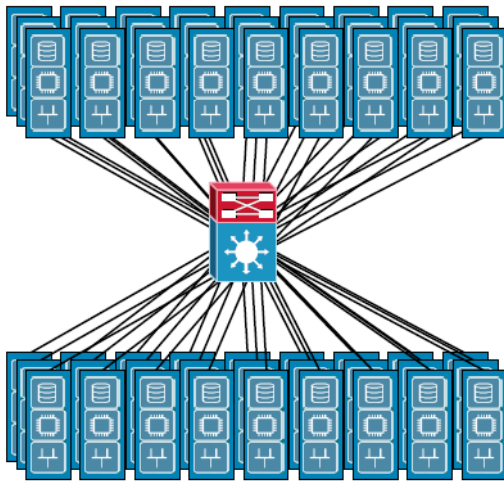


- Vektör işlemciler, aynı anda birden fazla veri üstünde matematik işlem yapabilen işlemcilerdir.
- Şu anda süperbilgisayar dünyasında vektör işlemciler çok az kullanılmaktadırlar.
- Ancak bugün çoğu işlemci vektör işleme komutları içermektedirler (Intel SSE).
- Vektör işlemciler, aynı matematiksel komutu farklı veriler üzerinde defalarca çalıştırmak yerine bütün veri yığını alıp aynı işlemi yapabilirler.

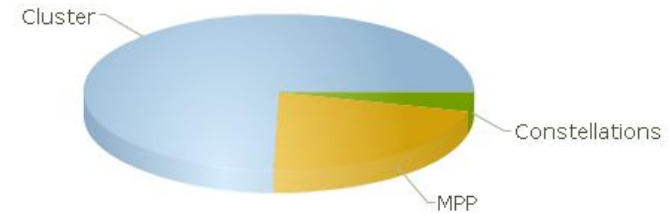
Processor Architecture / Systems  
June 2007



- Hesaplama küme bilgisayar kullanımı 1994 senesinde NASA'da Beowulf projesi ile başlamıştır. 16 Intel 486 DX4 işlemci ethernet ile bağlanmıştır.
- Yüksek performanslı hesaplama, artık küme bilgisayarlarla hesaplama halini almıştır.
- Küme bilgisayar, birlikte çalışmak üzere bağlanmış birden fazla sunucudan oluşur.
- En önemli dezavantajı kullanıcıya tek sistem arayüzü sunamamasıdır.

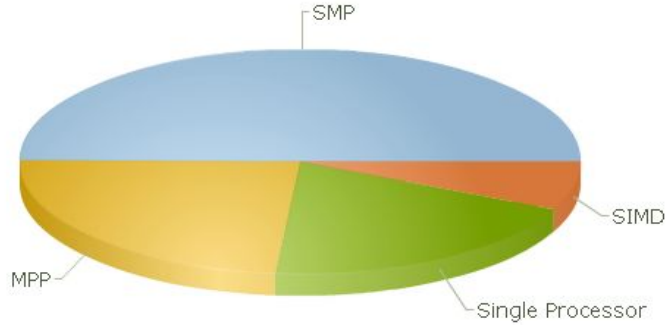


Architecture / Systems  
June 2007

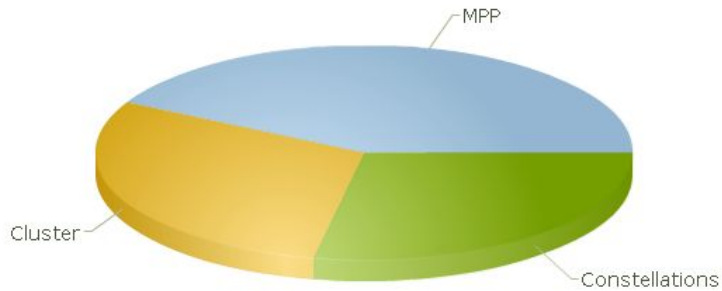


- TOP500 Listesine göre son 15 sene içinde süperbilgisayar sistemlerinde mimari değişimi

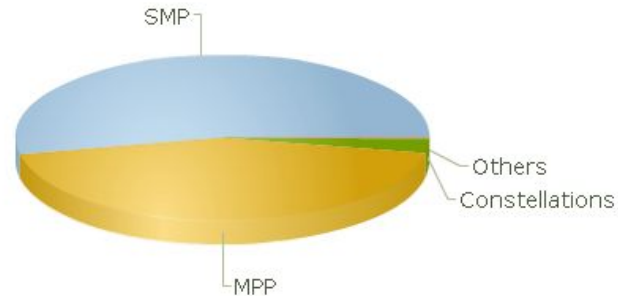
Architecture / Systems  
June 1993



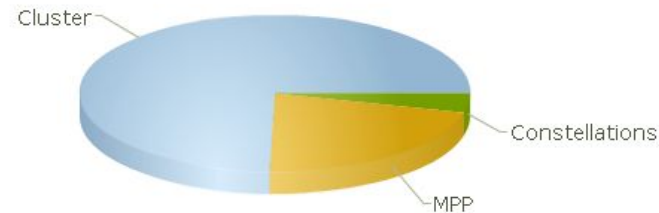
Architecture / Systems  
June 2003



Architecture / Systems  
June 1998

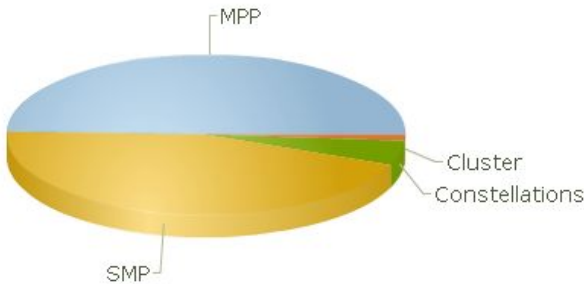


Architecture / Systems  
June 2007

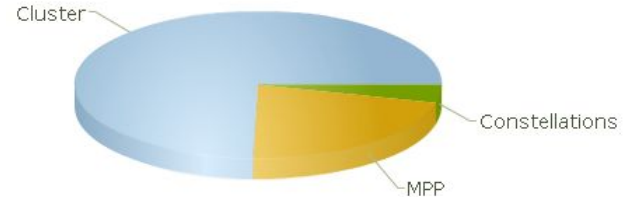


- Fiyat / performans
- Standardı oturmuş işletim sistemi, mesajlaşma gibi yazılım katmanları (Linux, MPI, OpenIB)
- Genişleyebilir, standardı oturmuş bağlantı teknolojileri (Gigabit Ethernet, Infiniband, 10 Gigabit Ethernet)
- Son senelerde süperbilgisayarların büyük bir kısmı küme bilgisayarlardan oluşmaktadır:

Architecture / Systems  
June 1999



Architecture / Systems  
June 2007

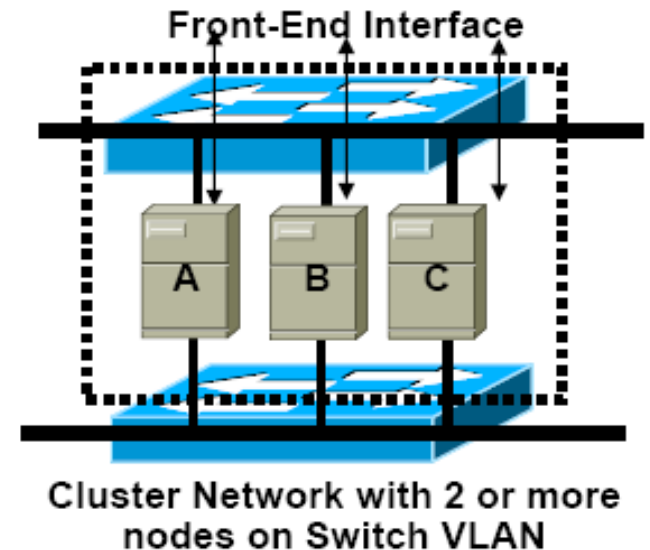
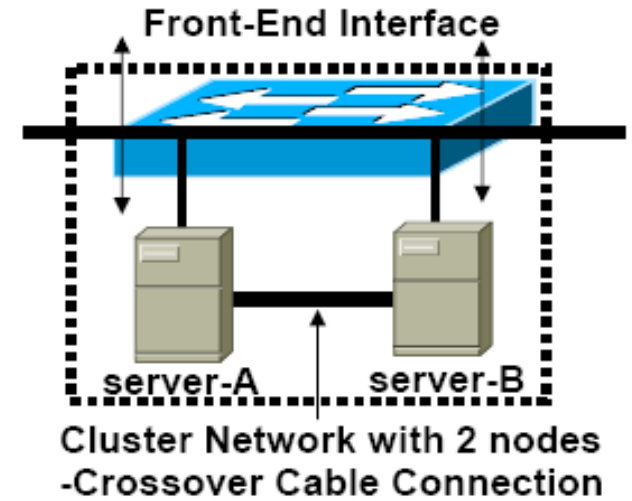


- Kümeleme iki veya daha fazla bilgisayarı:
  - Uygulama veya servis kullanılabilirliğini arttırmak için,
  - Yük dengelemek için,
  - Dağıtık ve yüksek başarılı hesaplama için ağ ile birleştirmektir.
  
- Kümeleme değişik sistem katmanlarında gerçekleştirilebilir:
  - Depolama: Paylaşılmış disk, ikizlenmiş disk, paylaşılmayan veri
  - İşletim Sistemi: UNIX/Linux kümeleri, Microsoft (?) kümeleri
  - Uygulama Programlama Arayüzü: PVM, MPI
  - Uygulamalar

- Küme bilgisayarların önemli mimari dezavantajları vardır:
  - Ortak bellek yoktur.
  - İletişim bellek okuma/yazma hızına göre yavaştır.
- Bu kısıtlamalar uygulama için önemlidir. Uygulamanın bunlara göre de geliştirilmesi gerekebilir.
- Güç ve klima için genelde daha fazla miktarda bütçe gerekir.
- Ölçeklenebilirlik yakalamak bazı uygulamalar için zordur.

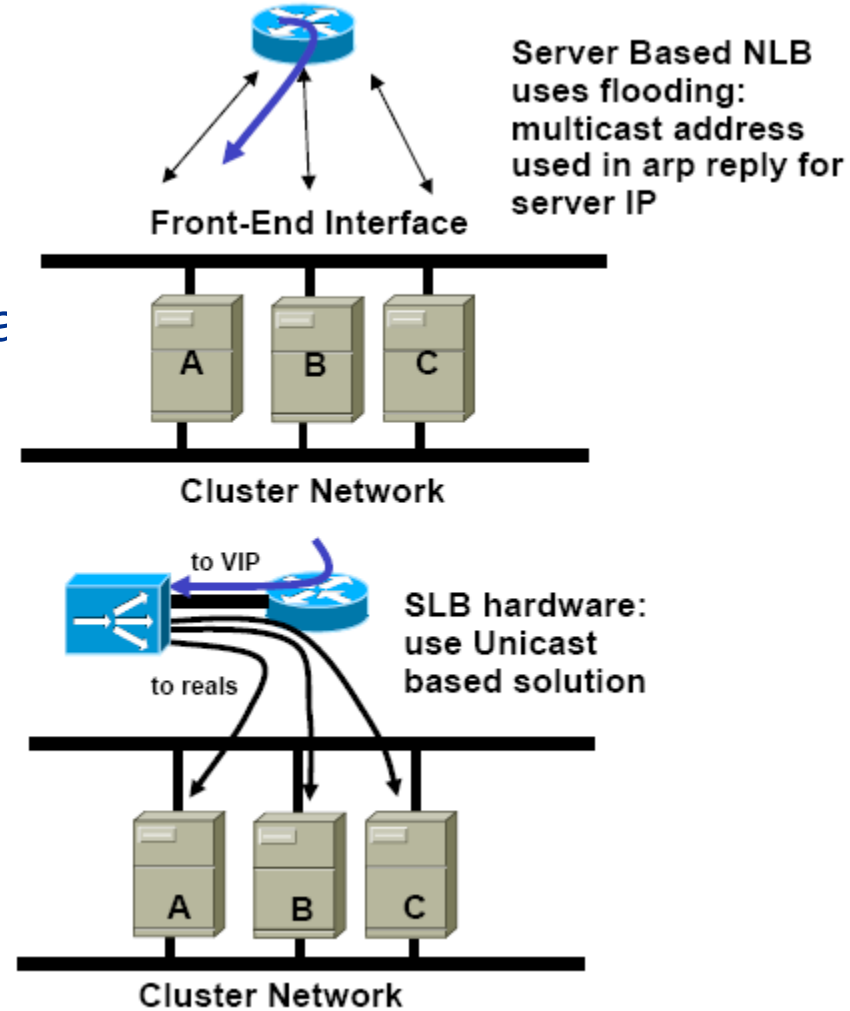


- HA kümeleri, servislerin ayakta kalma sürelerini arttırmak içindir.
- Aynı servisin birden fazla kopyası çevrimiçi veya çevrimdışı bekletilir. Serviste bir sorun olduğu zaman devreye alınır.
- Linux-HA projesi, sıklıkla bu amaçla kullanılan bir yazılımdır.

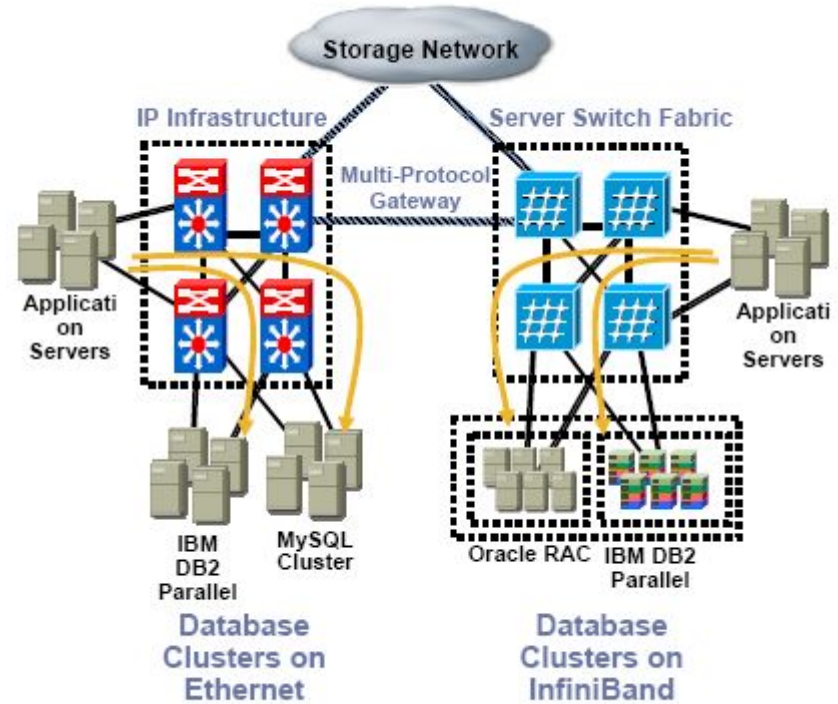




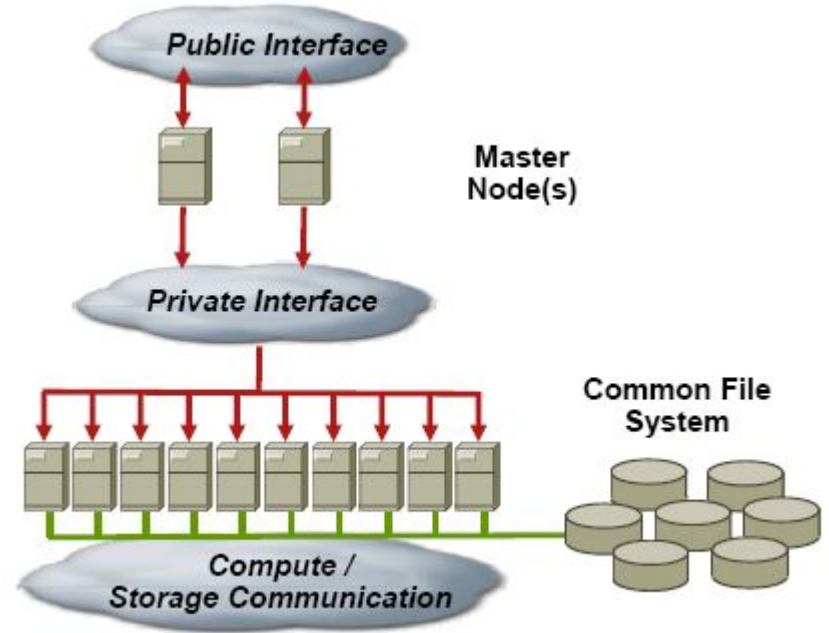
- Yük dengeleme kümeleri, ön arayüzden gelen bütün iş yükünü karşılayıp arkadaki sunuculara aktarırlar.
- Bu kümeler, sunucu çiftliği olarak adlandırılırlar.
- LSF, MAUI, Sun Grid Engine gibi birçok yük dengeleyici yazılım vardır.
- “Linux Virtual Server” projesi de oldukça sık kullanılan bir yük dengeleyici çözümüdür.



- Son senelerde birçok veritabanı üreticisi, yüksek kullanılabilirlik, genişleyebilirlik ve yüksek başarım için kümeleme teknolojilerini için ürün çıkarmıştır.
- Bu çözümlerin bir kısmı paylaştırılmış disk alanı, bir kısmı ayırık veri alanları ile çözüm sunmaktadır.



- Bu kümeler, zaman kritik paralel, seri veya parametrik hesaplama işlerini çalıştırmak için kullanılır.
- Normal bir bilgisayarda inanılmaz sürede bitebilecek işlemci kritik uygulamaları çalıştırır.
- Genellikle normal PC veya sunucular ve Linux ile oluşturulan kümeler Beowulf ismini alırlar.
- MPI, YBH kümelerinde en çok tercih edilen haberleşme kütüphanesidir.



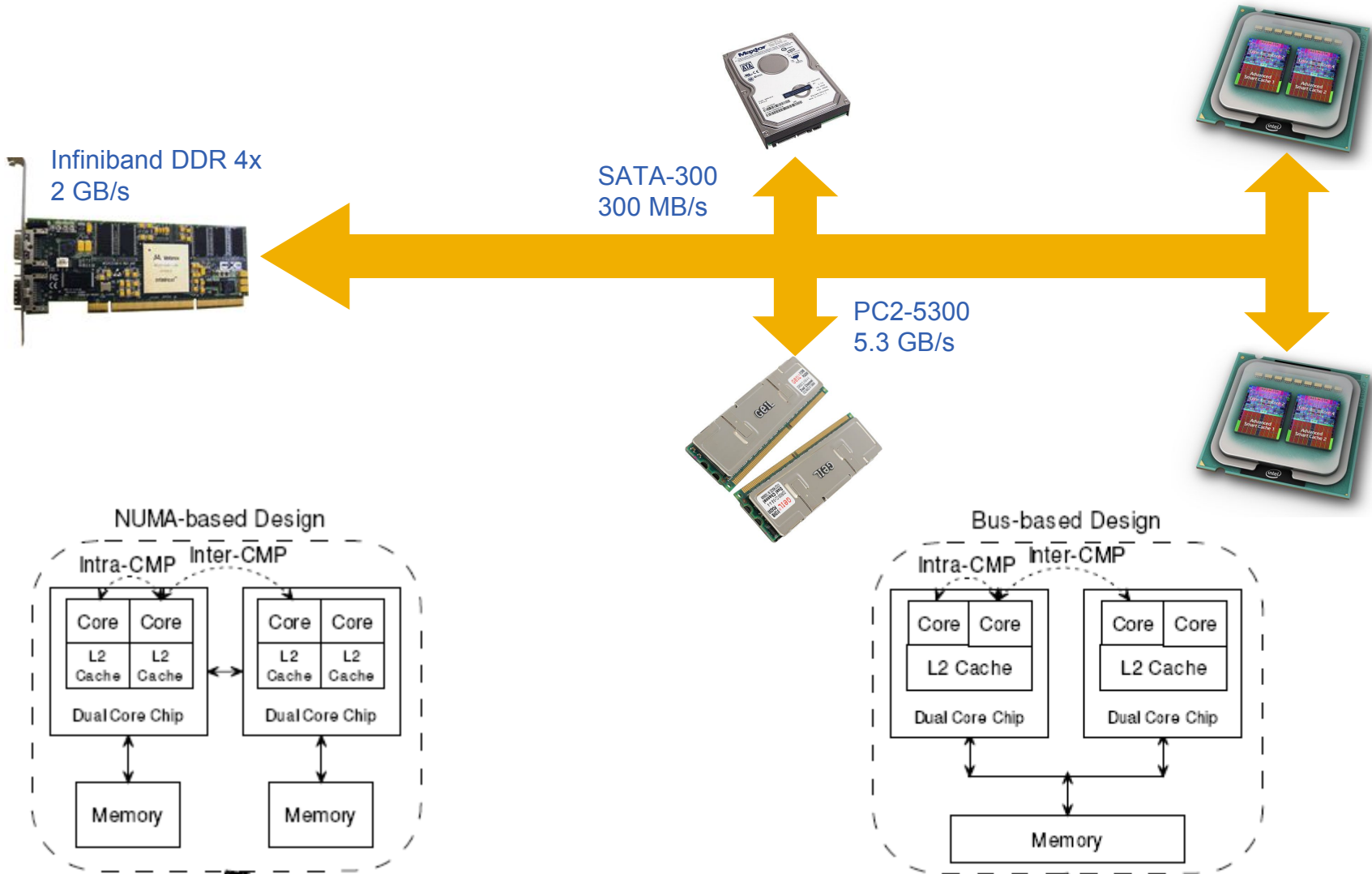
- Yüksek başarımlı hesaplama ihtiyacını karşılamak isteyen bir kullanıcının önünde iki seçenek vardır:
  - Uygulamasına göre küme bilgisayarı edinmek.
  - Erişebildiği küme bilgisayarın özelliklerine göre uygulamasını geliştirmek, değiştirmek veya optimize etmek.
- Her iki durumda da bilinmesi veya hesaplanması gerekenler:
  - Uygulamanın özellikleri, gereksinimleri (yüksek bellek, her sunucuda yüksek miktarda geçici disk alanı, özel kütüphaneler ...),
  - Kümenin büyüklüğü (işlemci, bellek, disk),
  - Ağ bağlantı biçimi (gigabit ethernet, infiniband),
  - İşletim sistemi (Linux, Microsoft (?) ...),
  - Birçok kullanıcı veya grubun birlikte çalışabilirliği,
  - Derleyiciler (GNU, Intel, Portland Group ...)

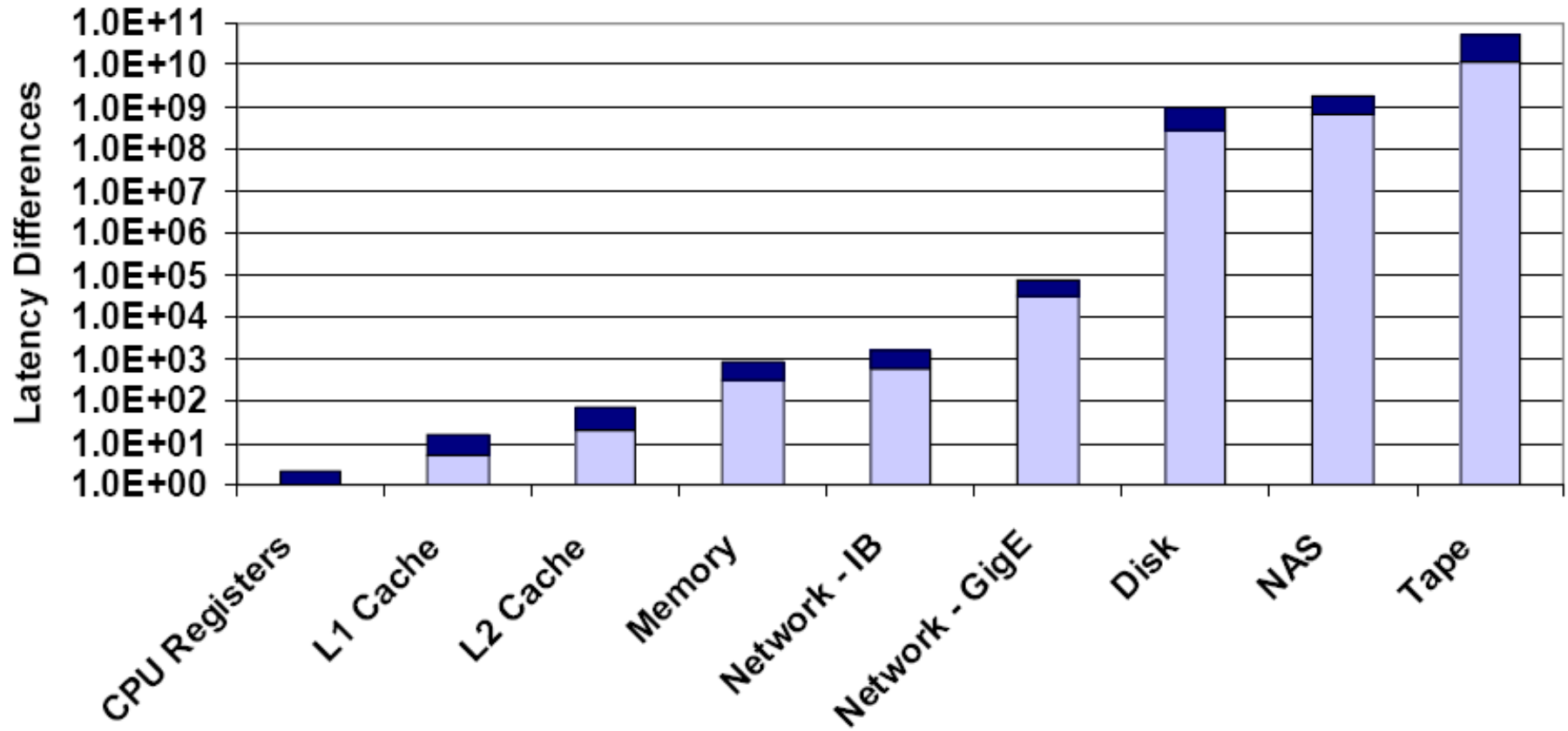




- Günümüzde 1U boyutta 16 çekirdekli sunucular almak mümkün olmaktadır.
- Küme bilgisayarlarda sunucu seçimi konusunda birçok faktör vardır:
  - İşlemciler : Tek çekirdek, çok çekirdek, çoklu işlemci soketi ...
  - Anakart : PCI-X, PCI-Express, HyperTransport ...
  - Sunucu form faktörü : Blade, rack monte, PC ...
  - Bellek : Boyutu, DDR-2, DDR-3, FBDIMM ...
  - Disk : Boyutu, SATA, SCSI, SAS ...
  - Ağ bileşenleri : Gigabit Ethernet, Infiniband, Quadrics ...

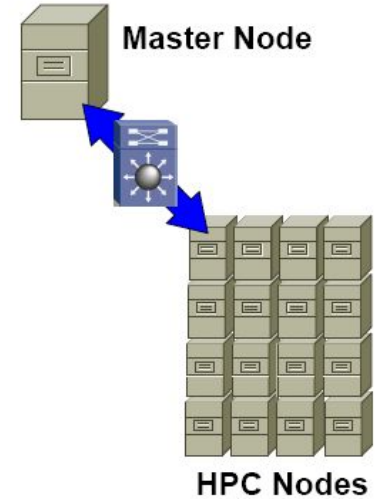
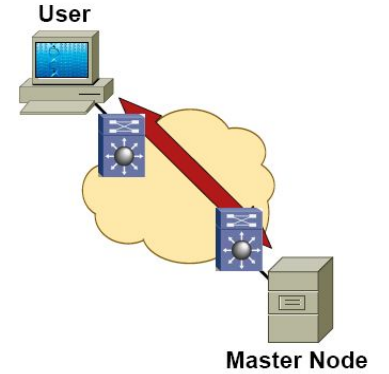






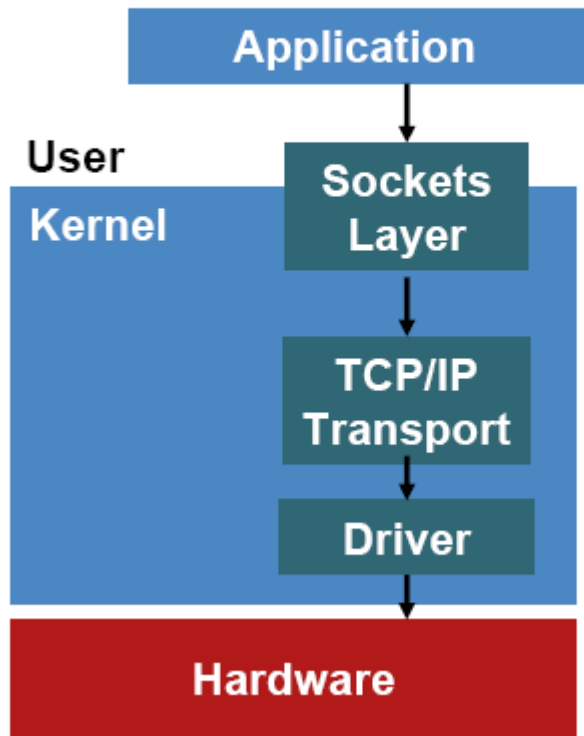


- Genellikle tek bir kümede birden fazla ağ bulunur:
  - Kullanıcı ağı:
    - İş göndermek, görselleştirme, sonuç görüntüleme için kullanılır.
    - Grid haberleşmesi için de kullanılabilir.
    - Kümelere bağlanmak için genellikle ssh kullanılır.
  - Yönetim ağı:
    - İş planlamak, sunucuları izlemek, kurmak için kullanılır.
    - Genellikle IP üzerinden çalışırlar.
    - Ganglia gibi yazılımlar multicast çalışırlar.

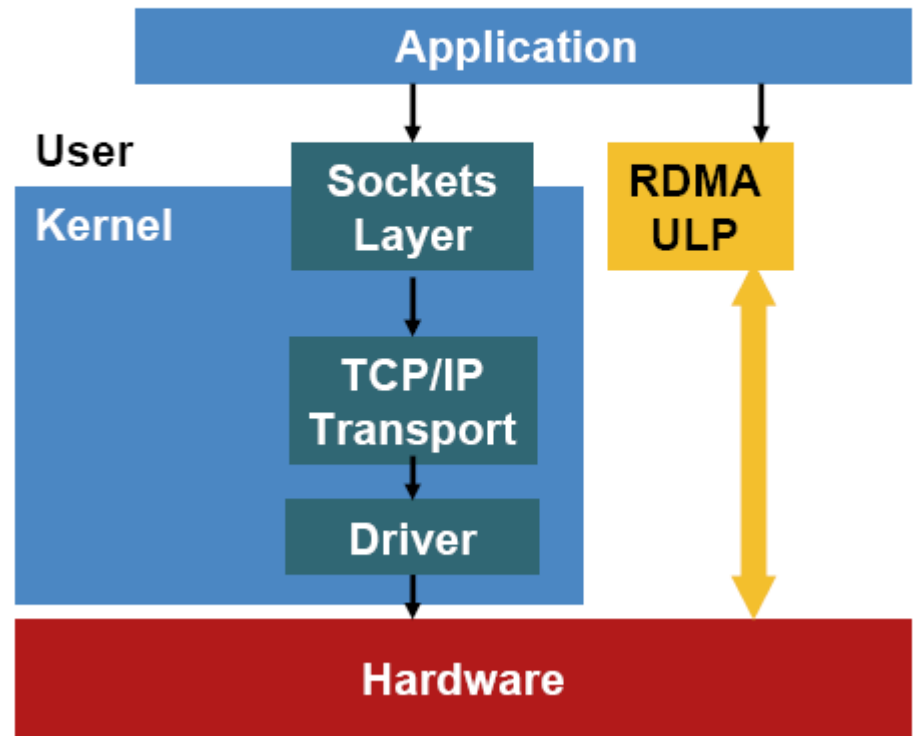


- Küme bilgisayar performansı ve verimi IPC ağı tarafından belirlenir. Haberleşmede harcanan her fazla süre daha az işlem zamanı demektir.
- Günümüzde küçük kümeler ve gevşek bağlı uygulamalar için gigabit ethernet ideal bir çözümdür.
- Büyük kümeler ve sıkı bağlı uygulamalar için Infiniband, Quadrics gibi çözümler vardır.
- Uygulama gereksinimlerini anlamak teknoloji seçiminde çok önemlidir.

## Traditional Model

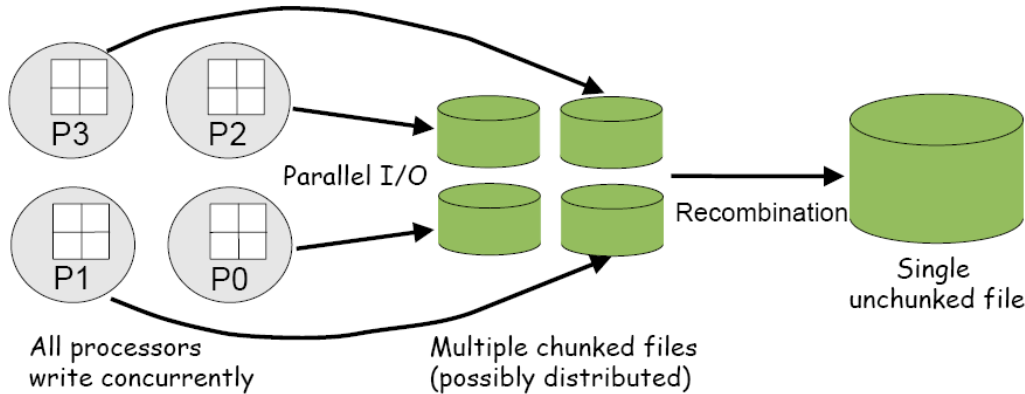
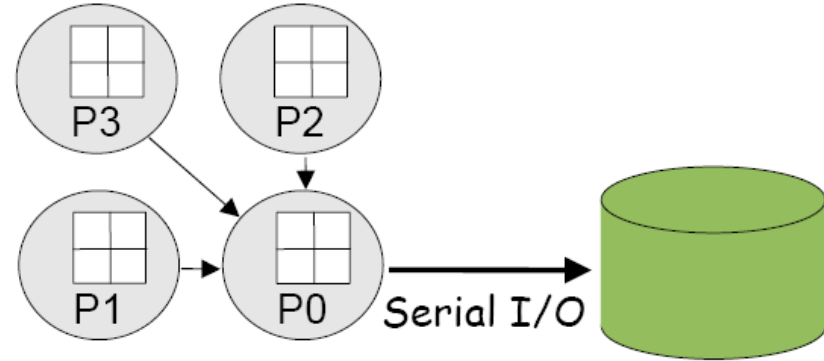


## Kernel Bypass Model

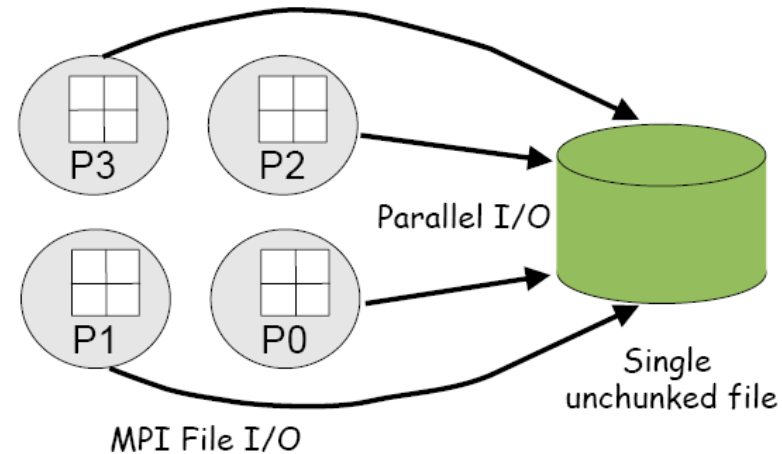


- Küme bilgisayarlarda çalışan kullanıcıların farklı depolama ihtiyaçları bulunur.
  - Ev dizini, uygulamalar için ortak veri alanı
  - Yığın veri saklamak için veri ambarları
  - Yedekleme ve yığın veriler için tape üniteleri
  - Bazı uygulamalar için sunucularda geçici paylaşılmayan disk alanları
- Küme bilgisayarlarda hesaplama yapılan sunucularda kurulum diski veya geçici disk alanı bulundurmamak gerekli değildir. Ancak çoğu durumda maliyeti düşüren bu çözüm tercih edilmemektedir.
- Uygulama performansı için özellikle paylaşılan disk alanlarının ihtiyaca uygun tasarlanması gerekir.

- P0 bütün işlemcilerden veriyi toplar.
- G/Ç işlemcisi darboğazdır.



- Sistemde açık dosya sayısı limitlidir.
- Parçalanmış dosyalar birleştirilmelidir.



- MPI G/Ç, paralel dosya sistemi ihtiyacı duyar.

- Paralel olmayanlar:
  - NFS, CIFS
- Paralel (“Metadata”)
  - Lustre : Ölçeklenebilir
  - Panasas : Ölçeklenebilir
- Paralel (“Metadata” olmadan)
  - XFS
  - IBM GPFS : Ölçeklenebilir
  - PVFS
  - Oracle Cluster FS

