

BMÜ-463 DAĞITIK SİSTEMLER  
FİNAL SINAVI - 11.01.2013

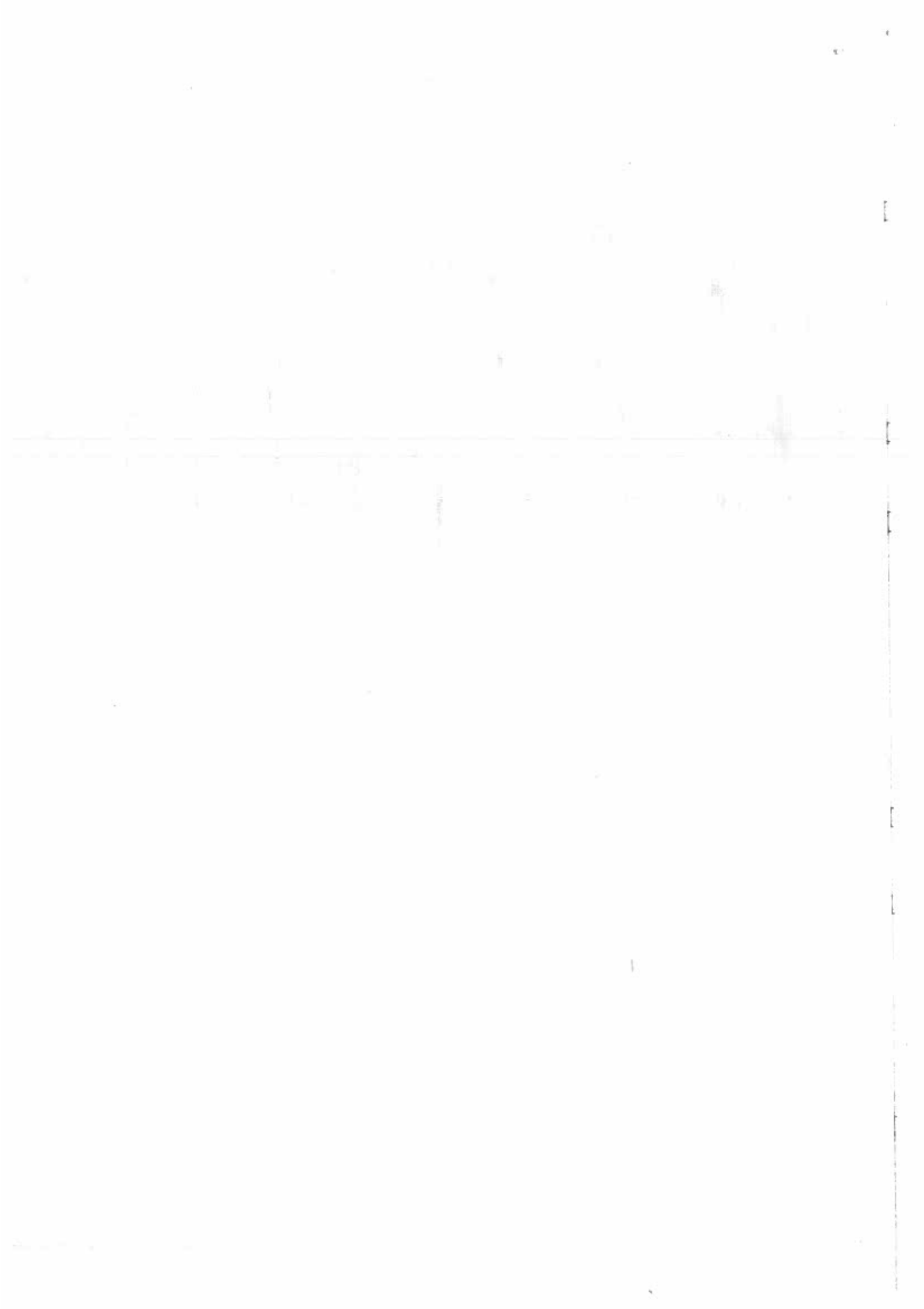
- 1) MapReduce nedir, nasıl çalışır açıklayınız?
- 2-a) Web Servisi nedir, neden ihtiyaç duyulmuştur?
- 2-b) Web Servis standartlarını yazıp kısaca açıklayınız.
- 4) Hizmet Odaklı Mimari nedir, kısaca tartışınız?
- 5) Bulut Hesaplama ve Grid Hesaplama kavramlarını kısaca tartışınız.

03260049  
Merve TELGEKEN

FIRAT ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ  
BMÜ-463 DAĞITIK SİSTEMLER  
ARA SINAV

- 1) Paralel Hesaplama nedir? Hangi amaçlar için yapılır?
- 2) Dağıtık Hesaplama kavramını kısaca açıklayınız.
- 3) Dağıtık Sistemlerde Orta Katman Yazılımı (Middleware) nedir, ne amaçla kullanılır? Şekille açıklayınız.
- 4) Dağıtık Sistemlerde ölçeklenebilirlik kavramını açıklayınız.
- 5) Dağıtık Sistemlerin temel özelliklerini ve karakteristiklerini yazıp kısaca açıklayınız.

Başarılar Dilerim  
Dr. Galip AYDIN



— Vize 27.11.2013 —

13:15

MULTI Core ---

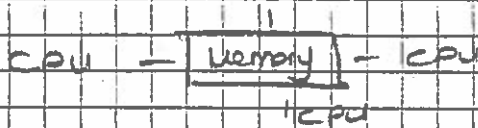
Cluster Computing: Super ölçekli evrensel kullanım  
PC'ler, telefonlar, tablet bilgisayarlar, simülasyon, veri tabanları  
paralel hesaplama yapma, belirli problemleri çözme.

Super Computing / High Performance Computing:  
Çok büyük ölçekli bilgisayarlar. 200.000  
CPU'ları var.

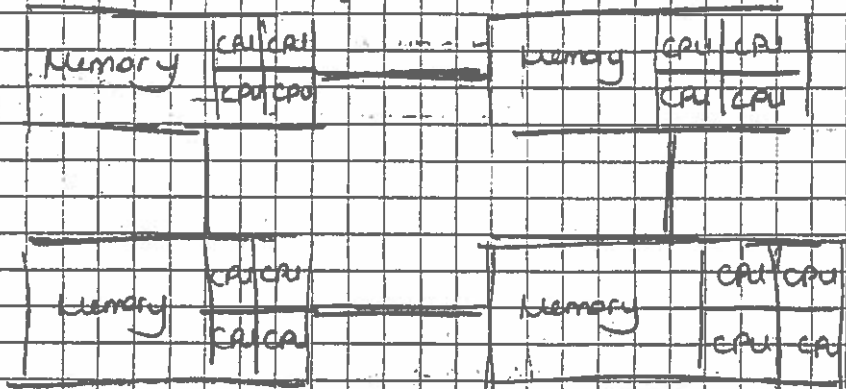
PARALEL HAFIZA YAPILARI:

- (1) Paylaşımlı Hafıza Yapısı: Cache coherency olabilir de  
olmayabilir de.  
Birden fazla işlemci aynı hafızaya erişebilir.  
Bir işlemci hafızada diğer işlemci tarafından  
değiştirildiği bilgileri tutar.

a) UMA Uniform Memory Access: Tüm CPU'lar aynı hafızaya erişim yapar.  
Cache coherency: Eğer bir işlemci paylaşımlı hafızaya bir  
update ettiğinde diğer işlemci haberdar olur.



- b) NUMA Non Uniform Memory Access: 4 farklı makine her makine  
kendi hafızasına erişir. Ama bu 4 makineye  
bir bus ile bağlanırlar. (4 laptopu düşün).



Böyle yaparak  
her CPU  
miktarını  
arttırıyoruz.  
Oysa  
her makine  
kendi hafızasına  
erişim  
yapabiliyor.  
Ayrıca diğer  
makineye de erişim  
yapabiliyor.

CPU diğer hafızalara erişimler (Access)

Cache coherency:

Her değışikliđi CPU'ya bildirilmek zorundur NUMA'ya çok zor  
ve masraflıdır. 16 CPU her CPU'dan cache coherency'i  
söylendiğinde durumda kullanılması için bir değışikli-  
ğin başka verileri etkilemediđi durumlarda.

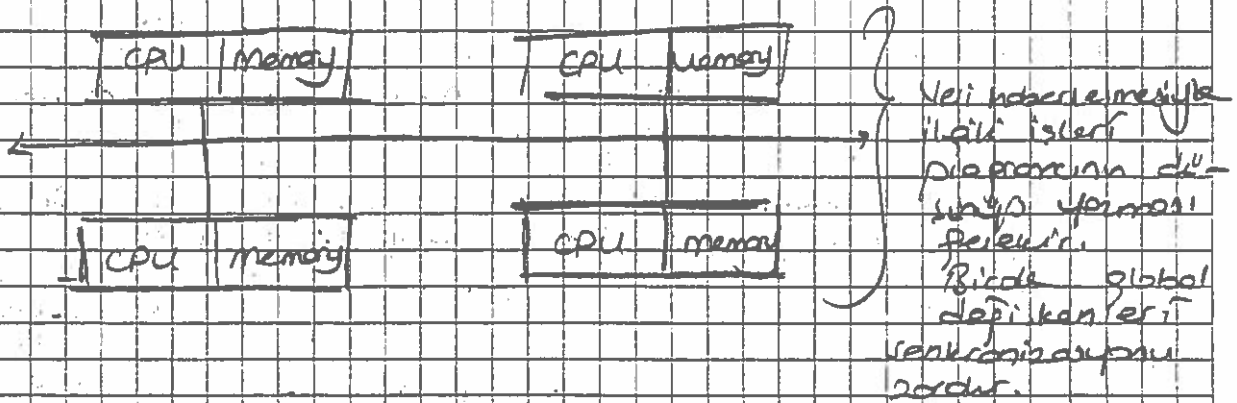
Hafızanın adres uzaylarını kontrol etmek kolaydır, fakat  
4 milyon değışiken containerın kontrol etmek zordur.  
Makine sayısı 4 is karmaşıktır.

## ② Distributed Memory

Her processor kendi local hafızasını kullanabilir Global  
adres uzayı yoktur. Bir local değışikliđi rapor etme  
bildirme durumu olur.

Fakat bir işlemi diğerinin bilgisine erişmek is-  
teğinde bu işlemiyle diğer hafızaya gidilir.  
internet veya başka iletişim ile işlem olabilir.

Avantajı: Ciddi olarak erişilebilir. İşleme  
sayısı arttıkça hesap hızı artar.  
Bus hızı hesap hızını belirler. Ucu ucudur.



Bir global değışikendeki değışikliđi diğerlerine bildir-  
mek zahmetlidir.  
Cache coherency yoktur.

Hibrit Distributed. Deđişik hafıza kullanarak çok hızlı  
işlem yapılabilir.  
2 hafıza mimarisinin olumlu yönlerini almış.

Meltem  
Jektir

SORU 1) Paralel Hesaplama nedir? Hangi amaçlar için yapılır?

CEVAP 1)

Birbirinden kısmen ya da tamamen bağımsız işlemlerin (process) birbirine paralel (koşut) biçimde çalıştırılması prensibine dayanır.

- Para ve zamandan tasarruf
- Daha büyük problemleri çözmek için
- Eş zamanlı çözüm üretebilmek için
- Yerel olmayan kaynakların kullanılması için
- Seri hesaplamanın sınırlarından dolayı

Paralel hesaplama yöntemleri

MPE modeli: Kaynaklar birbirleriyle mesajlaşarak paralel hesaplama yapılır. En çok kullanılanıdır.

Data parallel model: Data'lar parçalanır, sonuçlar birleştirilir.

SORU 2) a) Michael J. Flynn çok işlemcili bilgisayarları hangi kriterlere göre sınıflar? b) Flynn's Classical Taxonomy'ye göre mimari sınıflarını ve kısaca ne anlama geldiklerini yazınız.

Örneğin: Elimde büyük bir fotoğraf olsun. Bu resmin boyutunu 100 arttırmak için bunu parçalara böleriz, işlemler yapılır, sonuçlar birleştirilir.

Sei matrix toplamı kudu varsa, kod yazılırken paralel olarak hesaplanacağı belirtilir.

CEVAP 2)

a) Instruction ve Data

b) SISD : Single Instruction, Single Data → Eş zamanlı olarak aynı kısıttır, sei matrix'dir, sonuç tahmin edilebilir.

SIMD : Single Instruction, Multiple Data → farklı data'lar farklı CPU'lara gönderilip aynı işlem yapılır. Hızlı iş yapılır.

MISD : Multiple Instruction, Single Data → Bir işlemi farklı fiziksel birimlerde gerçekleştirip sonuçta bir işlemin yapılması.

MIMD : Multiple Instruction, Multiple Data → Doğru kullanım bilg. içerir. Aynı anda birçok program çalıştırılır. multitasking

SORU 3) Dağıtık Hesaplama kavramını kısaca açıklayınız. Ne amaçla gerçekleştirilir?

CEVAP 3)

Dağıtık sistemler fiziksel olarak dağılmış bilgisayarların bir araya getirilmesi ve koordine şekilde çalıştırılmalarıyla elde edilirler.

(Tanenbaum'a göre bir Dağıtık Sistem "Kullanıcısına tek bir bilgisayar gibi görünen ve bağımsız bilgisayarlardan meydana gelen sistemdir". Dolayısıyla bağımsız bilgisayarların tek ve tutarlı bir bilgisayarmış gibi çalıştırılmasını sağlayan sistemlerdir.

- Uzak kaynakların dağıtık şekilde kullanılabilmesi ve kolay erişim için
- Ölçeklenebilirlik için
- Sistemi oluşturan parçaların ayrıntılarından ve oluşabilecek hatalardan kullanıcının en az etkilenebilmesi için.
- İşbirliğini güçlendirebilmek için
- Daha emniyetli çalışma için
- Mobility

SORU 4) Merkezi (Centralized) ve merkezi olmayan (Decentralized) dağıtık mimarilerin özelliklerini kısaca açıklayıp bu mimarilere örnek veriniz.

CEVAP 4)

Merkezi mimariler klasik server-client mantığıyla çalışırlar. Frameworkler, çok katmanlı mimariler, web sunucuları, veritabanı sunucuları vb. request/reply mantığıyla çalışan sistemlerdir. Geleneksel örneği üç katmanlı modeldir. Bu modelde User/Interface/Processing layer/Data layer katmanları bulunur. Sıkı bağlı sistemler olduğundan ölçeklenebilirlik açısından sıkıntı yaşanabilir.

Merkezi olmayan mimariler genelde peer-to-peer denen sistemlerdir. Mantık olarak client veya server olmayan her iki rolü de üstlenen eş değerde bilgisayarlardan oluşurlar. Structured/Unstructured ve Hybrid P2P modelleri vardır. Bittorrent ağları başarılı örneklerdir.

Başarılar Dilerim  
Dr. Galip AYDIN



in için kısa sürede yapılması için 2 yöntem vardır.  
1. Kodun optimizeyenunun yapılması 2. Paralel sistemler

Paralel hesaplamalarda Bölme ve BİLEŞTİRME işlemlerinde zaman kaybı vardır.

Avantaj - Dezavantaj

- \* Aynı makina üzerinde hesaplama yapıldığında, aynı hafıza kullandığı için bu dezavantajdır.
- \* Paylaşılmış hafızalarda kullanıcının kolaylıkla elde edebildiği sonuçları elde edebiliriz.
- \* Farklı bilgi üzerinde aynı işlemi gerçekleştirir. Bu da işlemi yavaşlatır.
- \* Paralel hesaplama genelde bilimsel alanlarda kullanılır, ticari işlerde kullanılmaz.

##

Cloud: Elimizdeki kaynakları servis olarak kullanmaktır. Bir server sonrakılarınla çok server  
ide edilir. API: Servis çağırma ve durdurma işlemlerini gerçekleştirir. Ör/Amazon  
→ AWS: Donanım hizmet olarak kullanmayı sağlar.

MapReduce:

- Hadoop = Veri işleme, paralel olarak yapılmasıdır.
- Mahout = Makine öğrenme

Web Servisler: WSDL: Bir servis ve hizmetin ne olduğunu tanımlar. Bir XML dosyasıdır.  
SOAP: Başka bilgi. nelerle ilgili sorular.  
WDBZ: Web servis database.

Servislerin beraber çalıştığı sistemlere hizmet odaklı mimari (SOA) denir.

NOSQL:  
MongoDB  
Cassandra  
HBase  
Redis

İddiare katmanı: Çok sayıda makine üzerinde çalışan programları kullanıcıya  
göstermede tek bir sistem varmış gibi gösterir.

- \* Publish-Subscribe → orta katman yazılımının kullandığı yöntemlerden birisi.
- \* Yöntemde bir server var ve bu server kendisine ait kişilere mesaj gönderebiliyor.

Kline hesaplama nedir? Neden gerçekleştirilir? İçerik yazın. Dezavantajları  
yazın.

Tagora Hesaplama kavramını açıklayın. Neden geliştirilmiştir?

## Paralel Hesaplama Özellikleri

- 1- İş daha kısa sürede yapılır.
- 2- Az sayıda kaynak kullanarak, çok sayıda işlem yapılır.
- 3- Paralel hesaplama paralel sistemlerin omacını yerine getirmek için kullanılır.
- 4- Paralel hesaplama, kaç tane CPU varsa hepsi kullanılabilir.
- 5- Genelde bilimsel alanlarda kullanılır, ticari alanlarda kullanılmaz.

## Paylaşılmış hafıza - Dağıtık Hafıza

\* Paylaşımlı hafızada, hafıza üzerinde bir değişiklik yapıldığında tüm kullanıcılar aynı hafızayı kullandığında sonuçları hızlı bir şekilde görülebilirler.

\* Dağıtık hafızada, farklı makinelerde haberleşme network üzerinden haberleşmeyle olacağından yavaş olur.

# Paralel hesaplamanın Dağıtık Sistemler üzerinde kullanılması 2 aşamadan oluşur  
→ map aşaması  
→ reduce aşaması

## Dağıtık Sistemden İstenenler

- 1) Üzerindeki çokluğu kullanımdan gizlenmeli.
- 2) Ölçeklenebilir olmalı.
- 3) Sürekli çalışabilir olmalı.
- 4) Yeni kaynakları hizmete alırken eski kaynakların kesilmemesi gerekir.
- 5) Katmanlı mimari kullanılır. Herbir katmanda farklı işten sorumlu makineler vardır.

⇒ middleware (Orta katman): Çok sayıda makineyi tek bir sistem gibi görmeyi sağlar.

publish-subscribe: Orta katman yapısının kullandığı yöntemlerden biridir.

## Dağıtık Sistemler; (hedefler)

- Kaynaklar kullanıma açılır. (Grid)
- Kaynakların network üzerindeki dağıtıklığı gizlenmeli.
- Sistemin açık olması (Openness)
- Ölçeklenebilir olması



## Dağıtık Sistemin Güvenliliği

(2)

1. Access 2. Location 3. Migration 4. Concurrency (Tutulu) 5. Failure

### ÖLÇEKLENEBİLİRLİK

İhtiyaç duyulunca sistemin otomatik olarak genişleyebilmesidir. Çok sayıda istek aynı anda cevap verebilmesidir.

Cesitleri:

1.) Size scalability: Bir server kaç tane istek karşılar? Server'lar farklı yerlere durursa cevap hızına olur? gibi soruları yanıtlayan çeşitler.

2.) Administrative Scalability: Sistem ölçeklenebilirliği destekliyor mu? Sorunun köşklüdür.

İşletimsel ölçeklenebilirliği internet zaten yapıyor.

3.) Geographical scalability: Örneğin Gittigidiyor sitesinin web server'ları İstanbul ve Ankara'da bulunsun, farklı bir bölgeden istek gelince sorun olur mu?, arusuna cevap verir.

■ Ölçeklenebilirlik 3 noktada sorun yapıyor:

1.) Networkten kaynaklı haberleşme gecikmeleri

internet sitesinden cevap beklerken başka bir iş yapılabilmeli.  
gönderdiğimiz verinin boyutunu küçültmek (Javascript'ın tamamını post etmek yerine gerekli kısmı post etmek (Ajax))  
Server tarafında her istek için thread oluşturmak.

2.) Yükün dağıtılması

Yük = Hesaplama, veri

→ Hesaplama yükünü dağıtmak için farklı server'larda çalıştırılarak paralel programlara yapılır.

→ Veri yükünü dağıtmak için farklı mekanizmalar bulunur. Örneğin DNS. Her bölgede farklıdır ve isteklere karşılık verir.

3.) Verinin Çoklanması (Replication)

Veritabanı kopyasının aynı anda birden fazla yerde tutulmasıdır. Verinin çoklanması ölçeklenebilirliği hızlandırır.

Verinin kopyalanması web sitelerinde mirrored web ile yapılır.  
mirrored web = web sayfaları birbirine tutulur.

Ölçeklenebilirlik problemleri

1.) Tutarsızlık

2.) Çok büyük sistemler için ölçeklenmesi

3.) Global senkronizasyonu ölçeklenmede sorun yapması

## Dağıtık Sistemler Yapılandırılırken yapılan hatalar :

(3)

- Network homojen değildir.
- Network güvenli değildir.
- Topoloji her an değişebilir.
- Bant genişliği sınırsız değildir.
- Gecikme '0' değildir.

## Dağıtık Sistem türleri

- 1.) Dağıtık Hesaplama Sistemleri
- 2.) Dağıtık Bilgi Sistemleri
- 3.) Dağıtık Pervasive (Yaygın) Sistemleri

### Dağıtık Hesaplama Sistemleri

- Cluster hesaplama → Grid Hesaplama
- Cloud Hesaplama → Hizmet odaklı mimariler

Cluster Hesaplama : Birden fazla bilgisayarın network üzerinden birbirine bağlanarak toplam kapasitenin artırılmasıdır. Beowulf ve MPI paralel (kernel) hesaplama dağılımlarındandır.

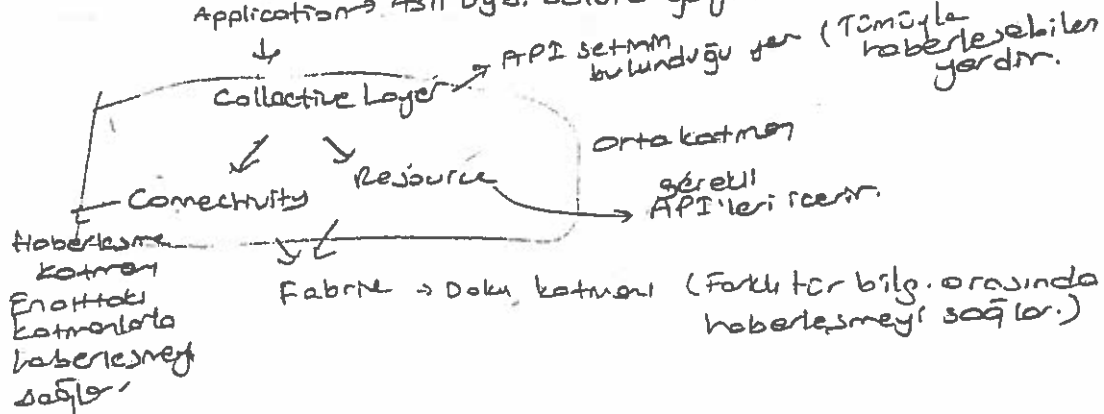
değerlendirir

- Ortak Bellek yoktur.
- Bazı uygulamalar için bellek erişim hızı yavaşlatarak azalır.
- İletişim hızı, bellek olumu/yatımı hızına göre yavaşlar.

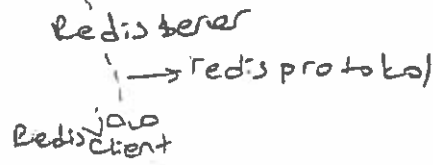


Grid Hesaplama : Clusterden daha büyük, heterojen sistemlerde kullanılır. Amaç; global bir işletim sistemi sağlamaktır.

Applications Asıl uyg. bulunduğu yerdir.



Redis (Replicated) Server



Redis bir NOSQL türüdür. Key value şeklinde çalışır.

Redisin redis.io'da kurulumunu yaptık.

Terminal üzerinden

src/redis-server yaparak

farklı terminalden redis-cli yaparak client oluşturduk.

Bilgisayarları bridge yapısını kullanarak sanal ip'leri bağladık.

Uçucu client ve pom yapılandırma dosyalarını oluşturduk.

pom.xml

Cache mekanizması : En hızlı erişilebilecek bellek alanına gelin. Server tarafında dinamik web sayfalarının cache'leri tutulur. (2)  
↳ JSP, ASP, PHP ...

Disoridan gelen istekler önce cache'e gönderilir eğer orda varsa direkt cevap döndürülür. Cache mekanizmasının en çok kullanıldığı yer web haritalarıdır. server'da bir değişiklik yapıldığında bunlar cache yenisiilir.

Image Proxy : En tepede bir resim varsa, altına bir çok resim vardır. Bunu file system'de tutulur.

Replication ve Cache'de en önemli sorun tutarsızlıktır.

Dağıtık Sistemler yapılandırılarak güvenliyle yapılan hatalar:

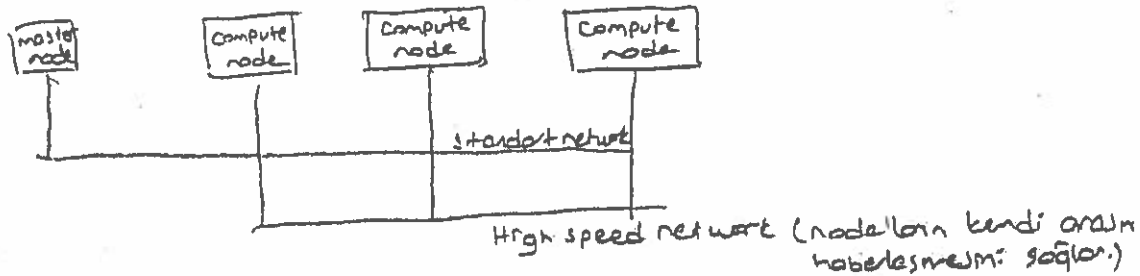
1. Network güvenli değildir.
2. Homojen değildir, çok sayıda server vardır.
3. Topoloji her an değişebilir.
4. Bant genişliği sınırlı değildir.
5. Gecikme sıfır değildir.

Dağıtık Sistem Türleri:

- Dağıtık Hesaplama Sistemleri (Distributed Computing Systems)
  - Dağıtık Bilgi Sistemleri
  - Dağıtık Pervasive (Yaygın) Syst.

Dağıtık Sistemlerde hesaplama yöntemleri:

1) Cluster Computing (Küme hesaplama) = Büyük bir problem dünün geniş çapta bir simülasyon yapacak çok büyük super bilgisayarlarla yapılabilir. Baska bir kelime küme hesaplama, MPI paralel programlamada sıkça kullanılır. Bir bilgisayar master seçilir. Bunlar kullanıcı ile iletişimi sağlar.



## İşleri dağıtık sistem kullanıyor?

- 1.) Hesaplama kaynaklarını çok kolay bir şekilde erişilebilir hale getirir.
- 2.) Sistem kaynakları ortak kullanılır. Örneğin: Seti@home ve grid
- 3.) Sistemin açık olması
- 4.) Ölçeklenebilir olması

## Dağıtık sistemden beklenenler

- 1.) Kaynak paylaşımı arttıkça güvenlik sorunu ortaya çıkar
  - 2.) Güçlü sistemler elde etmek için birçok kaynağın kullanılması gerekir. Kaynakların nerede olduğunu kullanıcılar bilmez, sadece yönetici bilir. Yani dağıtıklık gizlenmesi gerekir.
  - 3.) Erişim ile ilgili bazı noktaların gizli olması gerekir. Verinin başka tarafta hangi formatta bulunduğu gizli olması gerekir.
  - 4.) Serverların çalışıp çalışmadığını sürekli olarak kontrol edilmesi gerekir.
  - 5.) Bazen veritabanının başka bir yere taşınması gerekir. Buna migration denir. Bu taşıma işleminin gizlenmesi gerekir.
- Kısacası dağıtık gizliliği:
- Erişim, Location, migration, Tutarlılık, Failure
- Bir foto oluşturduğunda kulağın nerede gizlendiğini düşün.

3.) Çok sayıda kaynağı bir araya getirdiğimizde, çok sayıda platformu bir araya getirmiş oluyoruz. Bu sistemlerin birbiriyle haberleşmesi gerekir. Buda 3 şekilde gerçekleşir. a) Hardware b) Platform c) Language

4.) Çok sayıda işteği aynı anda cevap verebilmelidir. Ölçeklenebilirlik.

## Ölçeklenebilirlik çeşitleri

- Size scalability
- Geographical scalability
- Administrative scalability

## Ölçeklenebilirliği etkileyen problemler

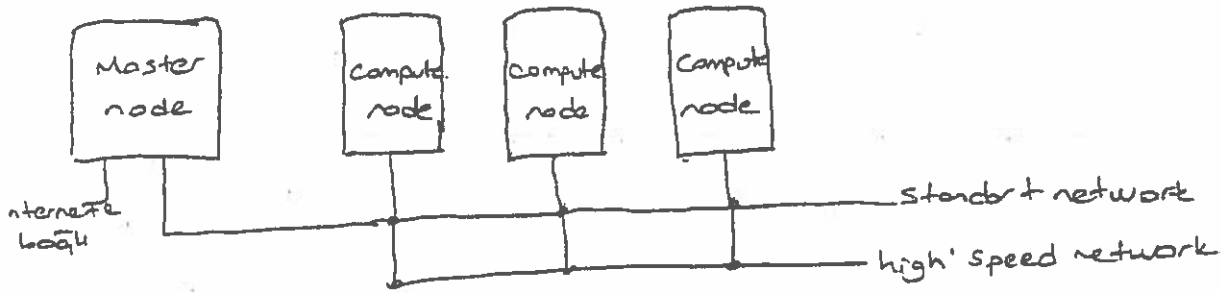
- a) Haberleşmedeki gecikmeler: Gecikmeyi önlemek için verinin boyutu küçük olmalı ya da kullanıcı başka sayfaya da yönlendirilmeli. AJAX kullanılır data hızlı olsun diye.
- b) Yükü dağıtmak: Yükü hesaplama ve veri. Hesaplama dağıtmak için farklı serverlarda çalıştırarak paralel programları gerçekleştirerek yapabiliriz. Veriyi bölmek için farklı mekanizmalarda bulunur.

c) Replication (Verinin Çoğaltılması): Örneğin, bir veritabanı'nın kopyası olsun. Bu kopya istedikler alıyor. Her yerde verilerin kopyalarını tutmamız gerekir. Bu kopyalama işlemi replication denir. Bu işlem mirror (ayna) yöntemiyle yapılır.

1

## Cluster Hesaplama:

- Birden fazla bilgisayarın birbirine network üzerinden bağlanarak, bilgisayarların toplam kapasitesinin artırılmasıdır.
- Depolama, işletim sistemi, uygulama program arayüzü, uygulamalar gibi değişik sistem katmanlarında gerçekleştirilebilir.
- Cluster homojendir. Yani bilgisayarların hepsi aynı türden aynı hızda çalışır.
- Beowulf ve MPI küme hesaplama yapımlarındandır.
- MPI'de bir bilgisayar master olarak seçilir, master bilgisayara bir işlem geldiğinde bölüp diğer bilgisayarlara gönderir.
- Küme hesaplama dezavantajları:
  - Ortak bellek yoktur.
  - Ölçeklenebilirliği yoklamak bazı uygulamalar için zordur.
  - İletişim hızı, bellek okuma/yazma hızına göre yavaştır.



- Master node'un görevli işleri paylaştırma ve toplamadır.
- Node'ların kendi hafızaları üzerinde işlerini yapmaları çok hızlıdır. Fakat 2 node haberleşmesinde iş yavaşlığına karşın high speed network kullanılır.

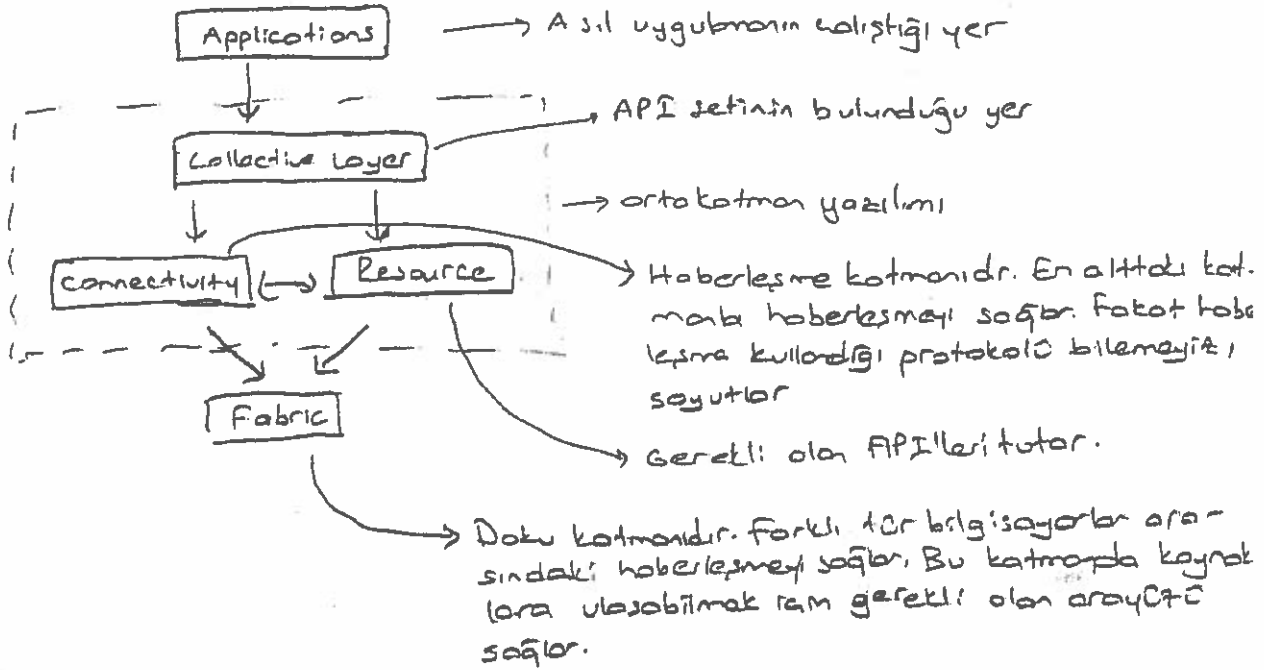
## CLOUD HESAPLAMA

- Server'in özelliğini bilmeden istediğimiz kaynağı aldığımız mimari türdür.
  - IaaS → Hizmet olarak Altyapı
  - PaaS → Hizmet olarak Platform
  - SaaS → Hizmet olarak yazılım.

## Grid Hesaplama

(2) . . .

- Cluster'lerden daha büyük, heterojen sistemlerde grid hesaplama kullanılır.
- Grid hesaplama yapılarak bilgisayarın özellikleri, işletim sistemleri birbirinden farklı olabilir.
- Grid Hesaplama amaç : Global bir işletim sistemi oluşturmak ve çok büyük hesaplamaları dünyanın her yerine yayılmış kaynakları kullanarak yapmaktır.



## HİZMET ODAKLI MİMARİLER

- SOA : Hizmet odaklı mimari örneğidir.
- Hizmet odaklı mimariler, dağıtık sistemlerin en başarılı ve en yaygın kullanılan mimarisidir.
- Birçok hizmet birleştirilerek büyük bir hizmet elde edilir.
- Web servislerden oluşur.

## Web Servislerin Çalışma Yapısı

- Servis sağlayıcısı web servisini oluşturur ve bu servisi bir web sunucusu üzerine koyar.
- Oluşturulan servisi bir UDDI kaydı olarak tanımlayabilmek WSDL teknolojisi kullanılarak bir adet WSDL uzantılı dosya hazırlanır. (Bu işlem IDE tarafından otomatik olarak gerçekleştirilir.)
- Servis sağlayıcısı servisi bir UDDI kaydı olarak kaydeder.
- İstemcinin başka bir web servisi olabilir, UDDI sorgulamasını gerçekleştirerek servisin bir özetini kendi tarafında oluşturur.
- İstemci tarafı oluşan bu servisi kullanarak karşı taraftaki web servise bir istek gönderir. Bu istek SOAP mesajları aracılığıyla gerçekleştirilir.
- Karşı taraf istenen bilgiyi SOAP aracılığıyla XML formatında geri döndürür.
- Alıcı taraf SOAP protokolüne kullanarak gelen XML'deki bilgiyi neye dönüştürür ve sonucu uygulamaya ramde istediği şekilde kullanır.

Hesap.java

Source package içinde

```
@WebService()
public class Hesap {
    public Integer toplamaYap(
        @WebParam(name="deger1") int deger1,
        @WebParam(name="deger2") int deger2) {
        return null;
    }
}
```

(Web Services içinde)

```
@WebMethod(operationName="toplamaYap")
public Integer toplamaYap(@WebParam(name="deger1")
    int deger1, @WebParam(name="deger2") int deger2) {
    return deger1 + deger2;
}
```



### Web Servis Özellikleri:

- HTTP protokolü üzerinden iletişim kurulur.
- XML web servisleri SOAP üzerinden iletişim kurar.
- Ağız protokollerini kullanarak haberleşirler.
- Diğer uygulamalar tarafından kullanılabilirler.
- XML tabanlı bir teknolojidir.
- SOAP basit ve mesaj tabanlı bir iletişim sağlar.

SOAP  
UDDI

## WEB SERVISLER

(3)

• Farklı bilgisayarlarda paralel Programlama yapmak için birbirleriyle habertmelerini sağlar.

RMI: Başka bir bilgisayardan bir programı, bir metodu çalıştırmak için geliştirilen ilk teknolojidir.



Proxy 2, Proxy 1'e anlayacağı şekilde mesaj gönderir. Topla fonksiyonu ile a ve b parametrelerini gönderir. Bu mesaj bir network mesajıdır TCP ya da UDP protokolüne kullanılır. Proxy 1 gelen mesajı çıkarıp Java nesnesine dönüştürür, PC1'e verir. Hesaplamayı yapıp sonra Proxy 1 üzerinden PC2'ye network mesajı olarak gönderilir.

mimarî karmaşık olursa RMI'yi burada kullanamayız.

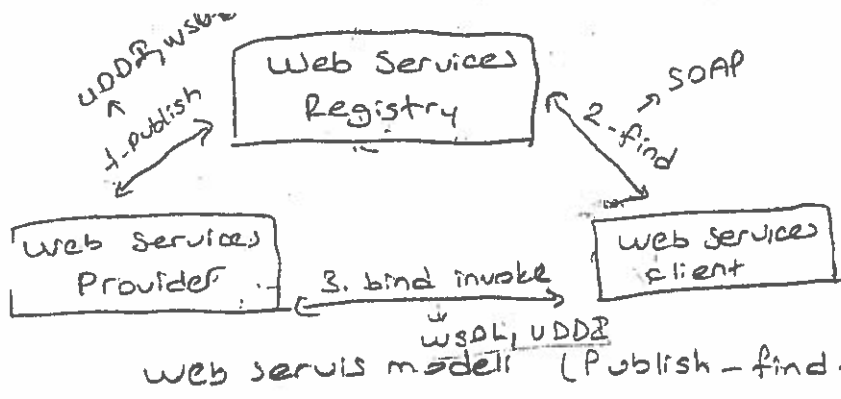
Web servis: Esnek, herkesin kullanabileceği, platformdan ve programlama dilinden bağımsız bir şekilde kullanılacak mimardır.

RPC, ORPC (DCOM, CORBA, JAVA RMI)

- Bunlar beraber çalışmazlar.
- Çoğunda firewall aşılamıyor.
- DCOM Windows'a, RMI Java'ya çok bağlı.
- Karmaşıktır.

Web Servisi standartları

- XML direkt gönderilmez, SOAP kullanır. SOAP üzerine XML mesajımları da ekler öyle yollar. Mesajlaşma HTTP → XML → SOAP şeklindedir.
- Web servisleri tanımlamak için WSDL dili geliştirilmiştir.
- Web servislerin keşfedilmesi için kullanılan register UDDI'dir. Web servisler kayıt ettirilmir, arayıp bulduğuna depodur.
- SOAP bir mesajdır, aynı XML tipinde mesaj barındırır. Bu mesajın ne tipte olduğunu WSDL belirler.



(4)

### Küme özellikleri

- Yüksek kullanılabilirlik
- Yük Dengeleme
- Yüksek Güvenilirlik
- Veritabanı

### Tanımlamalar

#### XML (Extensible Markup Language)

- web servislerinde veri tanımlamak için kullandığımız standarttır.
- HTML'de tag'lar bellidir, XML'de biz oluştururuz.
- XML oluşturma amacı; standart tipteki veriyi bir yerden diğer yere iletmektir.
- Uluslararası bir standarttır.

YBT  
YD  
YK

#### UDDI (Universal description discovery and integration) :

- web servisleri hakkında bilgi yayınlamak ve almak için kullanılan XML tabanlı bir teknolojidir.
- Web servislerin keşfedilmesi için kullanılan register'dir.

#### WSDL (Web Service description Language) :

- web servislerini, servislerin içerdikleri metodları, bu metodların da işlemlerini tanımlanarak kullanılır.
- WSDL dosyasında çeşitli UDDI girdilerine ve SOAP mesajlarına referanslar bulunmaktadır. Kendi taraf WSDL dosyasını aldığı anda bu referansları inceleyerek web servisinin içeriği hakkında tam bilgiye ulaşabilir.

#### SOAP (Simple Object Access Protocol) :

- web servisi ile web servisi istemcisi arasındaki haberleşmeyi gerçekleştiren XML tabanlı protokoldür.
- Bir mesajdır içinde XML tipinde bir mesaj barındırır.

#### SOA (Service Oriented Architecture) :

- Servis yönelimli mimari anlayışa gelen bu yaklaşım web servis tabanlı uygulamaların çalışma mantığını açıklar.
- Uygulamalar servislere halinde dağıtık bir şekilde yayınlanır. Kullanıcının girdiği taraftaki yazılımlar ise bu servislerden ihtiyaç duydukları bilgileri ve kullanıcılara çıktı verir.

• tekrar kullanılabilirlik  
• bağımlılığı azaltma  
• basitleştirme  
• sağlanabilirlik

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

# DAĞITIK SİSTEMLER

## Paralel Hesaplama

Dağıtık sistemler fiziksel olarak farklı alanlarda bulunan hesaplama araçlarının internet ile birbirine bağlandığı sistemlerdir.

Paralel hesaplama, işlemlerin farklı bilgisayarlarda gerçekleştirilip tek bir bilgisayara toplanmasıdır. Paralel hesaplama için program paralel yazılması gerekir. Paralel program kodu programın çalıştırılacağı her bilgisayara eklenmelidir. Kodlarda özel şekiller için farklılık olabilir.

Paralel hesaplamanın temel özelliği ilk başlarda farklı bilgisayarlarda yapmakken daha sonra aynı makinede birden fazla CPU'nun olduğu durumlarda gerçekleştirilmiştir. Bir birden fazla CPU'ya 1 hafıza koyup hafıza ortak olarak kullanılır. (Paylaşımlı hafıza) Paylaşımlı hafızanın avantajı:

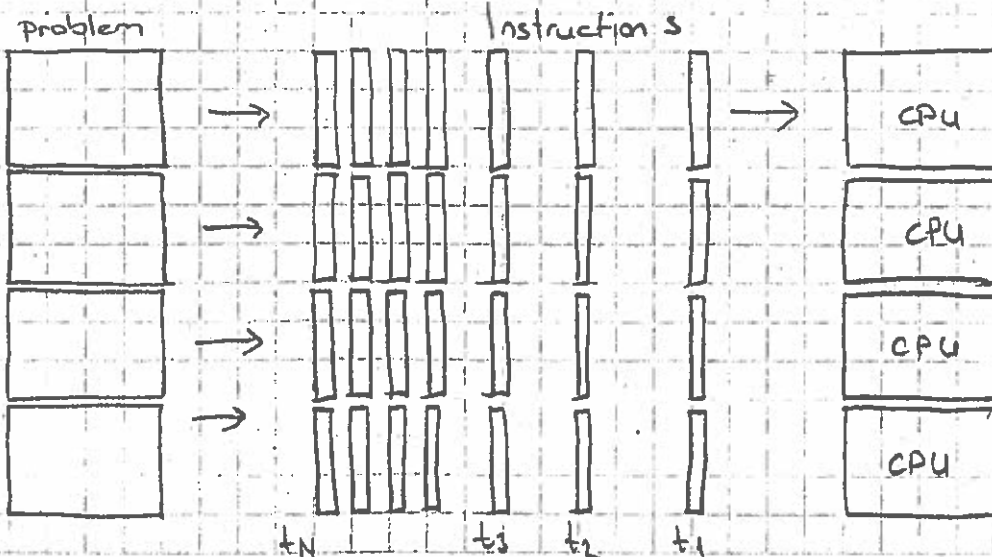
- Verinin parçalanmasına gerek kalmamıştır.
- Bilgisayarlar arasındaki haberleşmeden dolayı zaman kaybı engellenmiştir.

Paralel hesaplama elimizdeki bir problemi birden çok parçaya ayırarak daha hızlı çözebilme amaçlar. Ancak bazı problemler vardır:

- Birden fazla CPU olmalıdır. (Birden fazla komut çalıştırma)
- Bir problemin paralel olarak çözülebilmesi için eş zamanlı olmalıdır.

"Bir problemi elimizdeki sınırlı kaynaklarla en hızlı şekilde nasıl çözebiliriz?" sorusuna cevap olarak paralel hesaplama.

Parçalandığı zaman birbirini etkilemeyecek problemler olmalı. Komutlar eş zamanlı çalışmalı (Birden fazla CPU'ya ihtiyaç var)





## Paralel Hesaplama Kullanım Alanları:

- Atmosfer
- Fizik (Kuantum fizik, atom fizik)
- Biyoteknoloji, Genetik
- Kimya
- Jeoloji, Sismoloji
- Bilgisayar Bilimleri
- Database, data mining
- Web search

## Kullanım Sabepleri:

Temel sebebi daha kısa sürede bir problemi çözebilmek.

- Zaman ve para tasarrufu
- Büyük problemleri çözebilmek
- İş zamanlı (concurrency) çalışma imkanı sağlar. (accessgrid.com)
- Ede olmayan kaynakları kullanabilmek.
- Seri hesaplamalarda işlemlerin limitinin olması. Seri hesaplamada kullanılmamışın nedenleri:
  - \* Fazla bilgisayar yok
  - \* İşlem hızının sınırlı olması
  - \* CPU'da bir hesaplama limitinin olması. Herddisk'in sınırlı hızı (rpm) var. Veri transfer hızını, hesaplama hızı kısıtlar.
  - \* Ekonomik kısıtlama

## Nasıl Yapılır?

- Multiple execution units (Çoklu hesaplama üniteleri - birden fazla CPU)
- Pipelined komutlar
- Multi-core

## Paralel Hesaplama Sınıflandırma

Flynn's tarafından yapılmıştır. Bu sınıflandırmanın temelinde;

- Data
- Instruction

### Flynn's Classical Taxonomy

- SISD (Single Instruction, Single Data) → Tek CPU
- SIMD (Single Instruction, Multiple Data) → Birden fazla CPU
- MISD (Multiple Instruction, Single Data)
- MIMD (Multiple Instruction, Multiple Data) → CPU'lar bağımsız

## Single Instruction, Single Data

- Seri bilgisayarlar.
- Paralel hesaplamaların sonuçları belli.
- En eski, en çok kullanılan model.

CPU'ya bir veri alınır ve sadece bir sonuç döndürür.

## Single Instruction, Multiple Data

Bir den fazla CPU aynı anda bir komutu işler. Veriler farklıdır. Örneğin, bir fotoğrafın %10 renk düzeltme tek işlem ama birden fazla veri var. Ekran kartı işlemlerinde, GPU'lar da kullanılır. Yaygın değildir.

## Multiple Instruction, Single Data

Bir veri birden fazla CPU üzerinde işlenir. Örnek olarak; filtre kuma, ses filtreleme verilebilir.

## Multiple Instruction, Multiple Data (Super Computer)

En çok kullanılan, bilgisayarlarda var olan model. CPU'ların hepsi bağımsız çalışır. Seriklon veya asenkron çalışabilir. Networkte bağlı bilgisayarlar cluster yapılır. Bağımsız veriler üzerinde çalışılır.

Aynı bilgisayar üzerinde yapılırsa, Paralel Programlama,

Farklı bilgisayar üzerinde yapılırsa, Dağıtık Hesaplama adını alır.

SISD → Tek core'lu bilgisayarlar

SIMD → GPU'da kullanılır.

MIMD → Filtrelemelerde kullanılabilir. Su anda yolda.

MIMD → Su anda kullanılan bilgisayarlar. En çok kullanılan!

## Paralel Programlama Genel Terimler

Task: Hesaplama işinin birbirinden ayrılan parçaları, bağımsız çalışabilecek kısımlardır. Bir taskın diğer tasklarla ilişkisi yoksa paralel task (haberleşmesi gerekmiyorsa), diğer tasklarla haberleşmesi gerekiyorsa serial task dendir.

Pipelining:



## Shared Memory (Paylaşımlı Hafıza) :

Bilgisayar üzerindeki CPU'lar bütün işlemler için aynı hafızayı kullanır.

Dezavantajı : CPU'nun aynı hafızayı görüp üzerinde değişim yapabilmesi.

Bilgisayarın yapısını karmaşıktırması.

Avantajı : İşlemleri hızlandırması.

## Symmetric Multi-Processor (SMP) :

Birden fazla CPU'nun tek bir hafıza alanını kullanması ve tüm kaynaklara eşitli paylaşılan mimari.

## Distributed Memory :

Birden fazla board, her boardda bir CPU ve hafıza vardır. Bunların aynı amaçla kullanılmak istenmesi. Bunun için network gerekli. Normalde CPU kendi hafızasını görür, diğer hafızalarını göremez. Bu nedenle hafızayı diğerlerine açmak gerekir. Ek yazılım gerektirir.

## Communication (Haberleşme)

En fazla zaman kaybının olduğu yerdir. Network zaman kayıpları burada ortaya çıkar. En iyi durumda 1'in yarısında ve sonunda haberleşme gerekir. Paralel hesaplamaların haberleşme zamanı yükseltse yapılmaz.

## Senkronizasyon :

Birbirinden ayrı olarak çözülen problem parçalarının daha sonra birleştirilmesi gerekir. Bunun içinde parçaların birbirine senkronize edilmesi gerekir.

## Granularity (Hassasiyet)

Hesaplama içi bulunur. Eğer hesaplama > haberleşme ise

buna coarse denir. Hassasiyetin önemli olmadığı durumlarda kullanılır.

Fine : Hesaplama zamanı kısa aralarda haberleşme yapmak gerekir. Zaman çok habersizliğe girer.

## Observed Speedup :

Yaptığımız zaman tahmini ile gerçek zaman arasındaki farktır. Bir iş seride ne kadar sürede yapılır, paralelde ne kadar sürede yapılır.

## Paralel Overhead :

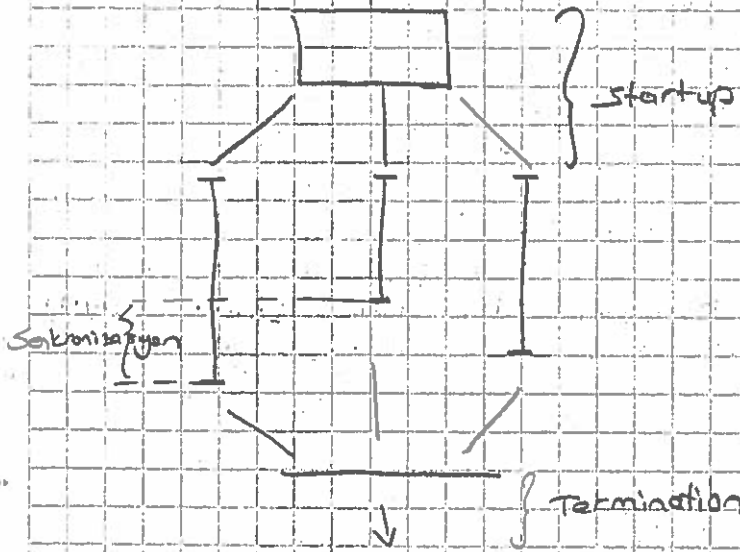
Ekstra yüklerdir Paralel Overhead'e neden olan faktörler.

→ Task start-up time : Başlatma zamanı uzun sürer. Değişkenlerin oluşturulması, değerlerinin atanmasından dolayı, ilk iterasyonda yüksek çıkar.

→ Synchronisation : 4 parçaya bölünen problemin 3'ü erken biter, 4'ü biterse 3'ü bekler. 4'ü biterse 3'ü bekler. 4'ü biterse 3'ü bekler.

→ Data communication

→ Task termination Time.



Programın ne kadarı paralelleştirilebilir? Türleri var.

## Massively Parallel :

Çok büyük, herşey paralel.

## Embarrassingly Parallel :

Problem sayısı az, haberleşme yok.

## Scalability (Ölçeklenebilirlik):

Web serverlerde ölçekleme en önemli şeydir. Bir program ölçeklenebilirse gelen işte göre kendini artırır. Ölçeklenebilirlik ihtiyaca göre büyütülmek. Örneğin; 3 makine 3 günde, 10 makine yarım saatte gözüyorsa ölçeklenebilirlik.

## Multi-core Processors:

## Cluster Computing:

Küme bilgisayarlar. Çok sayıda bilgisayar ama süper bilgisayar değil. Clusterla birbirine bağlayıp paralel hesaplama yapma. Bilgisayarları süper bilgisayar yapma.

## Super Computing / High Performance Computing:

Çok büyük kaynaklı pc. 100 petabyte RAM, 200.000 CPU gibi.

## PARALEL HAFİZA MİMARİLERİ:

### 1. Shared Memory:

Aynı makinede birden fazla CPU varsa tüm CPU'lar ortak hafızaya erişim sağlayabilirler. Yapılan değişiklik anında hepsi tarafından görülür.

Avantaj: Erişim süresi aynı, haberleşmeden dolayı zaman kaybı yok.

→ UMA

→ NUMA

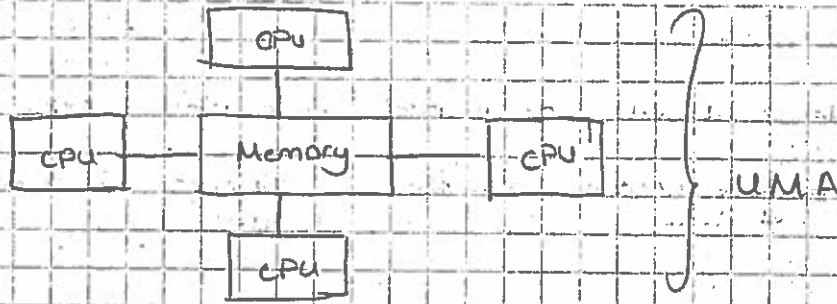
### → Uniform Memory Access (UMA)

Bütün CPU'ların özellikleri hafızaya uzaklığı (fiziksel olarak), erişim süreleri hepsi aynı.

Avantaj: Erişim süresi aynı olduğundan haberleşme kayıpları yoktur.

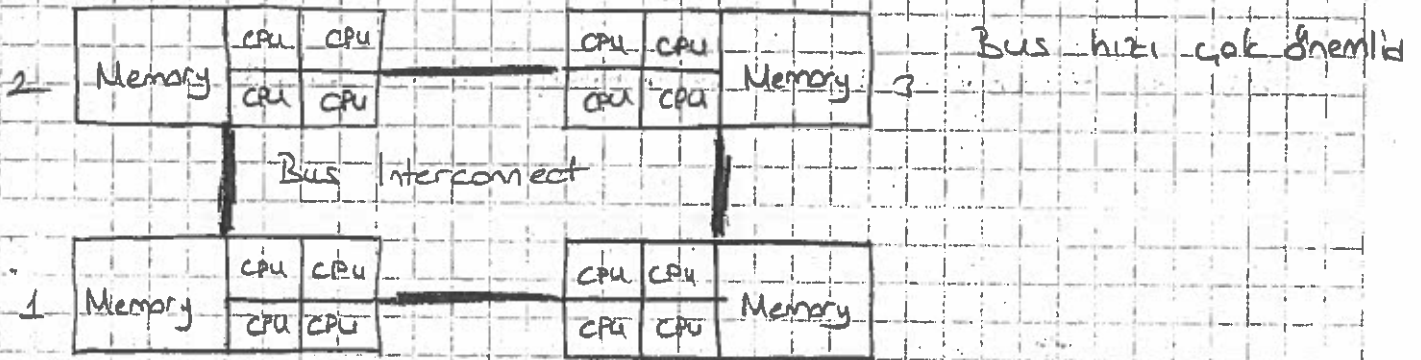
CC-UMA (Cache Coherent UMA) da denir. Cache coherent, cache tutarlılığıdır. Hafıza update edilirse hepsi anında haberdar olur.

CPU'lar tek hafıza kullanır. Simetrik Multi Processor da denirler (SMP).



## → Non-Uniform Memory Access (NUMA)

4 farklı makine, 4 tane CPU, her bir makinenin ayrı hafızası vardır ve bu hafızalar birbirinden ayrıdır. Hafızalar bir busla birbirine bağlanır.



Her makine birbirinden bağımsızdır, birbirlerine network ile bağlanırlar.

Her makine doğrudan diğer makinenin hafızasına erişebilir. Ancak 1. nolu pçnin 2. nolu pçye erişme süresi ile 3. nolu pçye erişme süresi aynı değildir. Yani erişim yavaşlığı vardır.

Tüm makineler diğer makinelerin hafızalarını adresleyebilir.

Burada cache coherent sağlanıyorsa CC-NUMA denir.

Cache coherent: Yanlış sonuçların önüne geçmek için gerekir. Donanımsal olarak herhangi bir update durumunda diğerlerine haber vermesi gerekir ya da belirli aralıklarda değerlerin alınması gerekir.

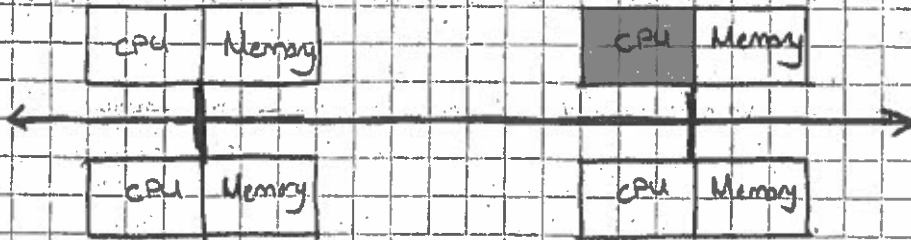
Avantaj: RAM ve hafıza arttı. Dolayısıyla çözülecek problem sayısı arttı. Global adres alanı sağlar.

Dezavantaj: Ölçeklenebilirliği zordur.

## 2. Distributed Memory:

İşlemciler sadece kendi local hafızasını adresleyebilir. Cache coherent yoktur.

Bir işlemci diğer işlemcinin adres uzayındaki bilgiye erişmek isterse işlemciye bildirir. (İstek bus'a çıkar.)



Ölçeklenebilir. Ucuz bir yapıdır. Kolaylıkla oluşturulabilir bir mimari.

### Avantaj:

→ Çok ciddi ölçeklenebilir. (Ancak makine ne kadar fazla olursa program karmaşıklaşır.)

→ Bus'ın hızı hesaplama hızını belirler.

→ Ucuz bir yapı.

→ Herkes kendi local hafızasını hızlı bir şekilde adresler.

### Dezavantaj:

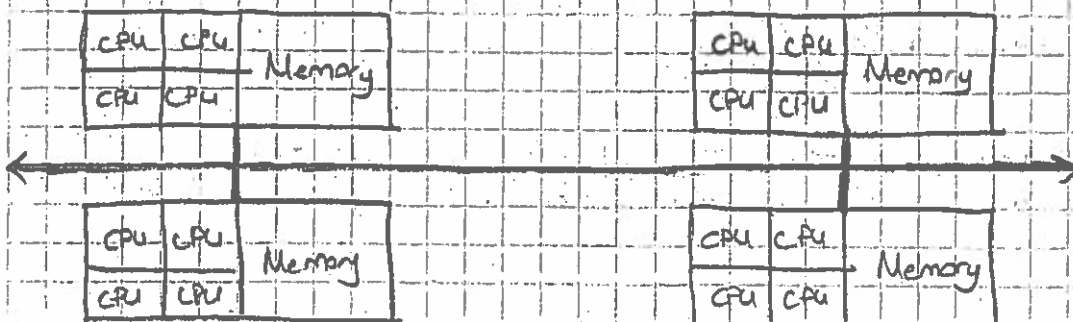
→ Programcı her şeyi düşünüp kendisi yapmak zorunda.

→ Global değişkenlerin senkronizasyonu zordur. Ne kadar global değişken varsa bus trafiği o kadar fazla olur.

→ Cache coherent yoktur.

## 3. Hybrid Distributed-Shared Memory

Ortak hafıza yok, herkes kendi hafızasını görüyor. İstek yapma lazım.





Avantaj → SMP kullanıldığında paralelleştirme yapılabilir.

→ Global değişkenleri senkronize edilmez. Gerekli global değişken için bus üzerinden istek gönderilir. Adres uzağı yönetimi yok.

Dezavantaj: Bus ne kadar hızlıysa paralel hesaplama o kadar hızlı olur.



Paylımlı hafıza modeli NUMA ile analarındaki fark, bu modelde diğer hafızadaki adreslerin bilinmesine gerek yoktur.

## PARALEL PROGRAMLAMA MODELLERİ

Paralel hesaplama için genelde C, C++ kullanılır. Ancak en iyi diye bir şey yoktur. Makine mimarisi ile programlama modeli birbiriyle uygun olmalı.

- Shared Memory
- Threads
- Message Passing
- Data Parallel
- Hybrid

### 1. Shared Memory :

Aynı adres uzağını paylaşırlar. Klasik programlamaya göre farklı hafızanın paylımlı olması. Problem parçalara bölünecek, her birini bir CPU işleyecek ve diğer CPU'larda bunu görmeli.

Semaför, Kilit → Paylaşılan veriyi korumak gerekir. (Aynı anda 2 CPU bir değişkenin değerini değümezin diye).

Çalışılan bilgisayar mimarilerine çok bğımlı. Global çözüm yok.

### 2. Thread Model :

Bir process'in birden fazla iş parçacığı vardır. Threadler processlerin global adres uzağını, global değişkenlerini kullanır. Kendine ait hafıza ve değişkenleri de vardır.

Mimariden bağımsızdır. Process içerisinde çalışır. Tek bir process gibi görünür.

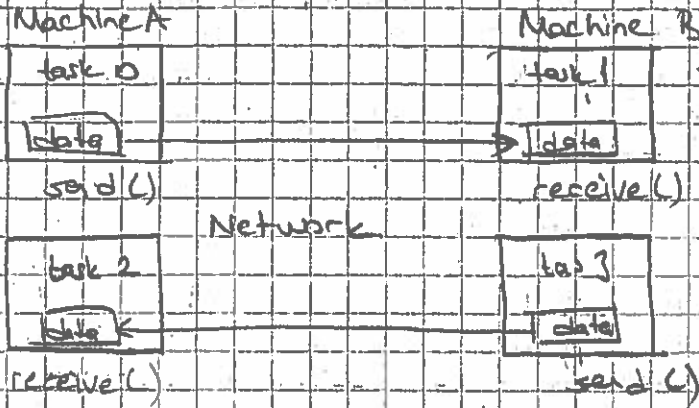
Avantaj → Threadler processlere göre daha hızlı oluşturulur. (Processler için kontrol blokları oluşturmak gerekir)

- Global adres üzerinden haberleştikleri için zaman kaybı yok
- Haberleşmeyi programcının düşünmesi gerekmez.

Threadlerin 2 farklı uygulaması : OPEN MP, Pthreads

### 3. Message Passing Modeli (Mesajlaşma)

Paralel hesaplama için en önemli, en çok kullanılan modeldir. Mesajlar bus üzerinden gönderilir. Mesajlar ne istendiği koyulur. İki tarafta da Message Passing ile çalışması gerekir. Karşıda alıcı olmalı.



Her zaman send() karılık receive() olmalı.

Bağımsız bilgisayarlar, kendi CPU, kendi hafızası vardır. Mesajları kullanarak haberleşiyorlar. CPU'ların local hafızaları var. İstenilen oranda gerçekleştiriliyor.

Mesaj ortalık bus varsa bus üzerinden yoksa network üzerinden yapılır. Paralel hesaplamanın en çok kullanılan modeldir. (sıradan bilgisayarlar üzerinde kurulabilir).

Programcı sadece paralelleştirmeyi yapar, mesajlaşma kütüphanesi ile MPI → En yaygın olanı)

### 4. Data Parallel Model

Paralel hesaplamanın en önemli kısmı verilerin eşzamanlı işlenmesi üzerinde gerçekleşir. Veriyi işleyecek programları paralelleştirilir. Veri parçalanarak işlem yapılıyorsa data parallel modeldir. Örneğin dizinin her elemanına 1 eklemeye gibi. Tek program var. Program döngüyü parçalıyor. Program yazarken döngünün yanına yazılır. Sistem onu CPU'lara dağıtır. (sistem-MPI)

Clusterda kaç makine varsa o sayıya bölür. İş otomatikleştirmek amaçtır. Paralel constructörler (parallel while, for)

Factor 90-95 → Paralel hesaplama da önemli yere sahip.

## PARALEL PROGRAMLAMA TASARIMI (3. Slayt)

Parallellendirme 2 şekilde olabilir

→ Otomatik

→ Manuel

→ Otomatikleştirme Problemin niteliklerinden dolayı önemli ölçüde gerçekleştirilemiyor. Kararı programcı verir. Programı paralelleştiren compiler'dır. Compiler kaynak kodu analiz edip paralelleştirilen yerleri otomatik olarak bulur. Bu işlem bir dizi şeklinde yapılır. Bu işlemden sonra paralelleştirmenin mantıklı olup olmasına bakılır. 10 milyon elemanlı olsa olursa daha sonra program derleyip çalıştırılır. İş karnajik değilse bu yöntem bazıları sağlar.

→ Manuelde paralelleştirilerek yerlere programcı karar verir. Programcı paralelleştirilecek yerlere kendi karar verir. Ne derece faydalı olacağını programcı hesaplamalıdır. Karar verirken programcı domaini iyi anlamalı. Parallellendirme uygun olur mu? Problemlerde buna bakılmalı.

Örneğin; bir maddelerin etrafındaki diğer maddelerle farklı bağ oluşturma ihtimalleri vardır. İşler birbirinden bağımsız ise %100 paralelleştirilebilir. (Paralel olur)

Örneğin; Fibonacci serisi hiçbir zaman paralelleştirilemez.

### Program ve Problemi Anlamak

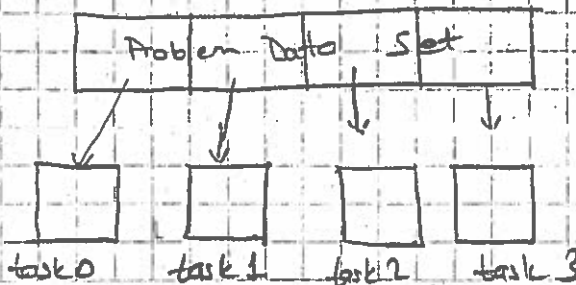
Hotspots: Programın profilini, temel noktalarını belirlemek. Profilers ve performans analizi yardımcı olur.

Bottlenecks: Programın yavaşladığı yerler. Bunlar paralel hesaplamam yapıldığı yerlerin altına atılmalı. Örn; I/O işlemleri programı yavaşlatır.

Partitioning: Programı ayırık parçalara bölmek 2 şekilde yapılabilir.

### → Domain Decomposition:

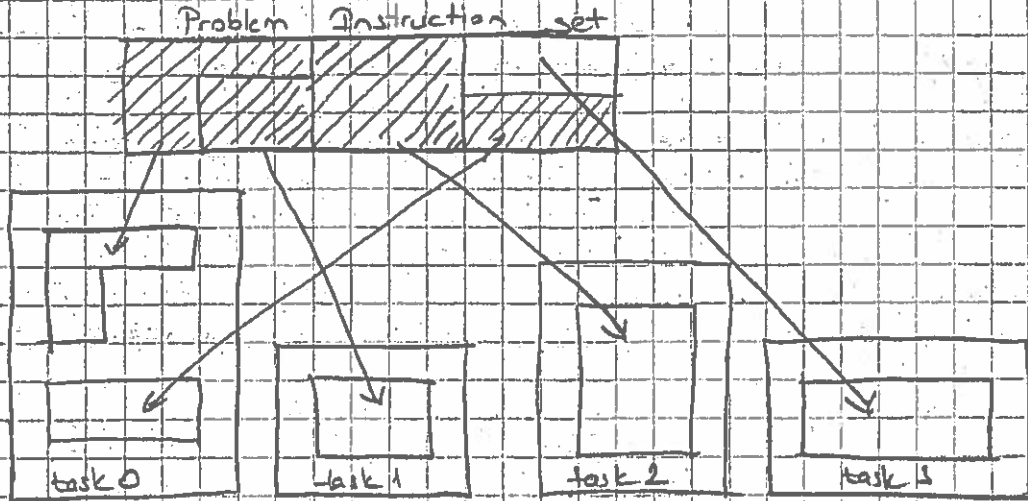
Veriyi parçalayıp parçaları işleme. İlk önce verinin bir kısmı işlenir, sonra geri kalan kısımları bölünür. Veri parçalara ayrılır. Her parçanın görevi veri bölgesinin bir kısmı üzerinde yapılır.





## → Functional Decomposition

Programın komutlarını parçalara bölme. Fonksiyonları, metodları, değişkenleri bölme.



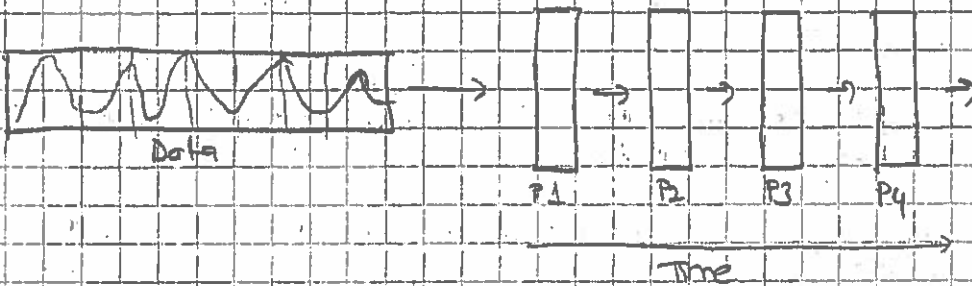
Problemi farklı görevlere parçalamak.

## → Ekosistem Modeli - Functional Decomposition

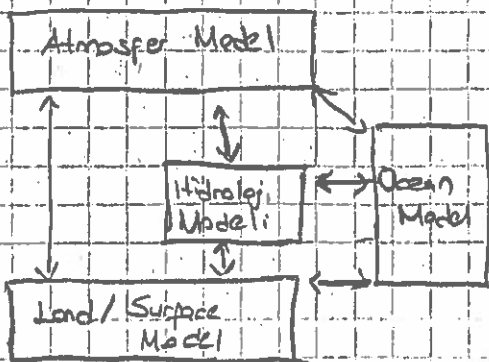
Başlangıç değerleri vardır. Bu başlangıç değerleriyle olguların durma noktaları vardır. Her biri içinde farklı modeller vardır. Fakat sonuçları bir önceki ve bir sonraki yapıya haber eder.

## → Sinyal İşleme - Functional Decomposition

4 filtre var. 10 sn'lik kumu 1. filtreye verilir. Daha sonra 2., 3. ve 4. filtreye verilir. 4 tane process her biri farklı iş yapıyor, ancak bir problemin parçası.



## → Klima Modeli - Functional Decomposition



Atmosfer model yağış hızını üretiyor. Olgular modeli alıp işliyor. Deniz yüzeyi modeline veriyor. Deniz yüzeyi modeli tekrar atmosfere veriyor. Birbirinden bağımlı işlemler.

## PARALEL PROGRAMLAMA TASARLAMA (6.SLAYT)

Haberleşme maliyetini farklı laborantlardaki bilgisayarlar artırır. İnternet, LAN gibi haberleşme yolları. Hesaplama 10ns iken haberleşme 100 ms olur. Paralel hesaplama yolları.

Paylaşımlı hafıza ile paralel hesaplama yapmak gerekir. (Super computer pahalı) İlk önce haberleşmeye bakmak gerekir.

Problem parçalara ayrıldığında ara sonuç yoksa bunlara embarrassingly parallel dır.

### Göz Önünde Bulundurulması Gereken Faktörler:

Haberleşme Maliyeti. En önemli olan şey görevler arası haberleşme maliyeti. Networku kaplıyor.

### Latency vs Bandwidth:

latency, A noktasından B noktasına mesaj göndermek için gereken zaman. haberleşme gecikmesi.

bandwidth, birim zamanda gönderilecek veri.

### Visibility of Communication:

→ Message Passing de mesajlar nerden gelmiş nereye gitmiş programcı bilir.

→ Data Parallel de programcı bilmez. Bant genişliğinin olması, latency bunda ortaya çıkar.

### Senkron ve Asenkron Haberleşme:

Senkron karşılıklı haberleşme dir. Karşı taraf mesajı alana kadar bekler. Dezavantajı bloke etmesidir.

Asenkron mesajı gönderip işine devam eder. Avantajı: non-blocking. Diğer işlerini yapmaya devam eder.

### Scope of Communication:

Haberleşme 2 şekilde olur.

→ Broadcast: Mesaj ağdaki herkese gönderilir. Hızlıdır, ancak güvenli değildir.

→ Scatter: Mesaj sahibinin posta kutusuna bırakılır. Güvenlidir.

## Haberleşmenin Verimliliği

Haberleşme zor ve masraflıdır. Bir kabloya kadar tolerans edilebilir. Örneğin; bant genişliği 100 Mbps ise bunun 300'ü haberleşmeye kaldırmak proses kilitler, iş yapamaz hale gelir. MPI bunu önceden hesaplama yapıyor.

Paralel hesaplamanın kayıp zamanı haberleşmeden kaynaklanır.

MPI - Tail

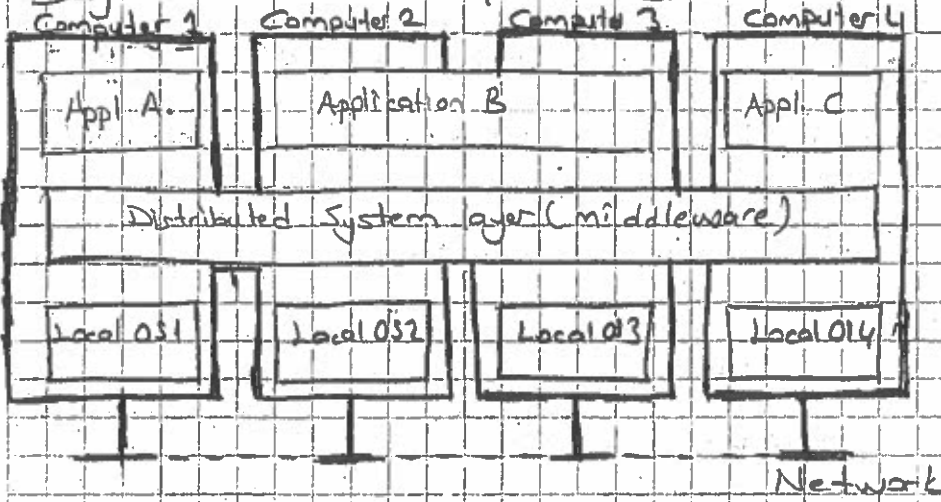
MPI - Rank

MPI - Size

## DAĞITIK SİSTEMLERİN TANIMLANMASI

Kullanıcılara tek bir sistem gibi görünür, ancak farklı coğrafi yerlerde olan pc'lerden oluşan sisteme dağıtık sistem denir. Tamamen bağımsız tek birim işlem yapabilen pc'ler Distributed memory ile yakın. 3 pc'yi bir araya getirip dağıtık sistem kurulabilir. Paralel hesaplama ile farklı, bu da paralel hesaplama ile işler abanına kalır. Dağıtık sistemlerde öyle değil.

## Dağıtık Sistem Oluşturma Yöntemleri



Middleware → Orta Seviye Katmanı

4 pc farklı yerlerde hepsinin işletim sistemi farklı. Bunları bir araya getirmek için haberleşmeyi sağlayarak bir katman lazım (Paralel hesaplama pc'lerin aynı olması lazım). Buradaki yapı gerçek (İstenilen pc sisteme eklenebilir).

4 pc'yi birbirine bağlayıp tek bir pc gibi kullanmak bağlama gerelde internetle olur. Middleware işletim sistemi üzerinde kurulur. Bu katman üzerinden haberleşme olur. Bu sistem sayı-

sinde networke bagli diger pc'lerden haberdar olunur.

Dagitik sistemin bir arayuzi olur. Bu arayuz arka tarafdaki 4 pc'yi tek gibi gosterir. Avrupa gridi isin ozelligine gore sistemde basta olan pc'ler uzerinde hangisinde calisacagina karar verir. (uygun olanı secer.) Her pc'nin onalik durumunu gorunur. Sistemdeki pc'lerden birinin gitmesi durumunda; digerleri onun gorevini devralir, sistem cokmez. Torrentler, ...

Public bir sisteme herkes girebilir. Degilse hesap gerekir. Ideal sisteme samon icinde yaklasilmistir.

Onemli olan middleware'in ozellikleri ve ne kadarini yapabildigidir. Kullanilacagi yere gore ozellikleri degisir. Bir veritabaninin yedeeginin başka bir yer uzerinde tutulmasi pahalı bir yapı.

Transparency: Bazı ozelliklerin gizlenmesi (Tek pc olarak gosterme)

→ Access: Erisimin gizlenmesi. Arka plandaki kaynakların gizlenmesi. Asıl amaç, tek bir pc gibi gosterme. Global isletim sistemi kurmak en uc noktada yapılmak isteyen. Hiçbir sekilde cokmayacak. Bilim adamları için global isletim sistemi olusturmak internet uzerine kurulmuş sistemler var. Örn: Grid, web arayuzi var. Kendi yaptığımız uygulamalar, işler girildiginde gorulebiliyor. Kaynaklara erişimin gizlenmesi gerekiyor.

→ Location: Kaynağın nerede olduğunu, konumunu gizleme

→ Migration: Bir kaynağın başka bir yere tarındığını gizleme. Kullanıcı görmemeli. Veritabanının elazigla taşınması.

→ Relocation: Servislerin yerini görmemeli ve bilmemeli.

→ Replication: Kaynağın yicklenmesini gitlemek.

→ Concurrency: Kaynağa eş zamanlı erişimin gizlenmesi. Alışveriş sitelerinde barkalarının siparişlerinin görümemesi.

→ Failure: Sistemdeki hataları gizlemeli. Örn: veritabanının cöktüğü san kullanıcıya söylenmez.

→ Flexibility: Esneklik. Sistemi oluşturunken yeni kaynak eklendiğinde ne olacak? sistem daraltılıp genişletilebilir mi? Mikro kernel modeli. Sistemin yeni kaynak eklenmesi, çıkarılabilmesi durumunda ne kadar esnek olursa kullanması o kadar kolay olur.

→ Reliability: Güvenilirlik. Sisteme bir iş yaptırırken sisteme güvenilmeli. Sistemin cökmeyeceğinden emin olunmalı.



→ Performance : Doğru sistem inşası sebebi performance. Elimizde olmayan kaynakları kullanma HPC → Hyper Performance Commu

### Scalability:

Ölçeklenebilirlik ihtiyaç olduğunda sistemin kaynaklarını artırma. 1000 kişi için oluşturulan sistemin gerektiğinde otomatik olarak 100000 kişiye yanıt vermesi. Gerekli donanımların otomatik sisteme alınması çok zor bir iştir. (3-6 yıl önceye kadar mümkün değildi) şimdi bulut bilişimle rahatlıkla yapılıyor.

→ Serverlarda ölçeklenebilirlik istenir. (web ve application server) web server kiralanır, sayfaya çok istek geldiğinde trafiği bölünmüş data centerde trafik artırılır. Bir noktadan sonra bu da problemi çözmeyi yük kontrol edilip otomatik server devreye alınması.

→ Merkezi Servis Yaklaşımı : herkes aynı yere erişip işlemi yapıyor. Tek server üzerinde işlemleri yapmak. Yine göre doğru bir yöntem (Üniversite için tek server). Tüm üyelerin otomasyonunu tek serverda tutmak mantıklı değil. Bunun için aynı anda birden fazla server olacak. 3 makineye kurup gelen istekleri paylaştırmak gerekir. (Merkezi olmaktan çıkarılır). Merkezi sistemden uzaklaşmak her zaman yapılabilecek bir iş değildir.

Mainframe'ler bankalarda önemli ölçüde işler üzerinde tutulur. Ölçekleme açısından temel problem geliştirilmesi kolay değildir. CPU RAM bir yere kadar artırılabilir. (Mainframe) Bir banka için çok fazla Centralized servislerin temel problemi ölçeklenebilirliktir.

→ Client-Server Modeli : Mainframe'lere göre daha esnek yük dağıtılabilir, belli bir sayıya kadar sunucu eklenebilir. Ölçekleme limiti networkle ilgili. Trafik belli bir limite ulaştığında ek server alınmıyor. En çok kullanılan model client-server modelidir. Client-server'da sunucu ve istemci açık şekilde tanımlanmıştır.

→ Peer-to-Peer Modeli : Tüm bilgisayarlar birbirinin eşi. Peer'ler hem sunucu hem istemci olabilir. Dosyalar peerlere dağıtılıyor. Herkese dayalı sunucular hem de client olarak diğerlerinden dayıya çekebiliyor. Hizmet sunucu için uygun değil. Ölçeklenmeyi kolaylaştırır. Depolama miktarını arttırabiliyor. Hizmet sunmada algoritmaları peer-to-peer paylaşıyor.

→ Kalıp ve İnce İstemciler : İnce istemci merkezi sunucuya bağlanıp sonuçları ince istemcide gösteriliyor. 50 kullanıcı servera bağlanıp uzak makineyi ile

→ Cok Katmanli Mimari: Degisik isler yapan sunucular var. Yuku mumkun olduđu kadar cok sunucuya dagitma. Web sunucular database sunucusu ayri yerlerde birbirleriyle haberlesiyorlar. Web serverlar 2 katmanli basliyor. 3 katmanli mimarilerde isi goster, hesaplari yap, veritabanini tut (yaygin olan) sunucular orka olunde dagitiliyor. (En yaygin olarak kullonulan hizmet sunumu icin) Sayi olarak peer-to-peer fazla. Bunlar merkezi modele baglanıyor. Bir tane merkezi sunucu var. Yuk miktarl artticta belli bir yere kadar baskleyebiliyoruz.

### Merkezi Olmayan (Decentralized) Sistemlerin Problemleri:

- Herkese kurtulmak dhenli - Birde fazla pedo tutulsun
- Hiçbir bilgisayar tam olarak sistem durumu hakkında bilgiye sahip degildir. Kismi olarak bilinir.
- Bilgisayarlar dosya paylasiirken sadece lokal bilgiye gore karar veriler.
- Bir parcin bozulmasi sistemi etkilemez. Bozulonin ustundeki veri bir diger parca da vardir. Ornegin, NoSQL, dagitik dosya sistemleri

NOT: Centralized sistemlerde cluster yapilsa bile veri parcalanmaz.

Olçekleme yapmak icin kullonulan methodlar:

Client - Server modellerinde formun kontrolu server tarafinda kontrol edilebilir. Burada gereksiz islemi yapilir. Form kontrolu client tarafinda yapılmalı.

DNS'ler her networkte bir tane bulunur. En yakin DNS sever cevap verir. Olçeklenebilirlige gicel bir dnmek. Network genisledikce genisleyen yere bir tane DNS koyup digerleri ile haberlesmesi saglanir.

### Dagitik Sistemler Tasarlanirken Bazı Varsayimlar Yapilir:

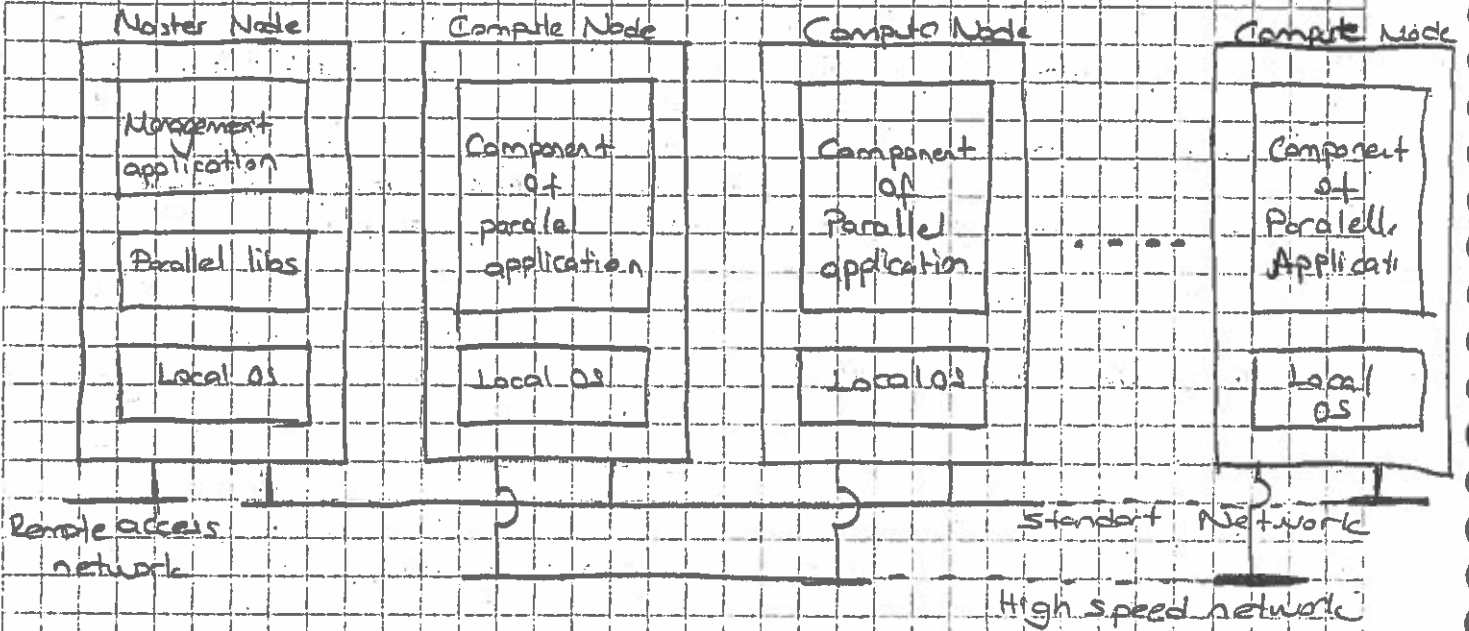
- Network'in guvenilir olduđu (isi sonuna kadar yapması)
- Network'in guvenli olduđu (kullonucu adi, sifre ile girilmesi)
- Network'in homogen olmasını kabullenme
- Topolojisinin degismedigini varsaymak
- Latency (ağı gecikmesi) sıfır kabul etmek.
- Bant genisliginin sonsuz olmasının varsayılması

→ Bir yere bir yere taşıma maliyetinin sıfır olduğunu kabul etmek

→ Tek bir yönetici olduğunu kabul etmek

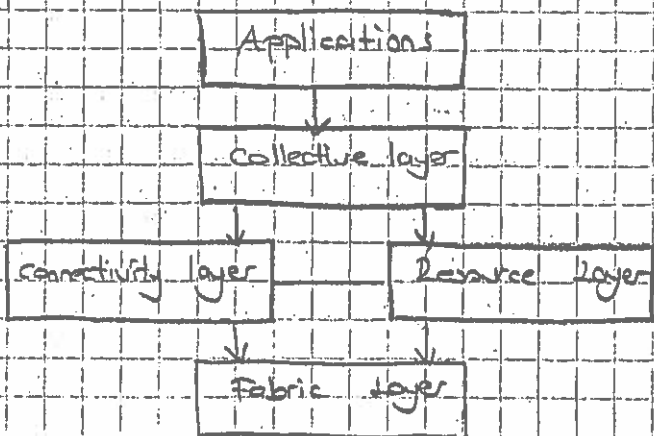
Bu noktalar mutlaka göz önünde bulundurulmalı. Bunlar göz önüne alınarak dağıtık sistem kurulmalı.

## Cluster Computing Systems (Küme Bilgisayar)



Aynı türden pc'ler kullanılmalı, pc'ler birbirine bağlı olmalı. High speed network veri değişimi için kullanılır. Standard network işi göndermek, iş almak, web arayüzü çağırır. Bir tane yönetici node (master node) işi parçalamak, sonuçları toplamak için vardır. Paralel hesaplamanın parçası, diğer pc'lerdedir. Management node işi parçalara böler, yapabilecek pc'lere gönderir ve geri alır. Pc'ler aynı mekandadır.

## Grid Computing System :



Pc'ler dünyanın her yerinde olabilir. Sanal organizasyonlar oluşturulur.

Fabric Layer: CPU, RAM, pc'yi kontrol eden ve yöneten: kalman

Connectivity Layer: Protokoller, nasıl erişileceği.

Resource Layer: Hangi pc'ne kaçır mesgul, ne kadar CPU kullanıyor.

## Collective Layer:

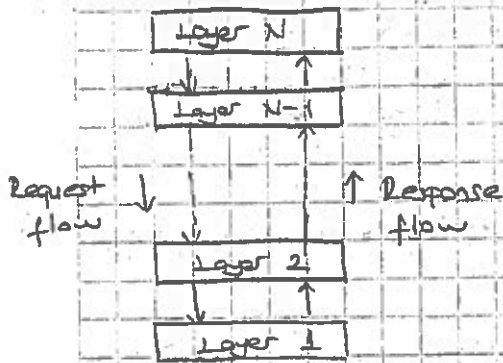
Application Layer: Bir programı collective layera verir. O dağıtır.  
Clusterden farklı peller farklı yerlerde. Hizmet odaklı mimari'nin temel adımı.

## MİMARİ MODELLERİ:

4 mimari var.

- Katmanlı Mimari
  - Nesne Tabanlı Mimari
  - Veri Merkezli Mimari
  - Olay Tabanlı Mimari
- } Farklı zamanlarda geliştirilen farklı kullanım alanları vardır.

### → Katmanlı Mimari



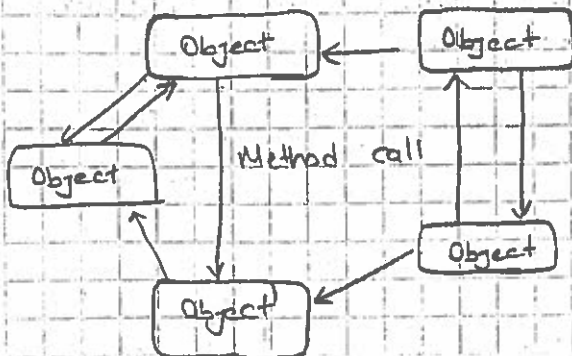
Asğıdan yukarıya istek, yukarıdan aşağıya cevapların geldiği model.

### Avantajları:

- \* Her katmanın ne iş yaptığının açık olduğu yerlerde kullanılır.
- \* Dağıtık veritabanının olduğu yerlerde kullanılır.

Dezavantajı: Katmanlar arası haberleşme tek yönlü. Aradaki katman atlanmaz.

### → Nesne Tabanlı Mimari



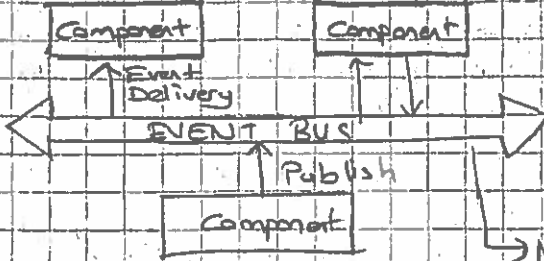
Eldeki programı başka palye kullanırmak. Amaç; başka makinelerdeki nesneleri kullanabilmek. RMI, RPC yazılımı ile yapılabilir.

Client-server mimarisine benzer. En yaygın kullanılan modeldir.



Veri paylaşım sistemleri için geliştirilmiş modeller (Veri Tabanlı, Olay Tabanlı)

### → Olay Tabanlı Mimari



Event bus var. Ortak veri yolu. Bağlı Peler busları geçen her şeyi görürler. Bir obje olduğunda busa mesaj üretilir.

Publish subscribe modelleri, son yıllarda geliştirilen en çok kullanılan sistem

→ Mesajın ulaşip ulaşmadığından emin olunmaz.

Bir broker çeşitli yerlerden veri gelebilir. Her biri farklı konuya yayınlanıyor. Müşteri olarak yayınlara üye olunuyor.

Kullanım alanı farklı (Web server çalışmaz).

Avantajı: Çok fazla mesaj üretiminin olduğu sistemlerde kullanılıyor.

### Problem:

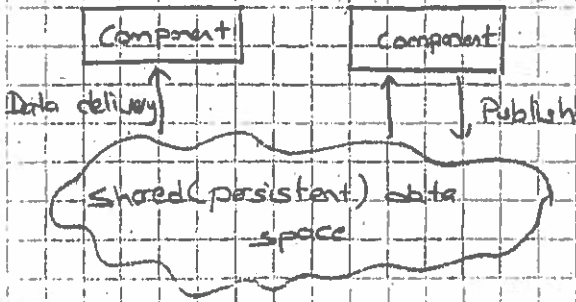
Publish yapan kişi subscribe'lardan haberdar değil. Mesajı gönderen yayına üye olan kimselerden haberdar olmaz. Buna dıcapılın denir. Bu sağlanırsa üyelerle ilgilenmek zorunda olmaz.

→ Mesajların yerine ulaşip ulaşmadığı bilinmez. Hassas işler yapılıyorsa bunun için ek bir şey eklemek gerekir.

→ Kayıp mesajları subscriber ile broker arasında gelen mesajların numaralarını bildirmesi gerekir.

Bu problemlerden dolayı event bus servera dönüşüyor.

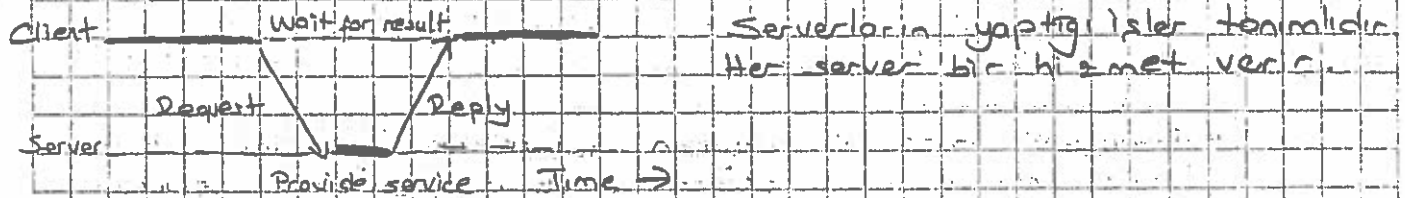
### → Veri Tabanlı Mimari



Haberleşme paylaşım ile veri uzayında yapılır. İşler dosyalar üzerinden yapılıyor. Veriler dosyalar üzerine yazılır.

4 farklı dağıtık sistem modeli. Bunları gerçekleştirirken 2 temel mimari var.

## 1. Centralized Architectures :



Server başka serverlara clientlık da yapıyor olabilir. Serverın yaptığı iş tam olarak belli ise server hizmeti sunabilmek için diğer serverlara clientlık yapabilir. (Decentralized mimari)

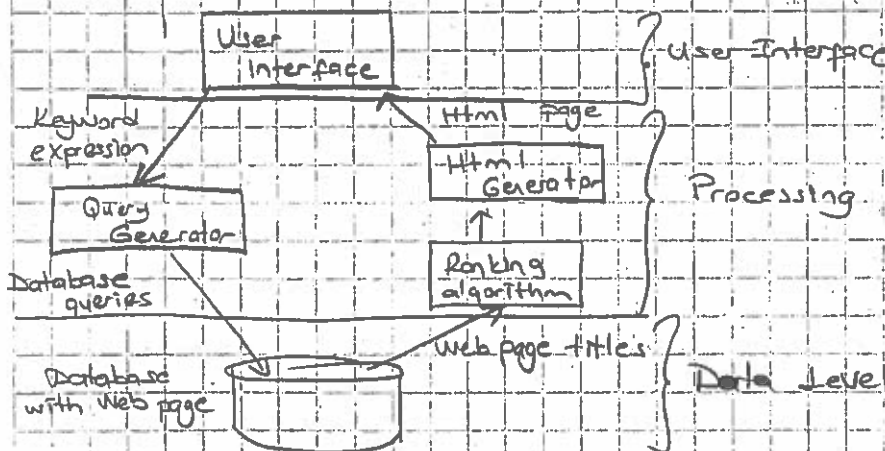
Client-server arasında bağlantının güvencili olduğu düşünülür. İstek yapılır, bağlantı kesilir. Server adrese göre cevabı gönderir. Bağlantısızdır. Avantajı, çok sayıda clientın servera bağlanabilmesi.

Problem: Hesaba para yatırma isteğini yerine ulaşıp ulaşmadığı kontrol etmek gerekir. Protokolle mümkün değil. Belli bir zaman içinde ulaştığı ile ilgili mesaj gelmezse tekrar gönderir isteği. (1. isteğin id'si ile 2. gönderilir) işleri karmaşıklaştırır.

Temel olarak kabul edilen 3 katman vardır.

→ User Interface Katmanı: Web arayüzünde görünenler, kullanıcılar var.

→ Processing Katmanı: En altta bir veritabanı. Örnek: google, arama motorları.



→ Data Katmanı:

Cok Katmanli Mimari: 2 Katman var. Biri client, diğeri server. Katman sayisini tam olarak belirtmek mümkün değildir. İstek cevap döngüsü vardır.

Döğrulama kullanıcı tarafında.

## 2) Decentralized Mimariiler:

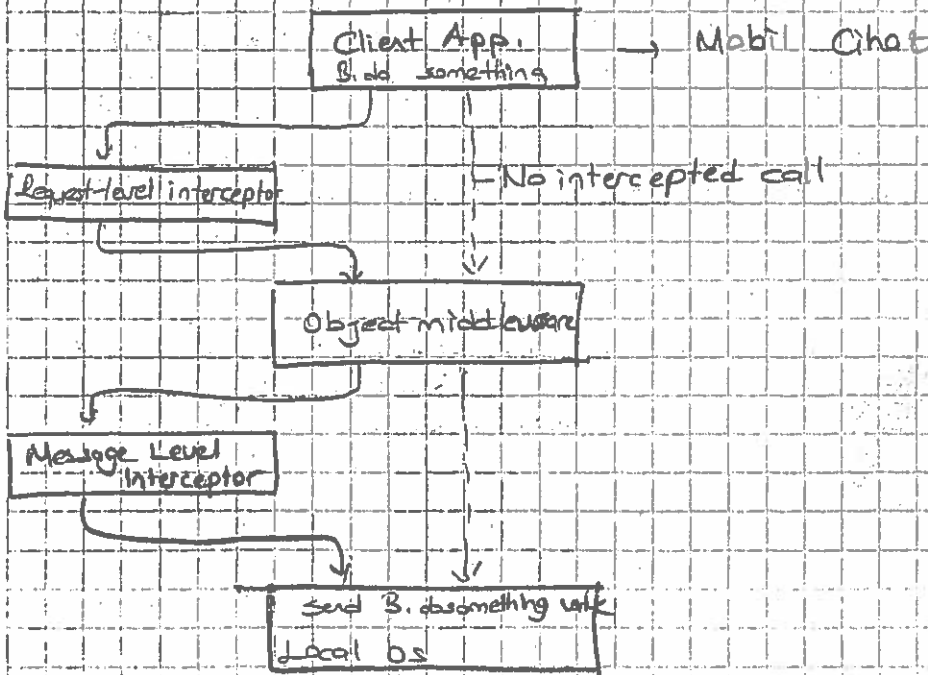
Server ile client arasındaki ayırım net olarak yapılamaz. Bir makine hem server hem client olabilir. Bunlara merkezi olmayan mimari denir. Örn:

Peer-to-Peer Networkler: Makine hem server hem istemci. Network'e üye olunuyor. Hem dosya paylaşıyor, hem de dosya indiriyor. Bir makine oluyor. Makineler eş olarak bakılıyor. Buna peer-to-peer (eşler arası paylaşım) denir. Asıl önemli olan distributed hashler. Hızlı çalışan mimarilerdir.

Yapısal olan ve yapısal olmayan yaklaşımlar var. Her node'un unick id'si var. Üye olunca otomatik id veriliyor. Dosyalarında unick id'si var. (Dosya id ile unick id eşleştiriliyor ve hashtable'de saklanıyor. Tek tek sorup bulma).

ÖDEV: Peer-to-Peer Networklerde structer-unstructer nasıl çalışıyor? Distributed hashtable nasıl çalışıyor, bunun için geliştirilen yaklaşım?

## Interceptors:



Buradaki amaç hizmet kalitesini yükseltmek.

Değişik sistem üzerinde elimizdeki kaynakları kullanmak için farklı yaklaşımlar var. Kodu yeniden kullanmak amaç.

Remote-object invocations : Mesajların arasına girmek. Katmanların arasına girmek.

Virtualde kullanılıyor. Amaç hizmetin kalitesini yükseltmek. Görselleştirme, simülasyon varsa bunları başka katmanlara yaptırmak lazım. Modülerliği sağlıyor. Server sadece simülasyonu yapmaz, sadece cevap verir. Interceptorlar kendilerini clienttan gizlerler.

Centralized ve decentralized in geliştirilmesi için bazı yaklaşımlar düşünülmüş. Adaptif sistemler, otomatik ayarlıyor kendini.

→ Separation of Concerns : Yapılmak istenen işler birbirinden ayrılır. İşlemler birbirlerinden ayırt edilecek katmanlara bânılır.

→ Computational Reflection : Runtime anında değişiklik yapma (reflection). Yansıma yöntemiyle yapılır. Diğer çağırarak yerine aynı nesneyi kullanan başka programdan çağırarak yapılıyor. Programı çalışırken modifiye etme. Ortama adapte etme. Java da var.

→ Component based Design : Programı parçalara ayırma. Veritabanı kısmını ayırma. Hesaplama kısmı ayrı model, user-interface kısmını ayrı bir model yapma.

Gelişim Nedeni :

Bu 3 şart tam olarak sağlanamamış. Reflection RPI tüm değişikliklere izin vermiyor. Problem birbirinden ayrılomayabiliyor her zaman.

Otonom sistem : Kendilerini otomatik olarak ayarlayan sistemler.

## PROCESSLER

Bir programın derlenip RAM üzerine yüklenmesi ve processlerin sırasıyla CPU'da çalıştırılması. Process kontrol bloğu olsun. Birden fazla CPU gibi kullanım multiprocessing denir.

## Avantaj :

Disadvantage: Her bir process'in olması için bazı şeyler gerekiyor. Ayrı bir hafıza zaman alıyor, hem de pahalıya mal oluyor.

Thread: Processleri küçük parçacıklara bölme. Threadler kendi processlerinin alanlarını, hafızalarını kullanır. Threadler 2 seviyede oluşturulur: User ve Kernel.

User: Kullanıcı her türlü şeyi gerçekleştirebiliyor.

Kernel: Her process için birden fazla thread oluşturuluyor. I/O olursa threadlerden birine yönlendiriyor, diğerleri bloke olmuyor. Bu her iki modeli kullanan birde hibrit sistem var.

Lightweight process: Hibrite benzer, sıklıkla kullanılmış. Kernel tarafında lightweight process denilen (biraz process, biraz thread) oluşturuluyor. I/O bunlardan birine yönlendiriliyor.

Thread paralelligi, paralel iş yapabilmeyi sağlar. Bir işin parçalarına bölünmesi için programlama dilinin bunu desteklemesi lazım.

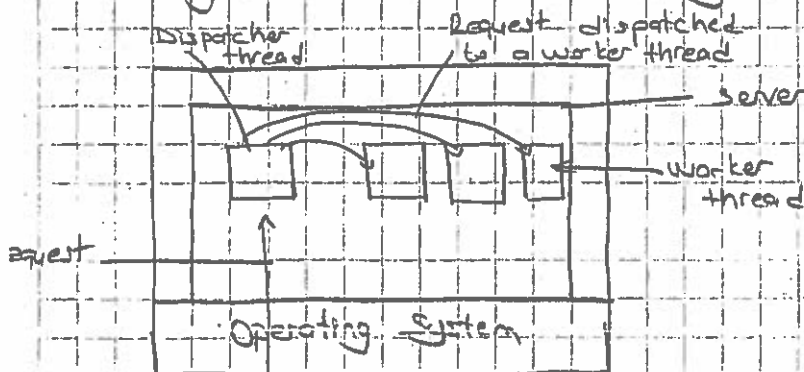
Örnek: Excelde yazı yazarken hesaplamaların yapılması diske kaydedilmesi, gramer kontrolü yapmak için birden fazla process çalıştırmak yerine bir process için belli bir alan ayrılıp bu işleri yapması için threadler oluşturuluyor. DOS'ta yazı yazarken başka birşey yapılmıyor.

Serverlar, clientlar multithreaded şekilde yazılıyor.

### Multithreaded Servers (!)

Clientten gelen isteğe cevap vermek. Daha önceki isteklere cevap vermesi, diğer istekleri düzürmemesi lazım. Bunun için serverda birden fazla thread (10) oluşturuluyor. Bu threadler gelen istekleri cevaplar. (Basit çözüm)

1.1. istek gelirse 10 thread cevap veremez. Thread sayısı sınırlı. Bunun yerine başka çözüm geliştirilmiştir.



Dispatcher thread istegi worker thread'e iletiyor. Kendisi dinlemeye devam ediyor. Network haberleşmesi yapan tek thread var. Dispatcher 80 no'lu portu dinleyip worker thread'e veriyor. Workerlar hepsi aynı birinin kopyası. Dağıtık sistemlerde birçok yerde kullanılır.



Client tarafında tarayıcı açıp sayfanın yüklenmesi sırasında html kodları haricinde java scriptler, csslerin yüklenmesi gerekir. (Tek thread olmasından dolayı böyle) Bunu çözmek için client tarafında multithreaded kullanılmaya başlandı. Head kısmında harici objeler varsa bir tane thread bunları getirmeye başlar. Gerici bölgeye yüklüyor. Tarayıcı bu resimlerin hepsini aynı anda yüklüyor. (Birden fazla thread) Bir sayfayı hızlı bir şekilde yüklüyor.

#### Process Çalışması

Process



RAM'de yer ayırılır.



Process Kontrol Bloğu Oluştur

↓ → CPU Dispatcher CPU'ya verir processi. Birden fazla process aynı anda

Processi olan CPU singula bunları çalıştırır. yönlendirir.

#### Sanallaştırma

İşletim sisteminin kullanacağı max kaynak miktarının bir adımı olması. Aynı anda birden fazla işletim sistemi donanımdan bağımsız olarak yazılım donanım üstünde çalıştırmak. (Virtualization)

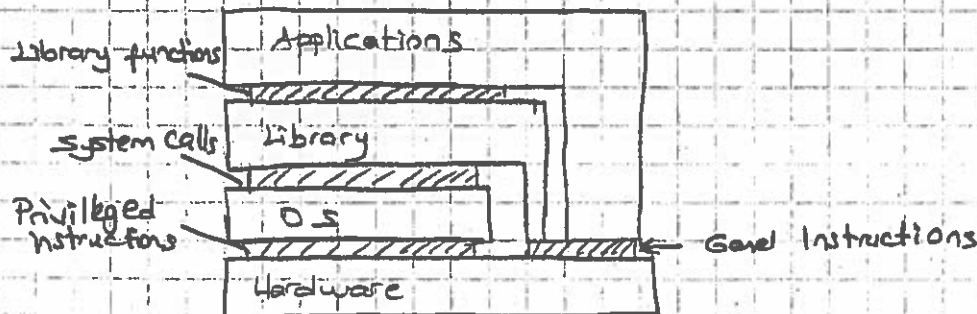
Sanal processler oluşturup, bunu tek bir process gibi çalıştırmak. Birkad farklı katmanda yapılabilir.

→ Birinci Seviye: Donanım ve yazılım arasında bir interface

→ İkinci Seviye: Sistem çağrıları, işletim sistemi sadece yapan işletim sistemi haricinde programların sistem çağrısı yapması için interface

→ Üçüncü Seviye:

Privileged instructions, kernel seviyesinde (Nüveler)



→ 4/5. Seviye: Virtual Machines Monitore, donanımı gizleyecek katman. Donanımla haberleşme yeteneği olacak. Üzerine yazılımlar kurulacak. Bu yazılımlar donanımı kullanıyormuş gibi olacak. Bir donanım gibi hareket edecek. İstedikçe kodlar işletim sistemi kumullarını. Virtual Machine Monitore ile ilgili ilk makale 2001'de yayınlandı.

## Alantıları :

→ Sistem çağrılar, interruptlar işletim sistemini boş bekleten şeylerdir.

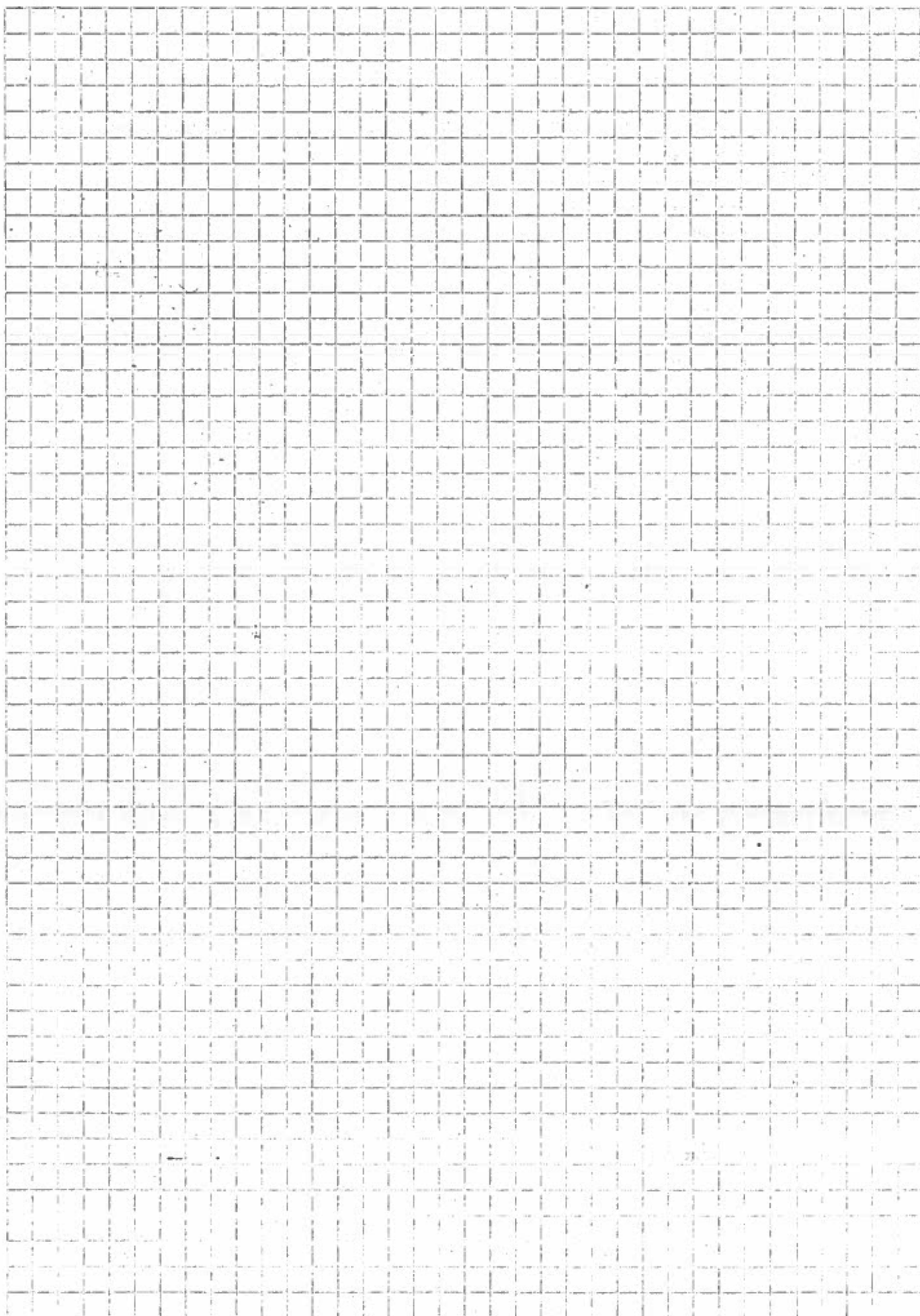
→ Bc makine kullanılabilecek kadar kullanma.

5'i gerçekleştirmek için donanım desteğine ihtiyaç var. 64 bit adresleme uzayı.

Hizmet olarak sunulabilir mi? (Cloud Computing) Alt yapıyı servis gibi sunmayı yapıyor. Donanımı hizmet olarak sunuyor.

Cloud Computing ihtiyaç olunca sanallaştırmayı otomatik yapıyor.

Amazon → sanallaştırma öncüsü.



# VERİ MADENCİLİĞİ YILIÇI SINAVI

Fen Bilimleri Enstitüsü

Yrd. Doç. Dr. Şule Öğüdücü

## Çözümler

1. Veritabanı ve veri madenciliği işlemleri arasındaki farkları yazınız. Örneklerle açıklayınız.

**Cevap:**

- Veritabanı sorgulamaları tanımlıdır. Veri madenciliği işlemlerinde aranan bilgi veri içinde gizli olduğundan tanımlı bir sorgulama yoktur. Veritabanı sorgulamaları için SQL sorgulama dili kullanılır. Veri madenciliği işlemlerinde kullanılan yaygın sorgulama dili yoktur.
- Veritabanı işlemleri canlı veri üzerinde yapılır. Veri madenciliği işlemleri, üzerinde işlem yapılmayan veri üzerinde yapılır.
- Veritabanı işlemlerinin çıkışı bellidir ve verinin bir altkümesidir. Veri madenciliği işlemlerinin sonucunda elde edilecek bilgi belirgin değildir.
- SQL sorgulamasına örnek: DVD satın alan tüm müşterileri bul.
- Veri madenciliği işlemlerine örnek: DVD birlikte sıkça satın alınan ürünü bul

2. Veri madenciliği uygulamalarında kullanılan modeller kaçına ayrılır? Açıklayınız.

**Cevap:** İkiye ayrılır.

- kestirime dayalı: sınıflandırma, eğri uydurma, zaman eğrileri
- tanımlayıcı: demetleme, özetleme, bağıntı kuralları, sıralı diziler

3. Bir veri madenciliği uygulaması tasarlanacaktır. Çalışılacak veri kümesindeki niteliklerden biri *yaş* bilgisidir. Bu niteliğe ait değerler şu şekildedir: 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 30, 45

- (a) *yaş* niteliğine ait veri için gürültüyü azaltmak amacıyla eşit aralıklarla bölme genişliği 3 olacak şekilde bölmeleme (*binning*) yapılacaktır. Bu işlemi yapmak için gerekli adımları gösteriniz. Yukarıdaki veri için bu işlemin etkilerini yorumlayınız.

**Cevap:**

Bölme No	Veri	Ortalama	Alt-Üst Sınır
1	13-15-16	14,7-14,7-14,7	13-16-16
2	16-19-20	18,3-18,3-18,3	16-20-20
3	20-21-22	21-21-21	20-22-22
4	22-25-25	24-24-24	22-25-25
5	25-30-45	33,3-33,3-33,3	25-25-45

Aynı değere sahip veriler farklı bölmelerde yer aldıkları için farklı değerler alıyorlar (Örnek: 16 hem 1. hem de 2. bölmede var. 1. bölmede ortalama kullanıldığında 14,7 değeri ile, 2. bölmede 18,3 değeri ile düzeltiliyor). Ancak alt ve üst sınır kullanılarak düzeltme yapılırsa bu problem ortadan kalkıyor.

- (b) Verideki aykırılıkları belirlemek için hangi işlemler yapılabilir?

**Cevap:** Demetleme

(c) Verideki gürültüyü düzeltmek için başka hangi yöntemler uygulanabilir?

Cevap: Eğri uydurma

4. Aşağıda verilen tablodaki verileri sınıflandırmak için karar ağacı oluşturulacaktır.

A	B	Sınıf
T	F	+
T	T	+
T	T	+
T	F	-
T	T	+
F	F	-
F	F	-
F	F	-
T	T	-
T	F	-

(a) A ve B nitelikleri için bilgi kazancını hesaplayın. Karar ağacı algoritması hangi niteliği kullanarak bölmeleme yapar?

Cevap:

A ve B niteliklerine ait sınıf dağılımları

	A=T	A=F
+	4	0
-	3	3

	B=T	B=F
+	3	1
-	1	5

Bölünmeden önce entropi:

$$E_{orig} = -0.4 \log 0.4 - 0.6 \log 0.6 = 0.9710$$

A niteliği kullanılarak bölündüğü durumda bilgi kazancı:

$$E_{A=T} = -\frac{4}{7} \log \frac{4}{7} - \frac{3}{7} \log \frac{3}{7} = 0.9852$$

$$E_{A=F} = -\frac{3}{3} \log \frac{3}{3} - \frac{0}{3} \log \frac{0}{3} = 0$$

$$Gain(S, A) = E_{orig} - 7/10 E_{A=T} - 3/10 E_{A=F} = 0.2813$$

B niteliği kullanılarak bölündüğü durumda bilgi kazancı:

$$E_{B=T} = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = 0.8113$$

$$E_{B=F} = -\frac{1}{6} \log \frac{1}{6} - \frac{5}{6} \log \frac{5}{6} = 0.6500$$

$$Gain(S, B) = E_{orig} - 4/10 E_{B=T} - 6/10 E_{B=F} = 0.2565$$

Veriyi bölmek için A niteliği seçilir.

- (b) A ve B nitelikleri için Gini index deęerini hesaplayınız. Karar aęacı algoritması hangi nitelięi kullanarak bölmeleme yapar?

**Cevap:**

Bölünmeden önce Gini index

$$Gini_{orig} = 1 - 0.4^2 - 0.6^2 = 0.48$$

A nitelięi kullanılarak bölündüęü gini göstergesi için kazanç:

$$G_{A=T} = 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 0.4898$$

$$E_{A=F} = 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0$$

$$Gain(S, A) = G_{orig} - 7/10Gini_{A=T} - 3/10Gini_{A=F} = 0.1371$$

B nitelięi kullanılarak bölündüęü gini göstergesi için kazanç:

$$G_{B=T} = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0.3750$$

$$E_{B=F} = 1 - \left(\frac{1}{6}\right)^2 - \left(\frac{5}{6}\right)^2 = 0.2778$$

$$Gain(S, B) = G_{orig} - 4/10Gini_{B=T} - 6/10Gini_{B=F} = 0.1633$$

Veriyi bölmek için B nitelięi seçilir.

- (c) Bilgi kazancı ve gini index deęerlerinin bölmeleme için farklı nitelikleri seçmesi mümkün müdür? Açıklayınız.

**Cevap:**Gini göstergesi ve entropi ölçütleri benzer özellikler gösterir, ancak bu ölçütleri kullanarak hesaplanan kazanç aynı yukarıdaki örnekte olduęu gibi aynı özellięi göstermeyebilir.

