

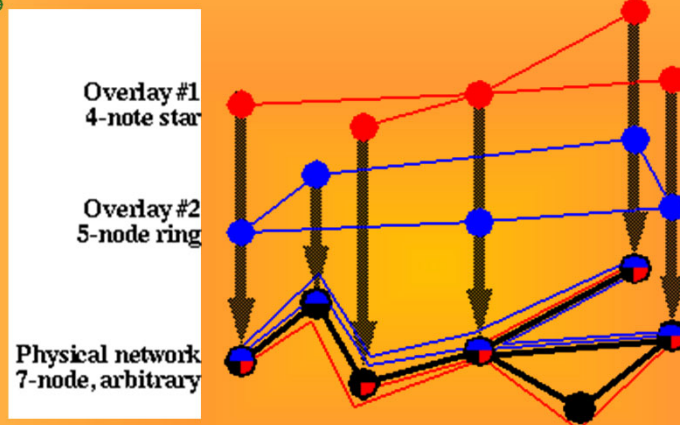


CDNs and FMB convergence

Everyone in the same network
Service content rules



What is an Overlay ?



What is the topology of this network?

WHICH network??



Overlay Networks: Overview

- Networks built using an existing network as substrate (*Virtual Networks*)

Internet

- Initially an overlay on the POTS (Plain Old Telephone System) network
- Overlays are a (quasi) structured virtual topology above the basic transport protocol level that facilitates deterministic search and guarantees convergence
 - Overlays could consist of routing software installed at selected sites, connected by encapsulation tunnels or direct links
- Examples of overlays:
 - MBone, 6Bone
 - P2P (Napster, FreeNet, Gnutella, Bittorrent)
 - Cooperating Caches
 - Server Farms
 - Content Distribution Networks (CDNs)



Content Distribution Networks

Client-Server and distribution models
Caching and load balancing

6

Learning outcomes

- Understand the purpose of content distribution on a network
- Discuss the rationale for CDNs and comment on different alternatives
- Describe the architecture of a CDN

© Rui L. Aguiar (rui.laa@det.ua.pt) - Uni. Aveiro

(recall FR): Web caches (proxy server)

Goal: satisfy client request without involving origin server

- user sets browser: Web accesses via proxy server
- browser sends all HTTP requests to proxy
 - object in cache: cache returns object
 - else proxy requests object from origin server, then returns object to client

```

graph LR
    subgraph Clients
        C1[client]
        C2[client]
    end
    subgraph Servers
        O1[origin server]
        O2[origin server]
    end
    P[Proxy server]
    C1 -- "HTTP request" --> P
    P -- "HTTP response" --> C1
    C2 -- "HTTP request" --> P
    P -- "HTTP request" --> O1
    O1 -- "HTTP response" --> P
    O2
  
```



More about Web caching

- Proxy server acts as both client and server
 - typically proxy server is installed by ISP (university, company, residential ISP)
- Why Web caching?**
- reduce response time for client request
 - reduce traffic on an institution's access link.



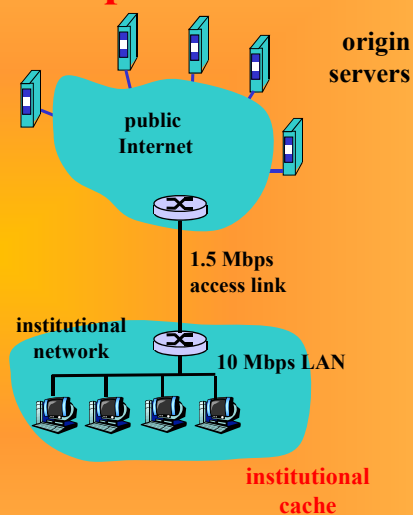
Caching example

Assumptions

- average object size = 100,000 bits
- avg. request rate from institution's browsers to origin servers = 15/sec
- delay from institutional router to any origin server and back to router = 2 sec

Consequences

- utilization on LAN = 15%
- utilization on access link = 100%
- total delay = Internet delay + access delay + LAN delay
= 2 sec + minutes + milliseconds





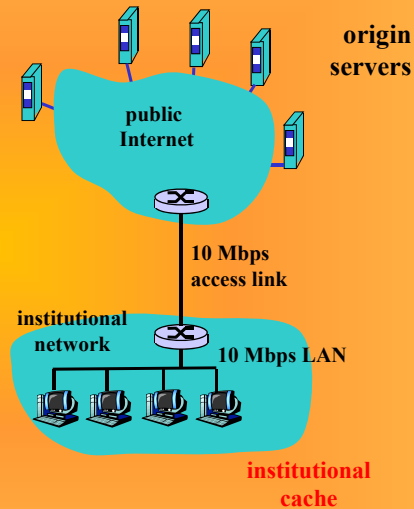
Caching example (cont)

possible solution

- increase bandwidth of access link to, say, 10 Mbps

consequence

- utilization on LAN = 15%
- utilization on access link = 15%
- Total delay = Internet delay + access delay + LAN delay
= 2 sec + msec + msec
- often a costly upgrade



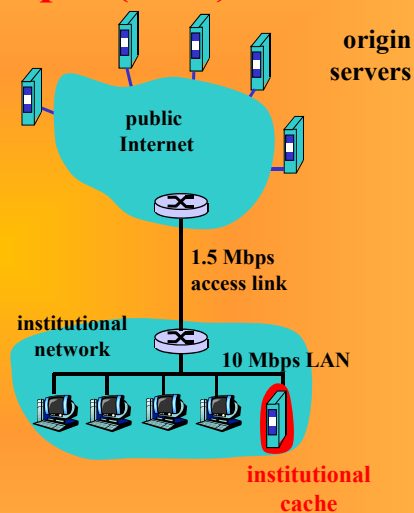
Caching example (cont)

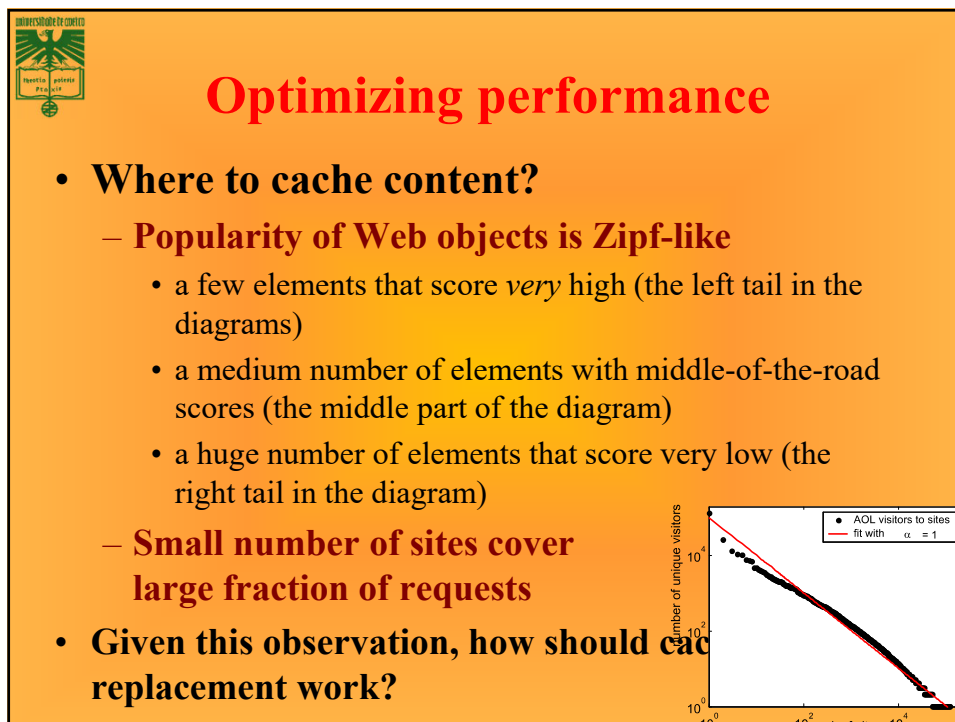
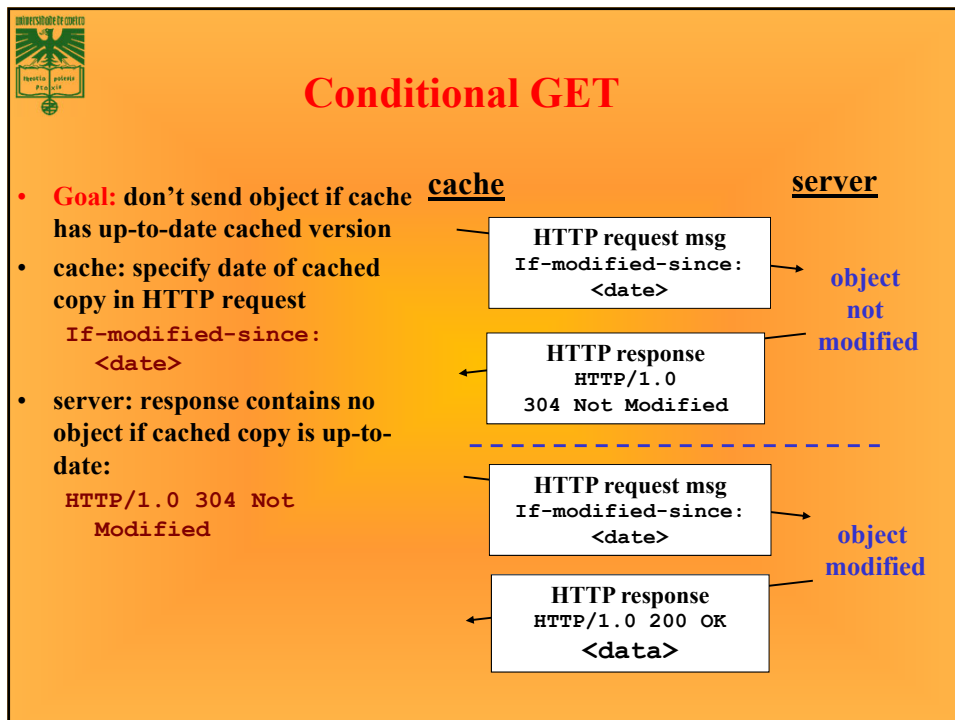
possible solution: install cache

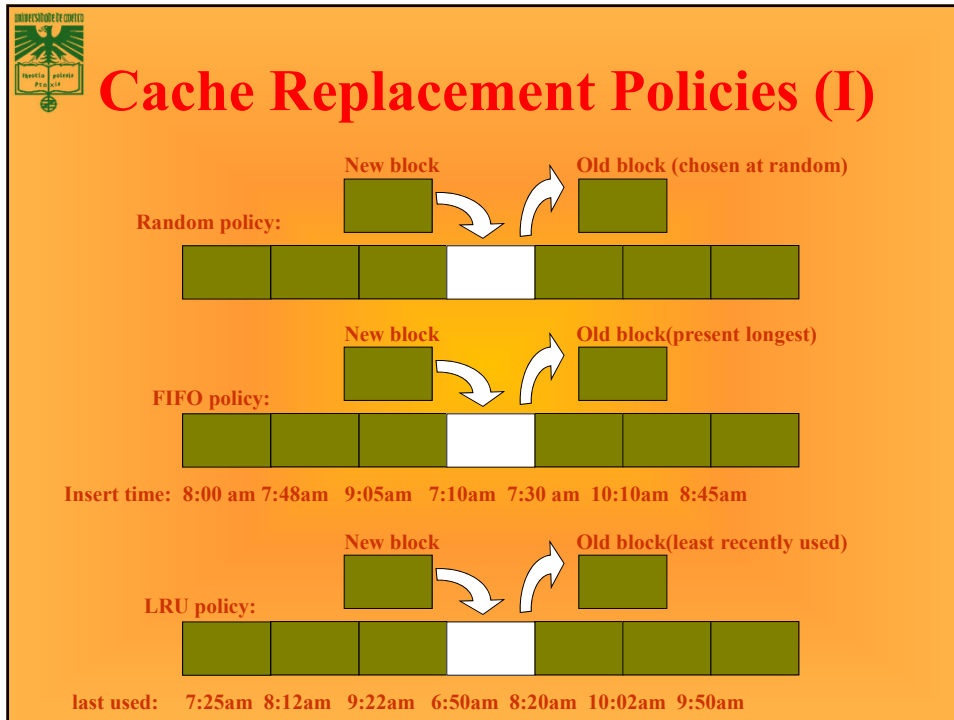
- suppose hit rate is 0.4


consequence

- 40% requests will be satisfied almost immediately
- 60% requests satisfied by origin server
- utilization of access link reduced to 60%, resulting in negligible delays (say 10 msec)
- total avg delay = Internet delay + access delay + LAN delay
= $.6 * (2.01 \text{ secs}) + .4 * \text{milliseconds} < 1.4 \text{ secs}$

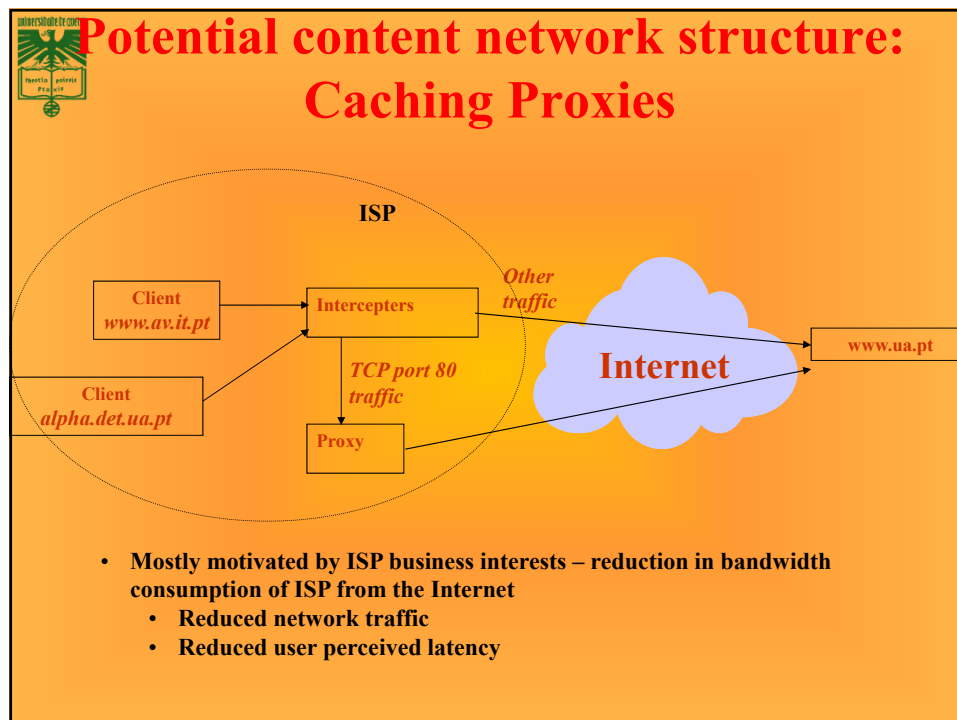
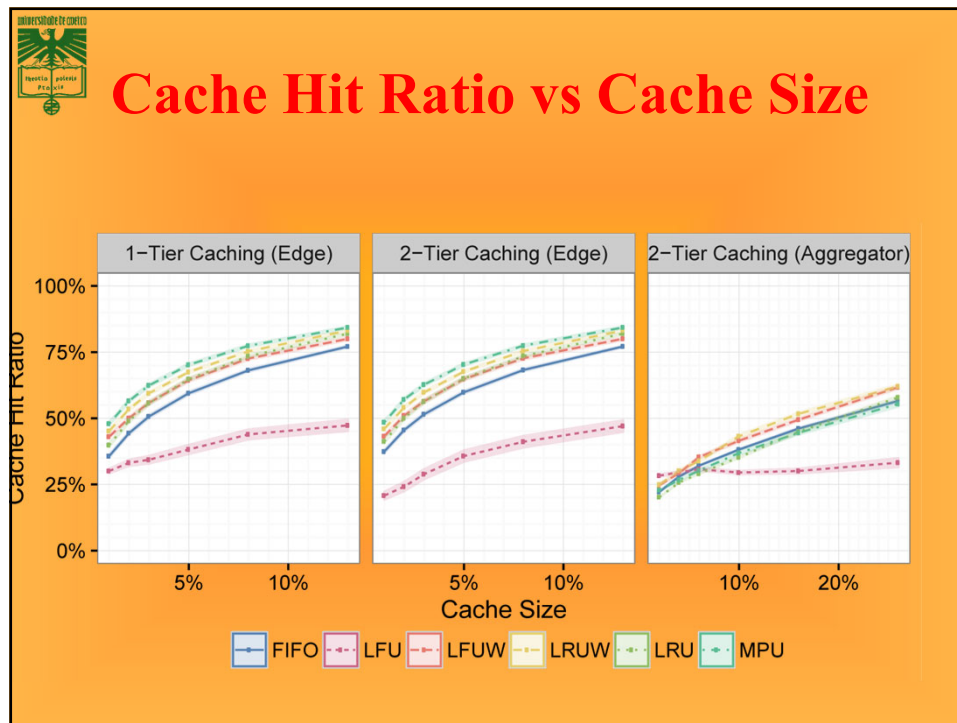


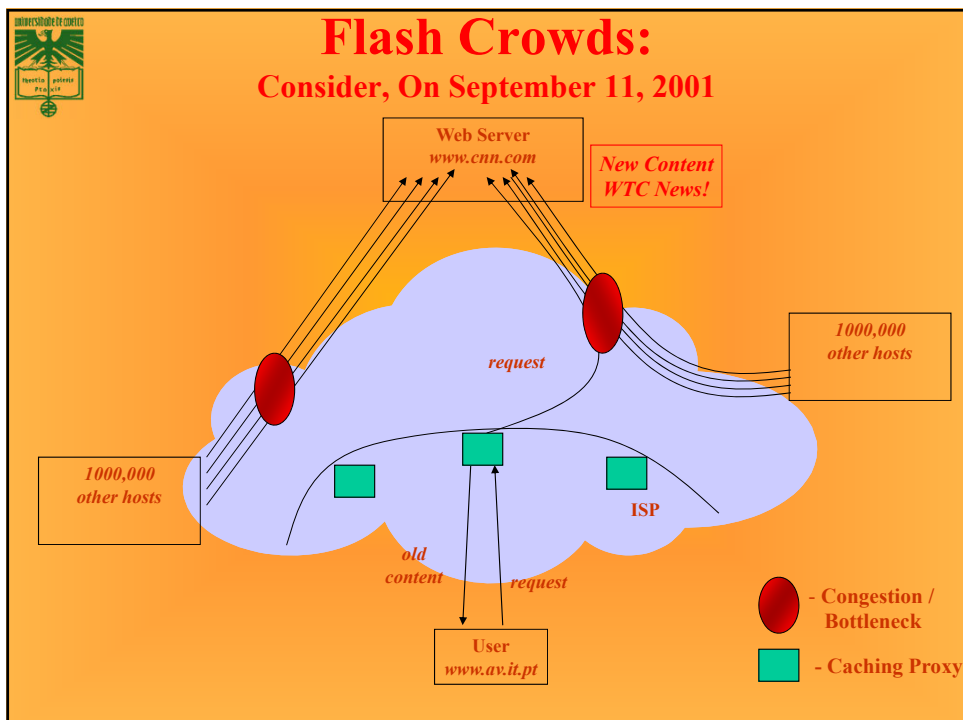
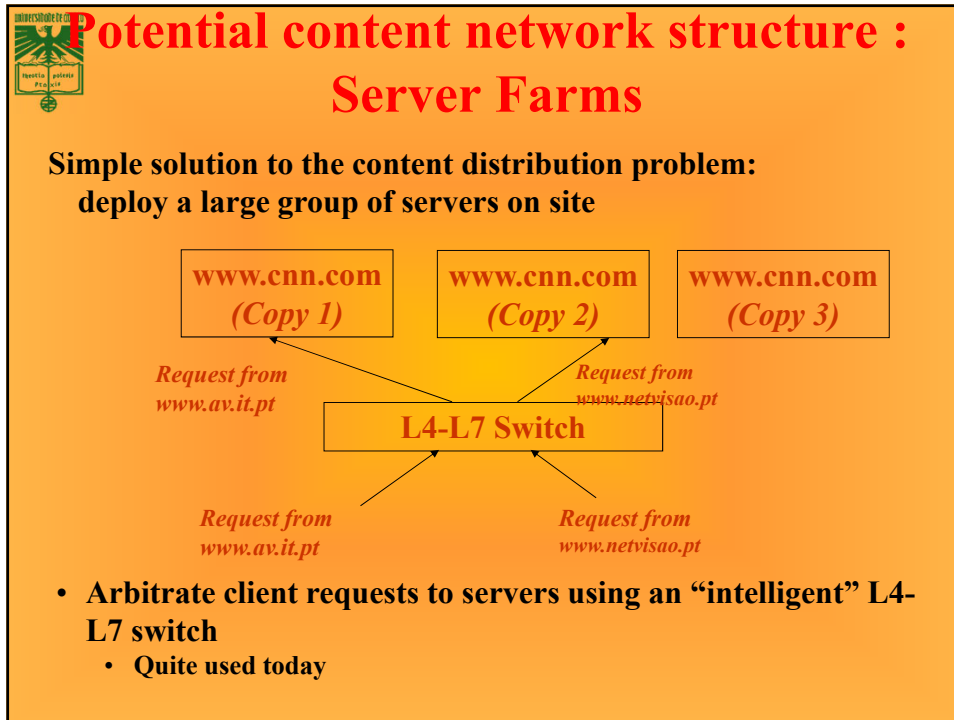




 **Cache Replacement Policies (II)**

- **LFU: Least Frequently Used**
- **MPU: Most Probably Used**
- **LFU and LRU weighted (give a weight to each page)**







21

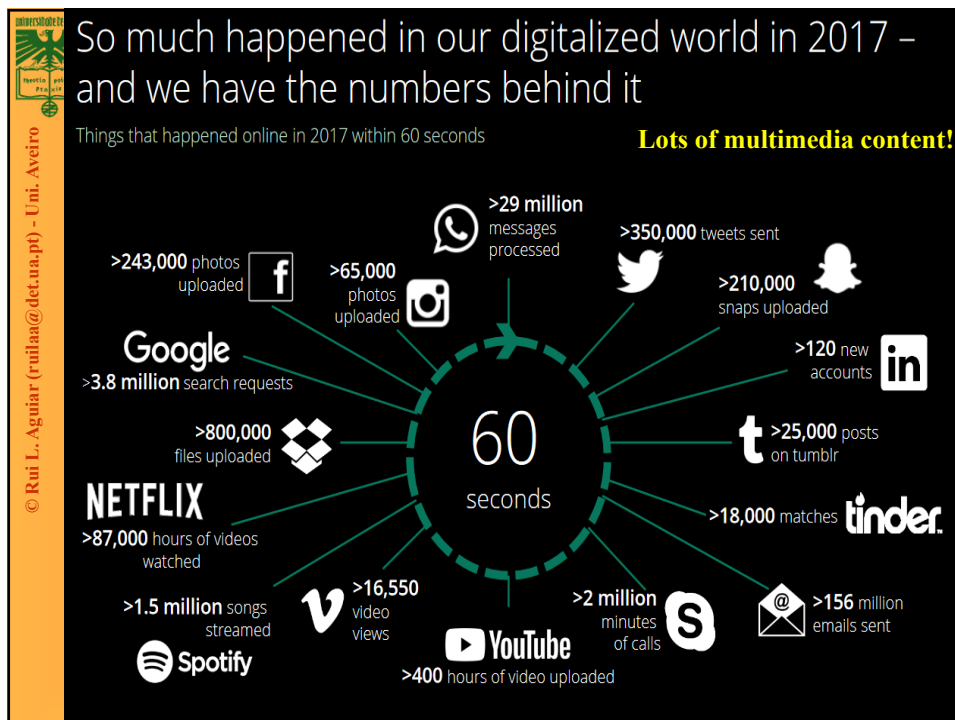
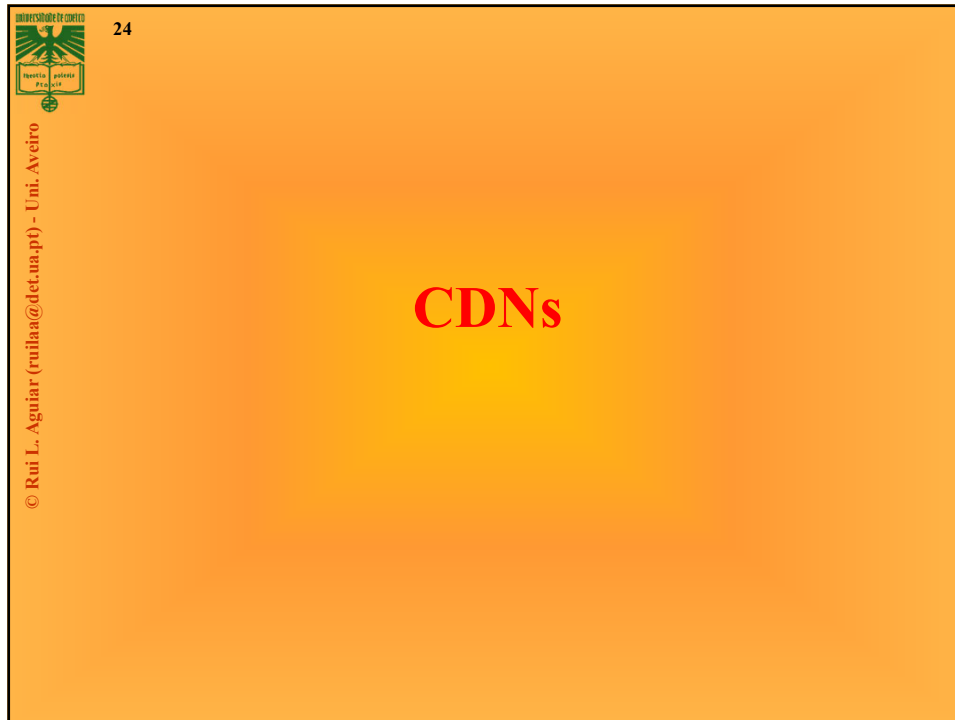
Why Not Web-only approaches for content networks?

- **Integrating file caching in proxies**
 - Optimized for 10KB objects
 - $10\text{GB} = 1.000.000 \times 10\text{KB}$
- **Memory pressure**
 - Disk access is 1000 times slower
 - Working sets do not fit in memory
- **Waste of resources**
 - More servers needed
 - Provisioning is a must



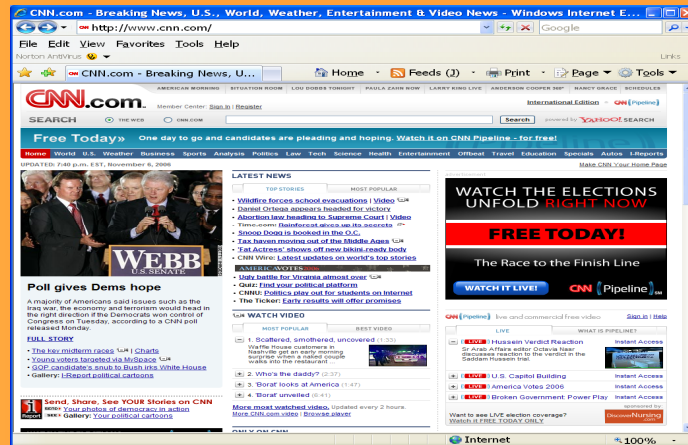
Problems with *Server farms and Caching proxies*

- Server farms do nothing about problems due to network congestion, or to improve latency issues due to the network
- Caching proxies serve only their clients, not all users on the Internet
- Content providers (*say, Web servers*) cannot rely on existence and *correct* implementation of caching proxies
- Accounting issues with caching proxies.
For instance, *www.cnn.com* needs to know the number of hits to the webpage for advertisements displayed on the webpage



27

CDNs Target environment?




**Most Web files are small (1KB ~ 100KB)
(initially....)**

28


Motivation

- IP based networks
- Web based applications have become the norm for corporate internal networks and many business-to-business interactions
- Large acceptance and explosive growth
 - Serious performance problems
 - Degraded user experience
- Improving the performance of networked applications
 - For a large set of applications, including VIDEO access
 - Use many sites at different points within the network
 - Stand alone servers
 - Routers



CDNs basics

- **What is a CDN?**
 - A network of servers delivering content on behalf of an origin site
 - A number of CDN companies well established now
 - E.g. Akamai, Digital Island, Speedera, CDN77, Cloudflare, Stackpat
 - Many companies are exploring CDNs
 - Avoid congested portions of the Internet
- **Consist of**
 - Edge servers deployed at several ISP (Internet Service Provider) access locations and network exchange points
- **Large-file service with no custom client, no custom server, no prepositioning**
- **Improve the response time of an Internet site**
 - Offloading the delivery of bandwidth-intensive objects, such as images and video clips
- **Intelligent Internet infrastructure that improves the performance and scalability of distributed applications by moving the bulk of their *computation* to servers located at the edge of the network**
 - Applications are logically split into two components
 - Executed at an edge server close to the user
 - Executed on a traditional application server

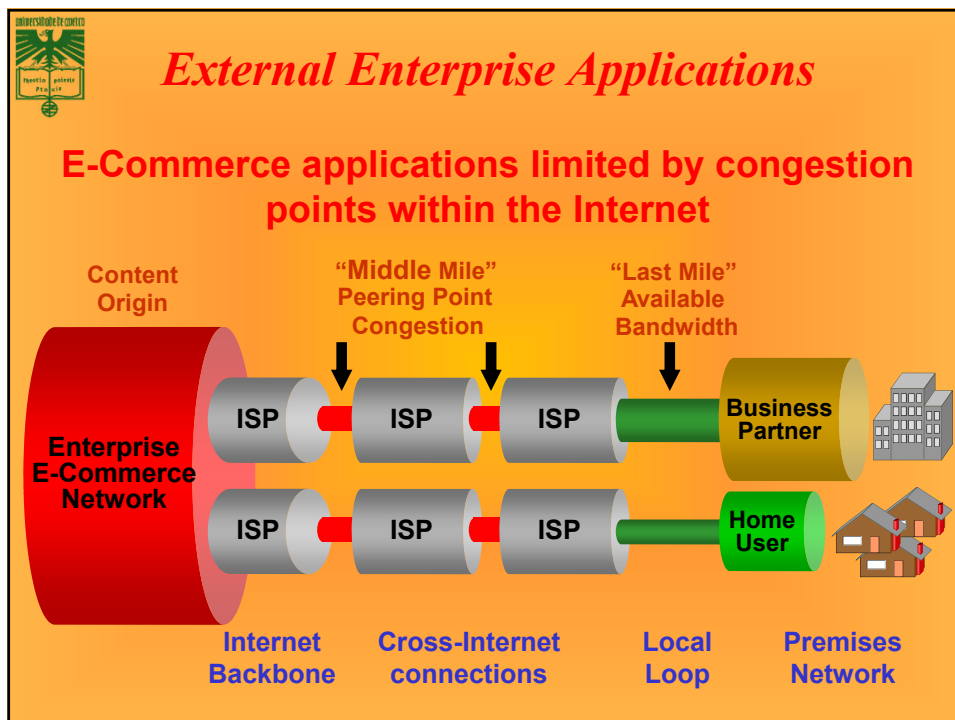
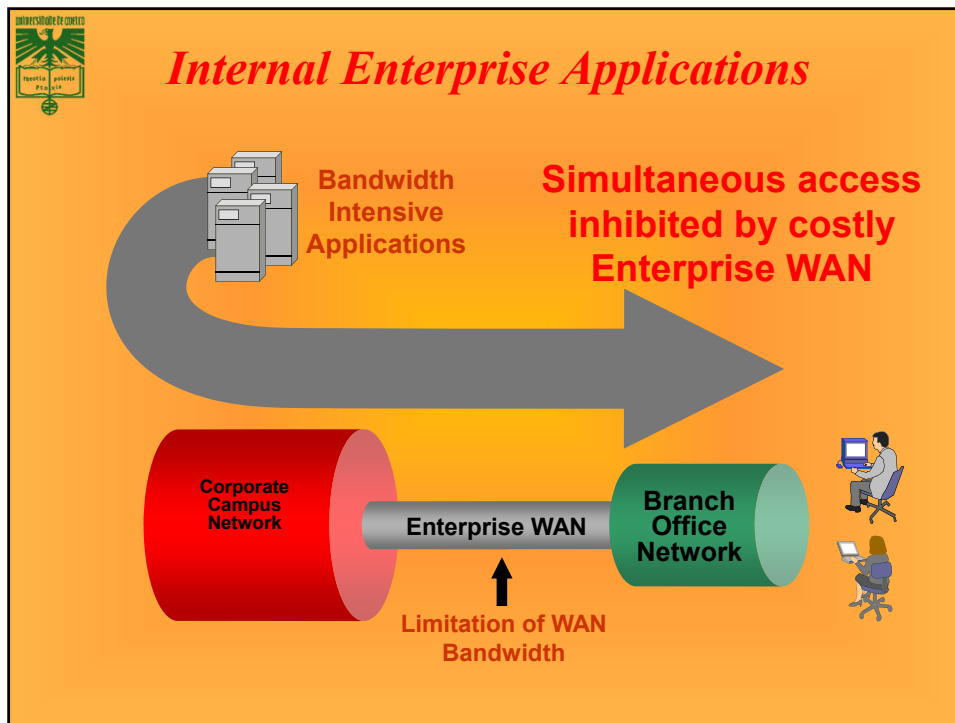


© Rui L. Aguiar (rui.aa@ua.pt) - Uni. Aveiro

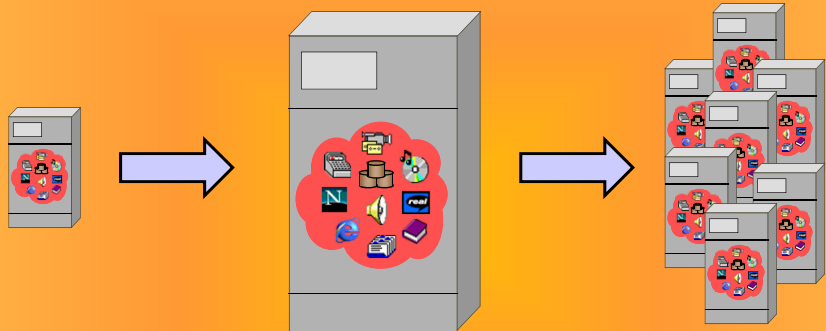
30

CDN Generations

- **First generation (early 90ies)**
 - Accelerate the performance of web sites
 - Support increasing volumes of traffic
 - Key disruption event: 9/11
 - Akamai technologies created
- **Second generation (early 2000ies)**
 - Support high volumes of multimedia traffic
 - Audio/video intensive networks
 - All ISPs developed/used CDNs
- **Third generation (2010+)**
 - Cloud computing
 - Amazon cloud (2008)
 - UGC (user generated content)
 - P2P and interactivity
 - AT&T distributed data centers (2011)
 - Mobile support, and device adapted content



Content Scaling

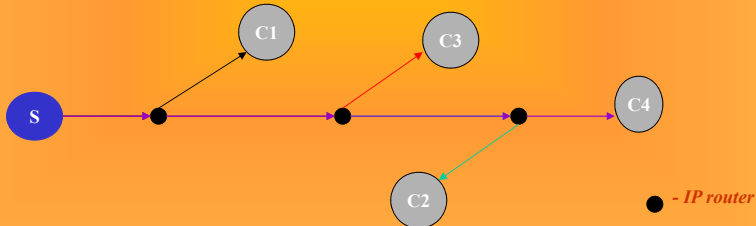


- **Need to scale content to handle numerous clients**
 - One can only scale 'vertically' to a point

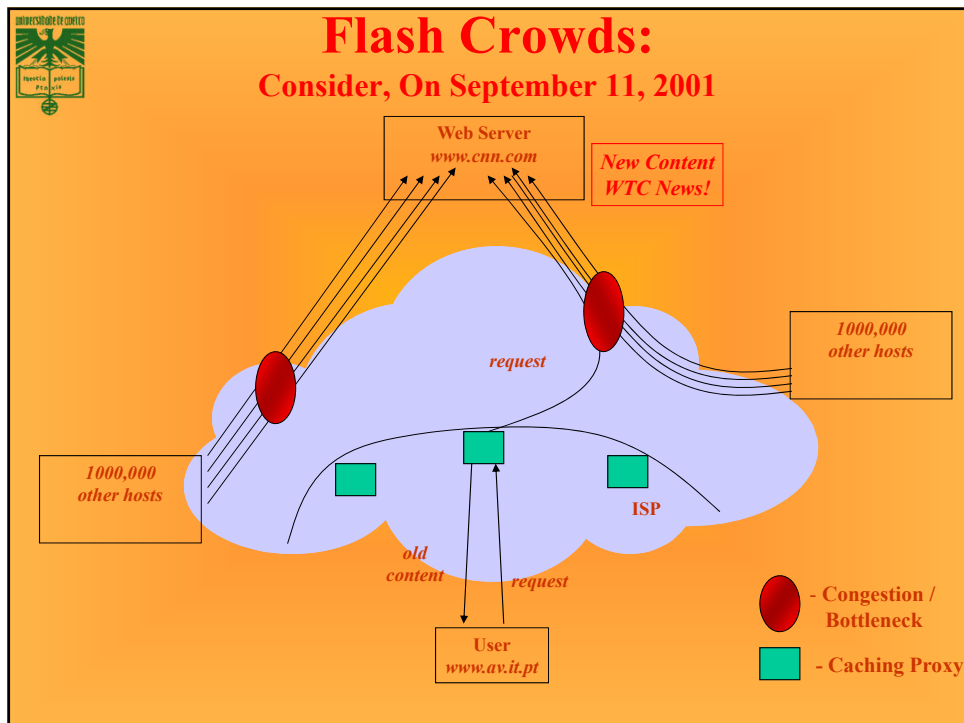
Multiple servers/locations introduces new issues

Early Motivations for Content Networks (1st generation)

- More hops between client and Web server => more congestion!
- Same data flowing repeatedly over links between clients and Web server
- Origin server is bottleneck as number of users grows
- Flash Crowds (*for instance, Sept. 11*)
 - *The Content Distribution Problem:* Arrange a rendezvous between a content source at the origin server (*www...com*) and a content sink (*users*)



● - IP router



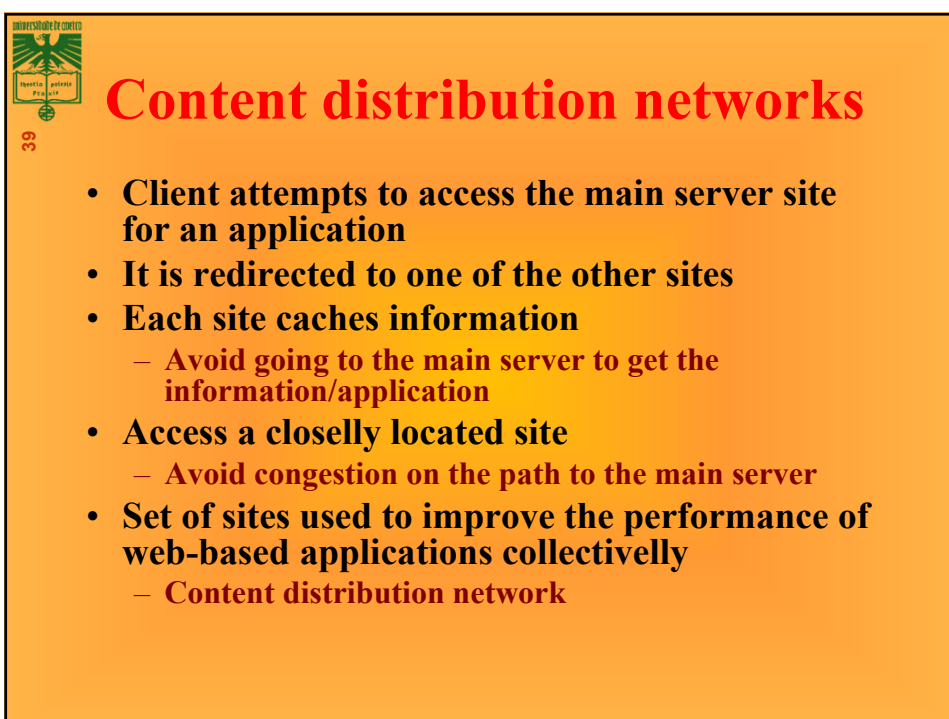
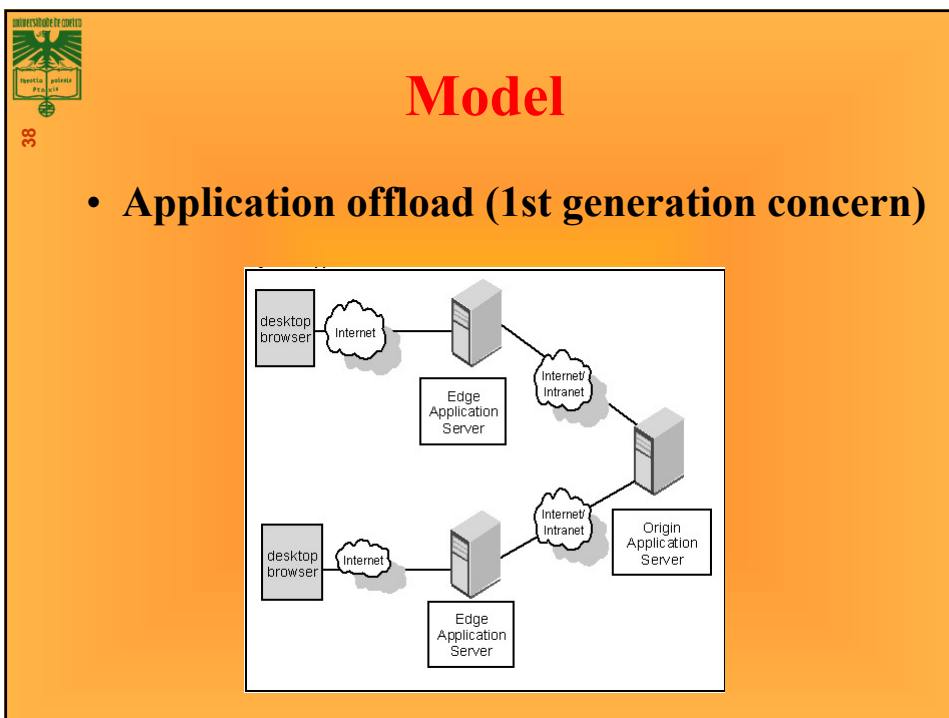
37

Flash crowd solution: CDNs..

What is a CDN?
A network of servers delivering content on behalf of an origin site

Large-file service with

- No custom client
- No custom server
- No prepositioning
- No rehosting
- No manual provisioning





Inside a CDN

- **Servers are deployed in clusters for reliability**
 - **Some may be offline**
 - Could be due to failure
 - Also could be “suspended” (e.g., to save power or for upgrade)
- **Could be multiple clusters per location (e.g., in multiple racks)**
- **Server locations**
 - **Well-connected points of presence (PoPs)**
 - **Inside of ISPs**



41

Advantages

- **Better scalability**
 - **Higher availability**
 - **Improved response time from a centrally managed solution**
 - **Nodes constituting the distribution network are designed to be**
 - **Self-configuring**
 - **Self-managing**
 - **Self-diagnosing**
 - **Self-healing**
- to ensure easy management and operational convenience**



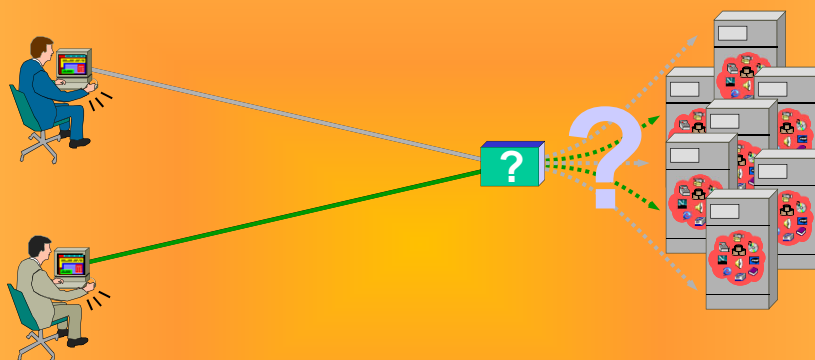
42

Challenges

- **Keep consistency among the enterprise data hosted by the offloaded applications**
- **Share session state among edge and origin application servers**
- **Distribution, configuration, and management**
- **Develop programming models consistent with current industry standards such as J2EE**
- **Application security.**
- **There is active research into general frameworks to be used to support distributed applications, as well as prototyping the ideas for specific application instances**

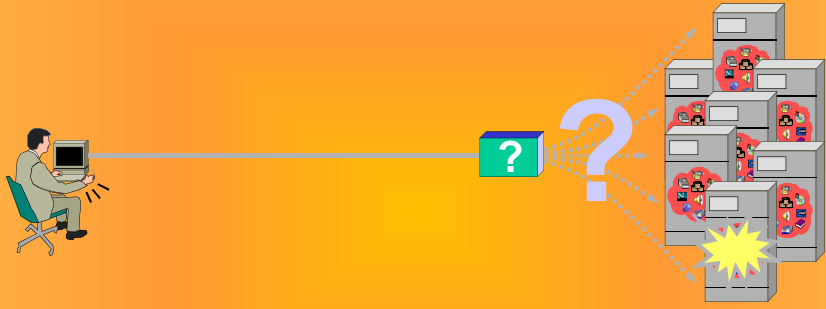


Load-Sharing Content



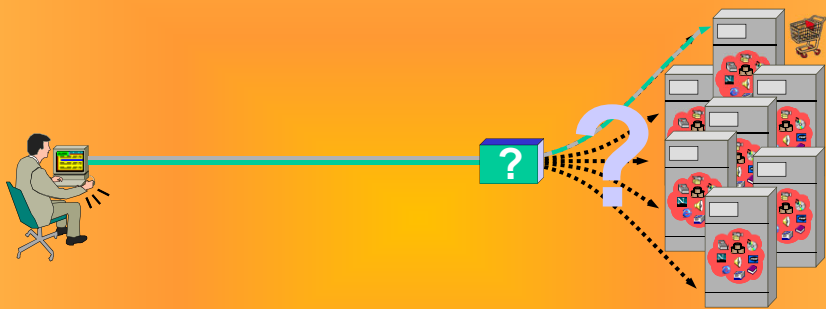
- **Handle requests fairly amongst servers/sites**
- **Easily add servers/sites to content service**
- **Adjust connections based on server/site load**

Content Availability with multiple servers?



- Synchronize content amongst servers/sites
- Avoid faulty servers/sites
- Faulty servers/sites includes invalid/dated content

Persistence with multiple servers?



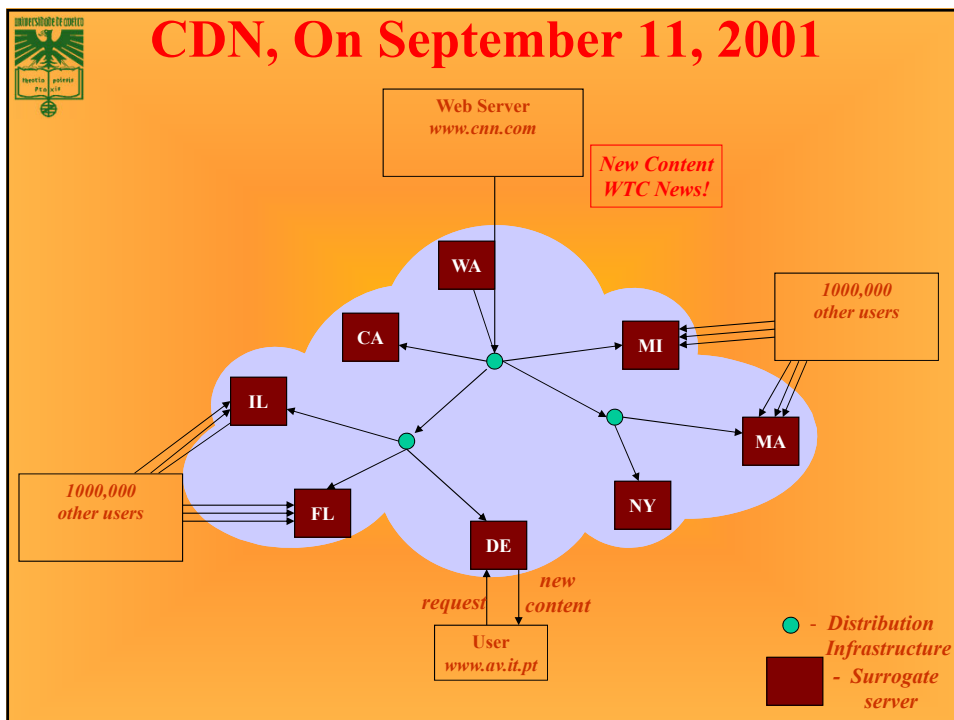
- Handle applications which use 'state'
 - Need to learn client ID to satisfy state requirement
 - Need to maintain state for period of time - variable

46

© Rui L. Aguiar (rui.laa@det.ua.pt) - Uni. Aveiro

Outline

- Overall context
- Challenges
- Potential alternatives?
- Architecture





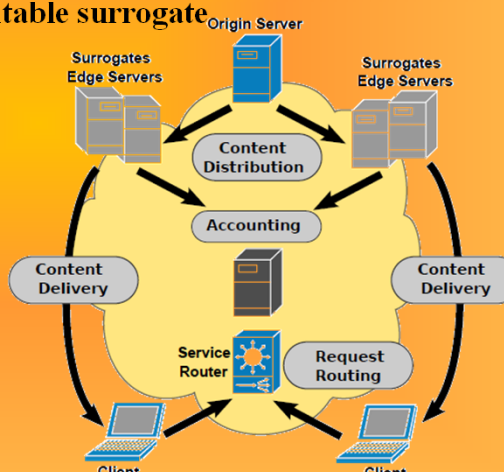
With CDNs

- **Overlay network to distribute content from origin servers to users**
 - Avoids large amounts of same data repeatedly traversing potentially congested links on the Internet
 - Reduces Web server load
 - Reduces user perceived latency
 - Tries to route around congested networks
- **CDN is not a cache!**
 - Caches are used by ISPs to reduce bandwidth consumption, CDNs are used by content providers to improve quality of service to end users
 - Caches are reactive, CDNs are proactive
 - Caching proxies cater to their users (web clients) and not to content providers (web servers), CDNs cater to the content providers (web servers) and clients
 - CDNs give control over the content to the content providers, caching proxies do not



CDN Components

- **Content Delivery Infrastructure:** Delivering content from producer to clients by surrogates
- **Request Routing Infrastructure:** Steering or directing content request from a client to a suitable surrogate
- **Distribution Infrastructure:** Moving or replicating content from content source (origin server, content provider) to surrogates
- **Accounting Infrastructure:** Logging and reporting of distribution and delivery activities





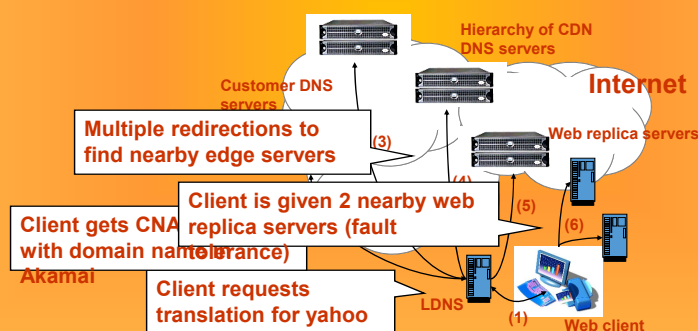
Mapping clients to servers

- **CDNs need a way to send clients to the “best” server**
 - The best server can change over time
 - And this depends on client location, network conditions, server load, ...
 - What existing technology can we use for this?
- **DNS-based redirection**
 - Clients request www.foo.com
 - DNS server directs client to one or more IPs based on request IP
 - Use short TTL to limit the effect of caching



DNS Redirection

- **Web client's request redirected to 'close' by server**
 - Client gets web site's DNS CNAME entry with domain name in CDN network
 - Hierarchy of CDN's DNS servers direct client to 2 nearby servers





DNS Redirection Considerations

• Advantages

- Uses existing, scalable DNS infrastructure
- URLs can stay essentially the same

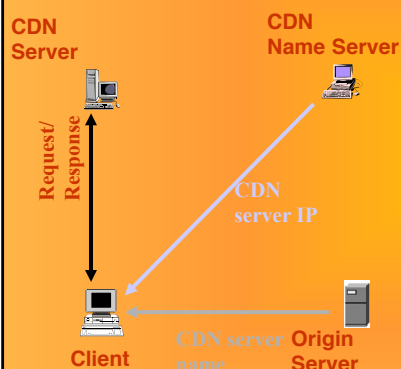
• Limitations

- **DNS servers see only the DNS server IP**
 - Assumes that client and DNS server are close. Is this accurate?
- **Content owner must give up control**
- **Unicast addresses can limit reliability**



53

What other CDN techniques are being used?



- **DNS redirection (DR)**
 - Full-site delivery
 - Partial-site delivery
- **URL rewriting**
- **Hybrid scheme**
 - URL rewriting + DNS redirection
- **Manual hyperlink selection**
- **HTTP redirection**
- **Layer 4 switching**
- **Layer 7 switching**
- **Anycast**



54

Offloading a portal

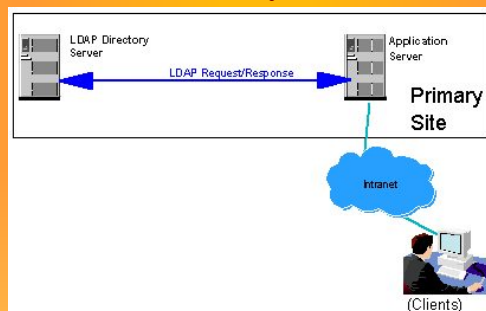
- **Portal servers allow users to access content and applications from a single access point**
 - Users can create persistent, customized views of applications and content chosen from the set of applications and content by the portal administrators
- **Portal server pages are personalized**
- **Often include dynamic content**
- **Significant amount of computation required for page assembly**
 - **Application offload**

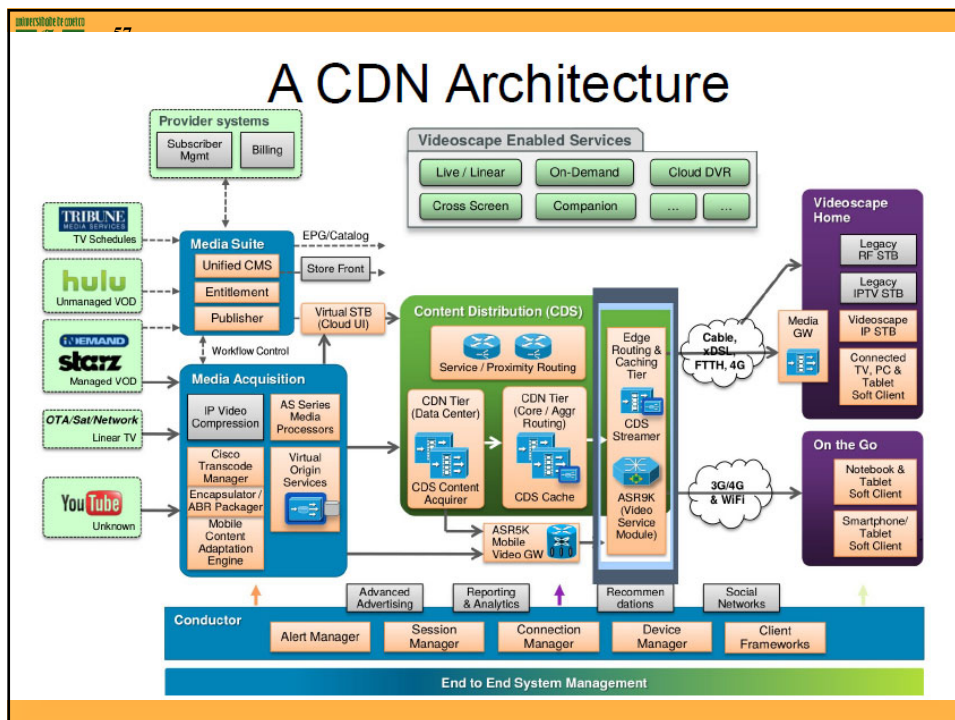
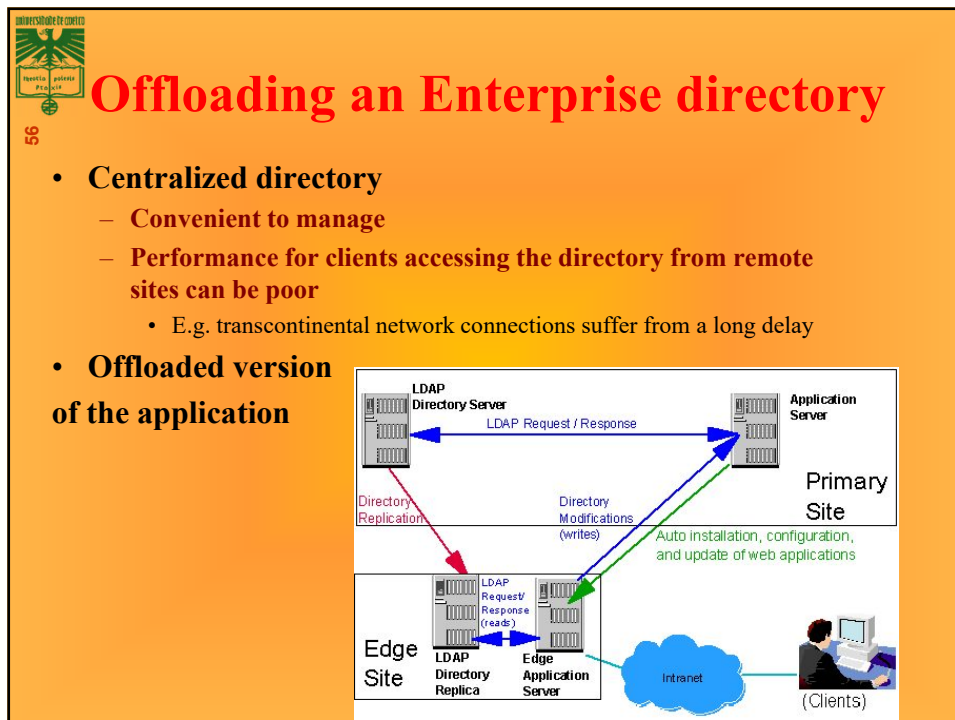



55

Offloading an Enterprise directory

- **E.g. a common e-Workplace tool**
- **The employee data is often stored in a central LDAP directory**
 - **Separate web-based application providing the interface to the directory**








Mobile Networks and the FMBC

The arrival of common services



60

Learning outcomes

- Understand the appearance of mobile networks
- Perceive the technologies underlying the integration of services
- Realize the triple and quad-play concepts.

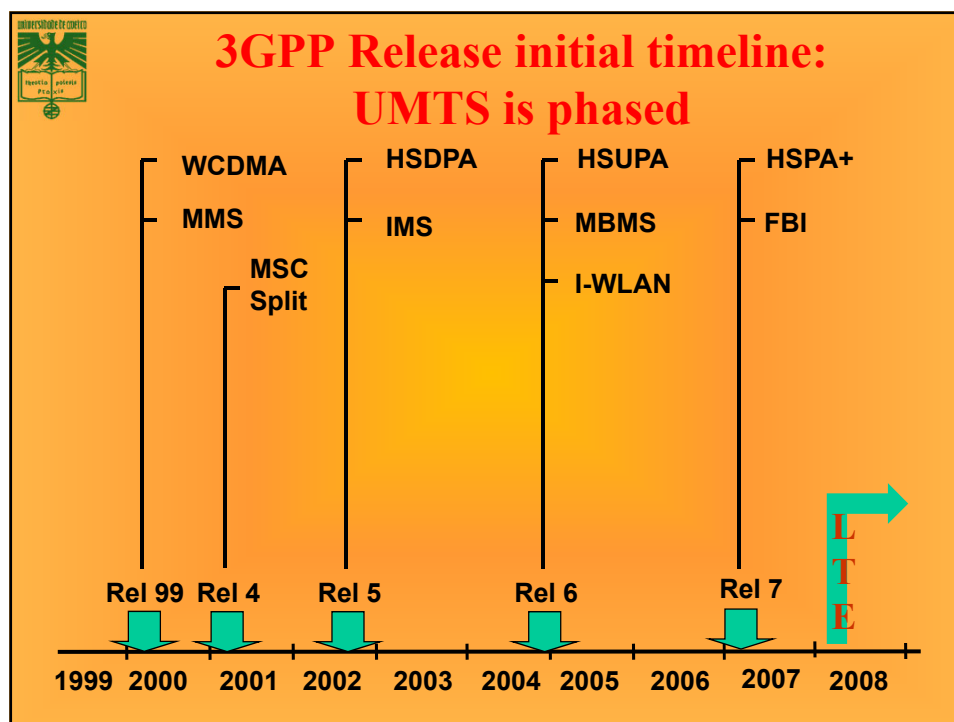
© Rui L. Aguiar (rui.laa@det.ua.pt) - Uni. Aveiro

Wide Communication technologies: Cellular

Comparison	2G	3G	4G	5G
Introduced in year	1993	2001	2009	2018
Technology	GSM	WCDMA	LTE, WiMAX	MIMO, mm Waves
Access system	TDMA, CDMA	CDMA	CDMA	OFDM, BDMA
Switching type	Circuit switching for voice and packet switching for data	Packet switching except for air interface	Packet switching	Packet switching
Internet service	Narrowband	Broadband	Ultra broadband	Wireless World Wide Web
Bandwidth	25 Mhz	25 Mhz	100 Mhz	30 GHz to 300 GHz
Advantage	Multimedia features (SMS, MMS), internet access and SIM introduced	High security, international roaming	Speed, high speed handoffs, global mobility	Extremely high speeds, low latency
Applications	Voice calls, short messages	Video conferencing, mobile TV, GPS	High speed applications, mobile TV, wearable devices	High resolution video streaming, remote control of vehicles, robots, and medical procedures

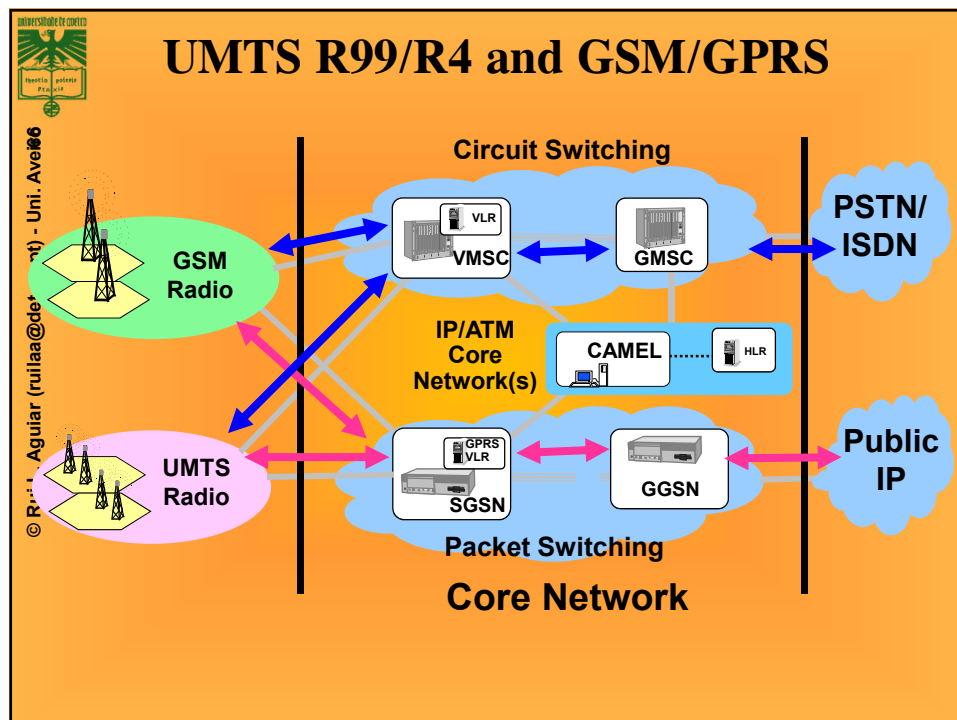
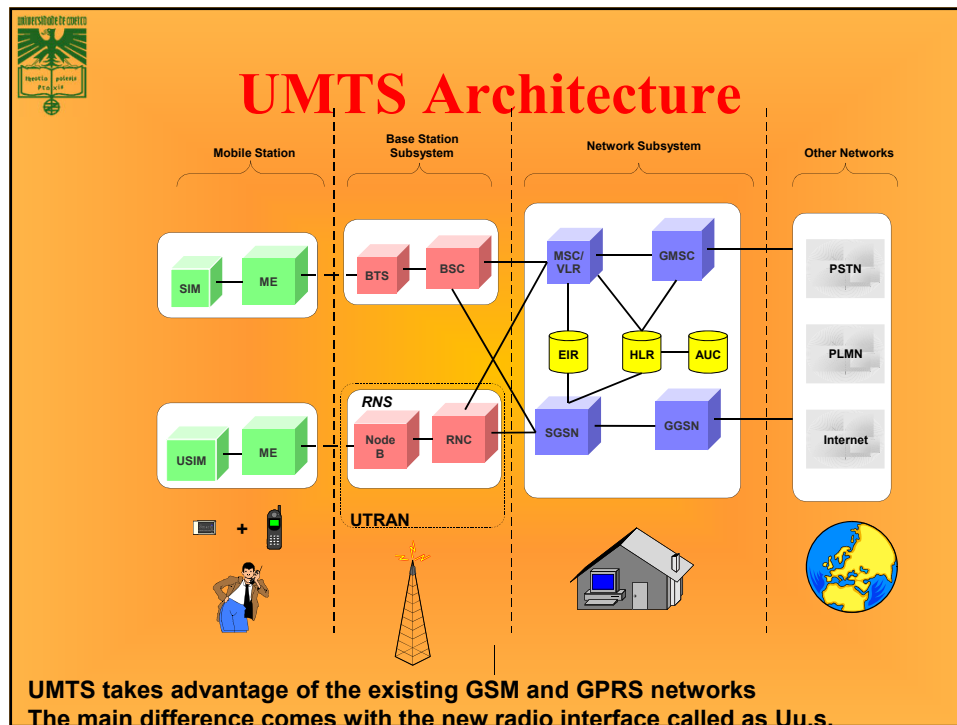
Early cellular systems

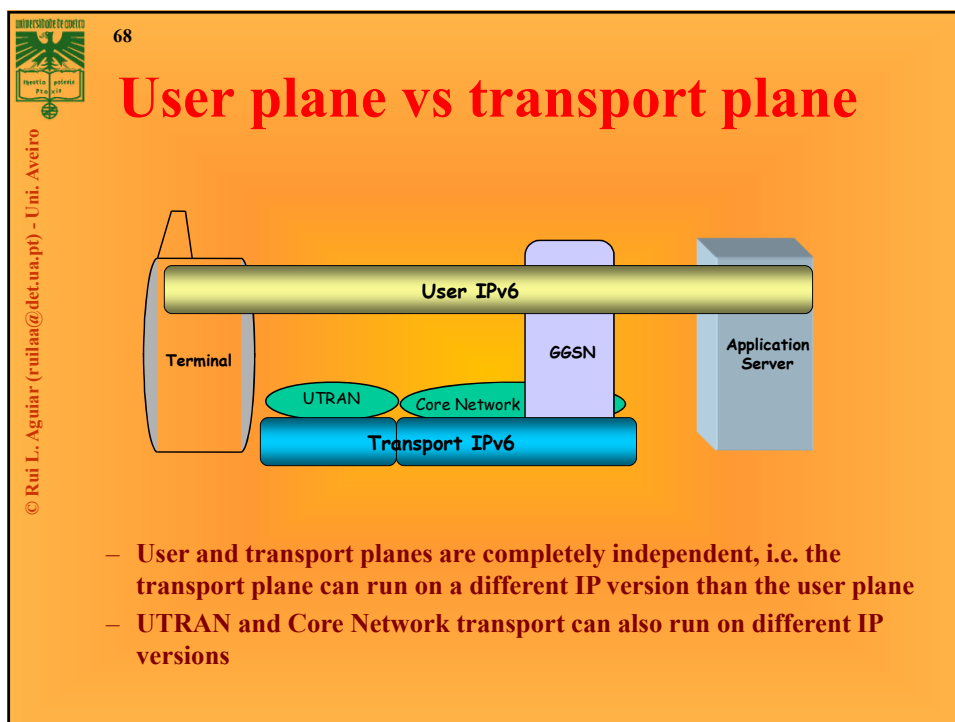
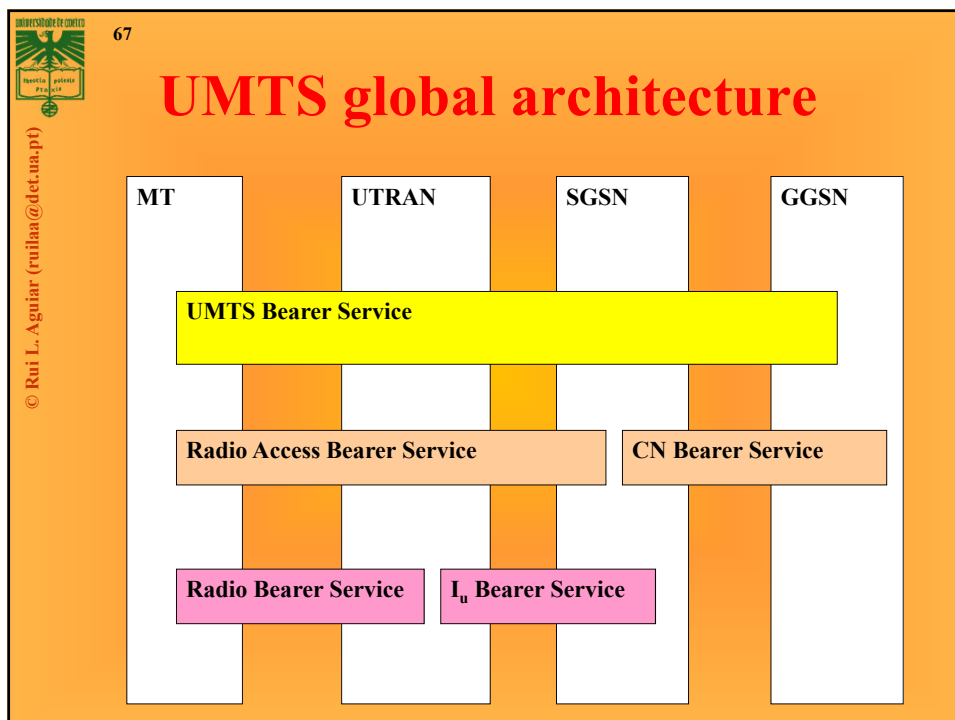
- **1G: analog systems (450-900 MHz)**
 - Signalling: FSK
 - Share of medium: FDMA
 - NMT (Europe), AMPS (US)
- **2G: digital systems (900, 1800, 1900 MHz)**
 - Share of medium : TDMA/CDMA
 - Circuit switching
 - GSM (Europe), IS-136 (US), PDC (Japan)
- **2.5G: extensions for packet switching**
 - Digital: GSM → GPRS
 - Analog: AMPS → CDPD
- **3G: networks for data applications**
 - High rates, data, Internet
 - Share of medium : TDMA/CDMA/CDMA
 - IMT-2000 (Europe: UMTS)

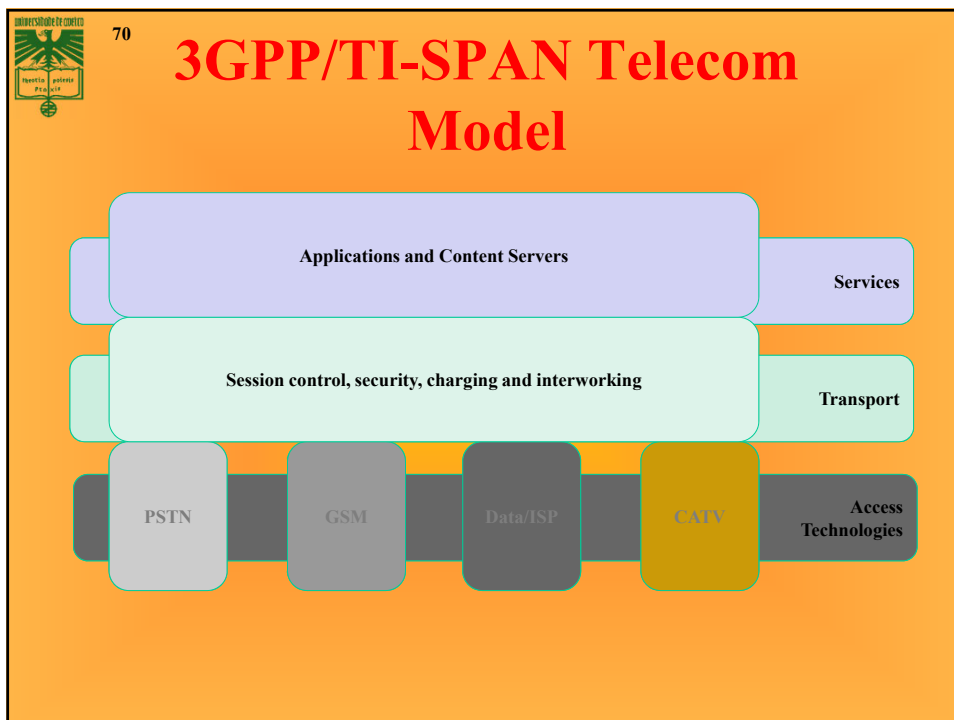
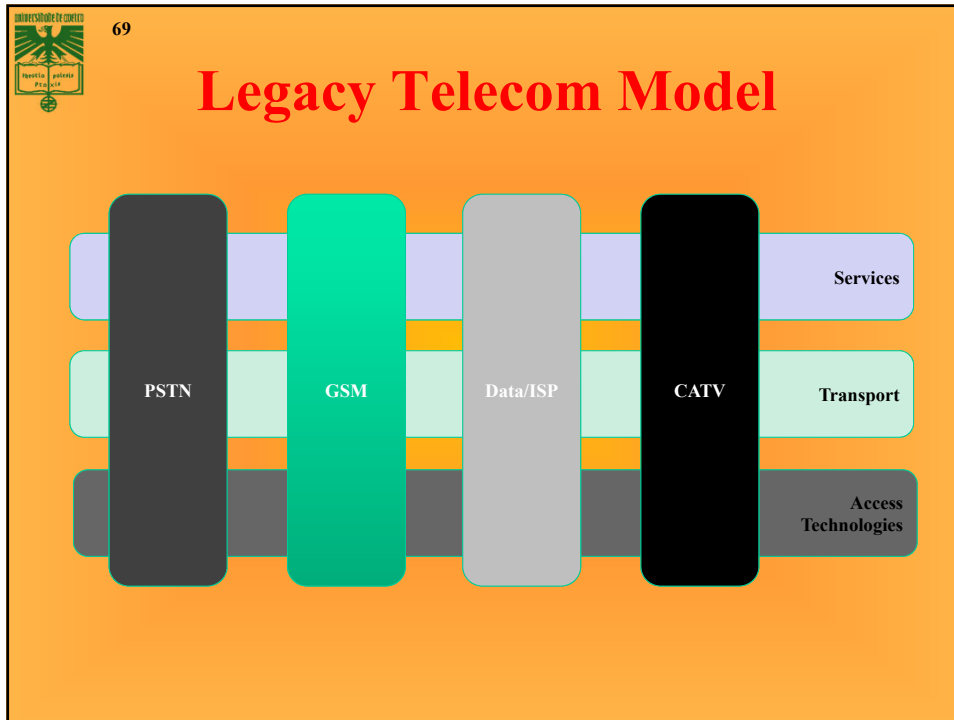


UMTS: first universal cellular data system

- 3G system
- Oriented to generalized service diffusion and its future users trends
 - Combines cellular, wireless, paging, etc. functions
- “multimedia everywhere”
- Developed as an evolution path of 2.5G systems
 - Progressive evolution (GPRS-EDGE-UMTS)
- (Initial) Data rates of UMTS were:
 - 144 kbps for rural
 - 384 kbps for urban outdoor
 - 2048 kbps for indoor and low range outdoor
 - Large rates later, progressively increased









71

Implications

- Any Device
- Any Access Technology
- Any Where

ALWAYS BEST CONNECTED

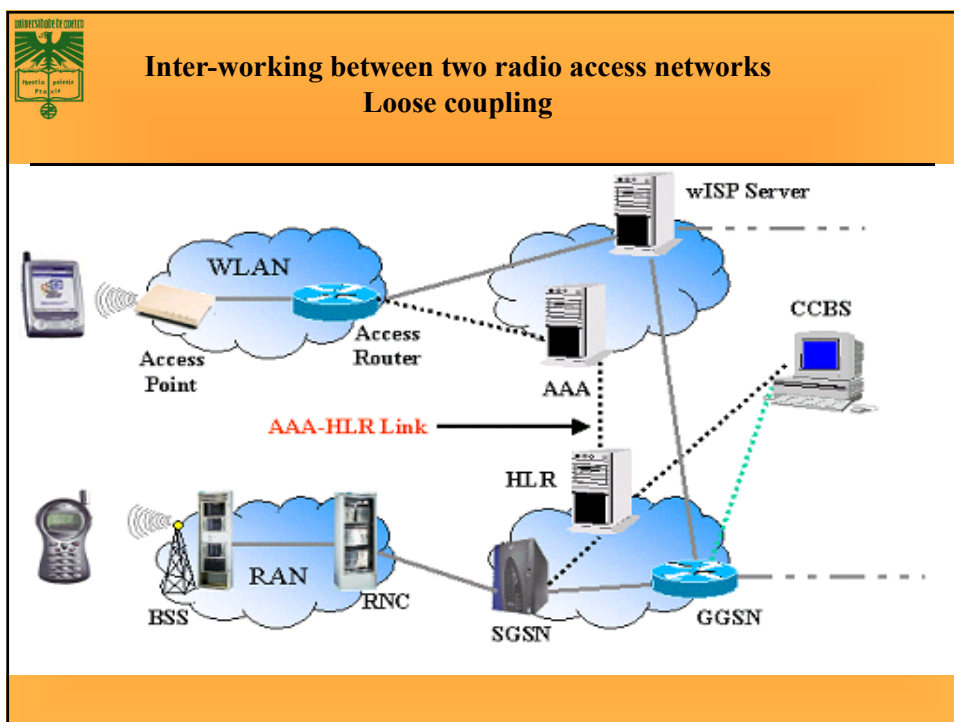
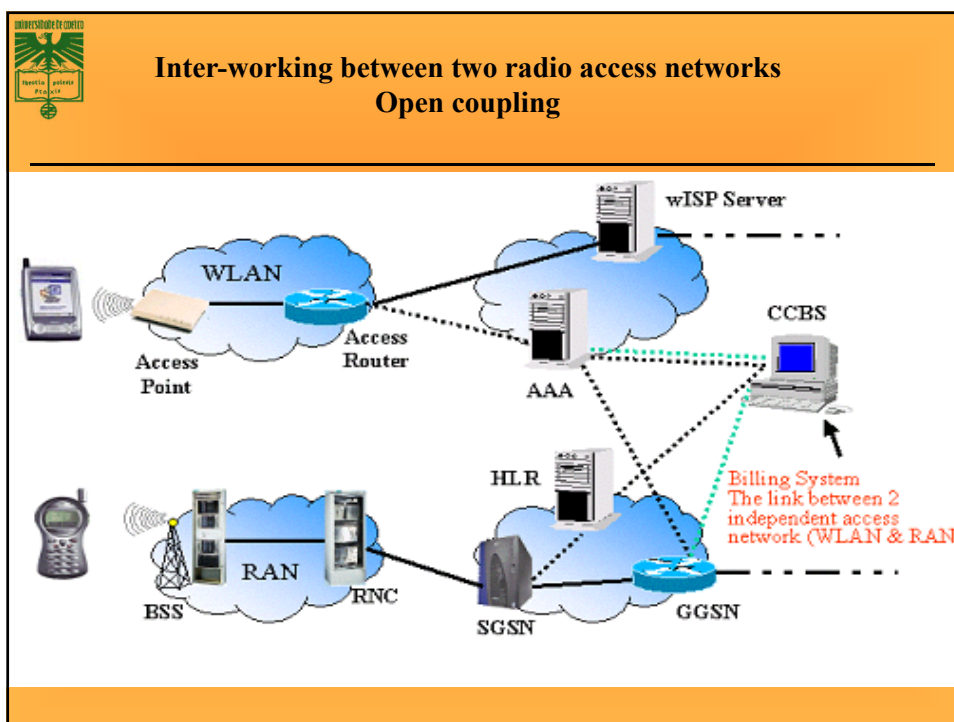
- One Network, multiple access technologies
- Common Session Control
- Generic Application Servers
- Single set of services that apply network wide
- Consistent user experience
- Operational efficiency
- New services/applications

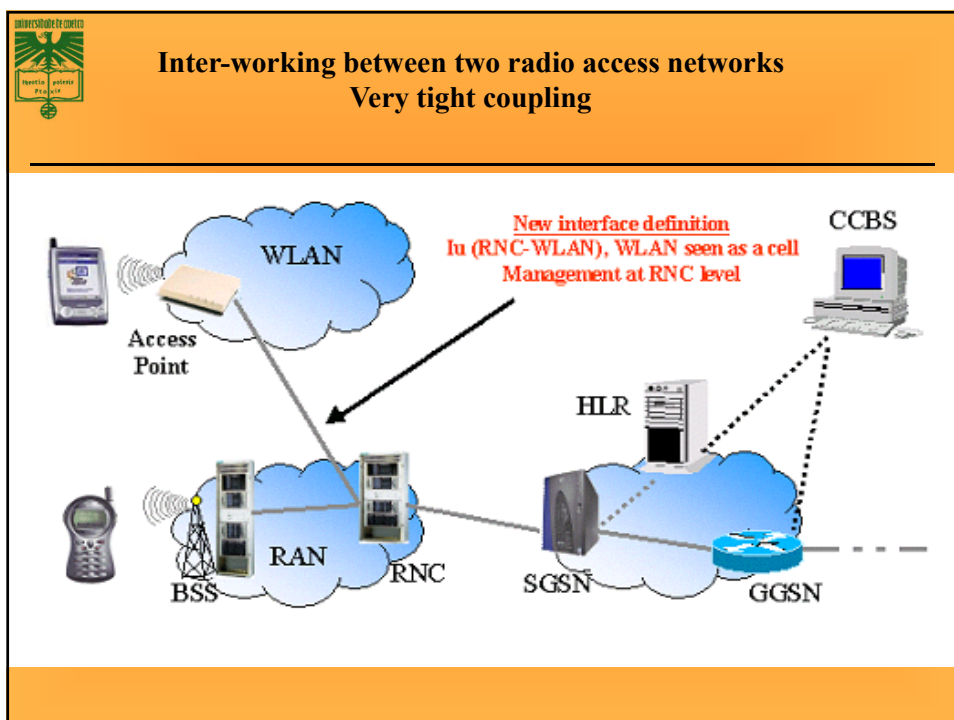
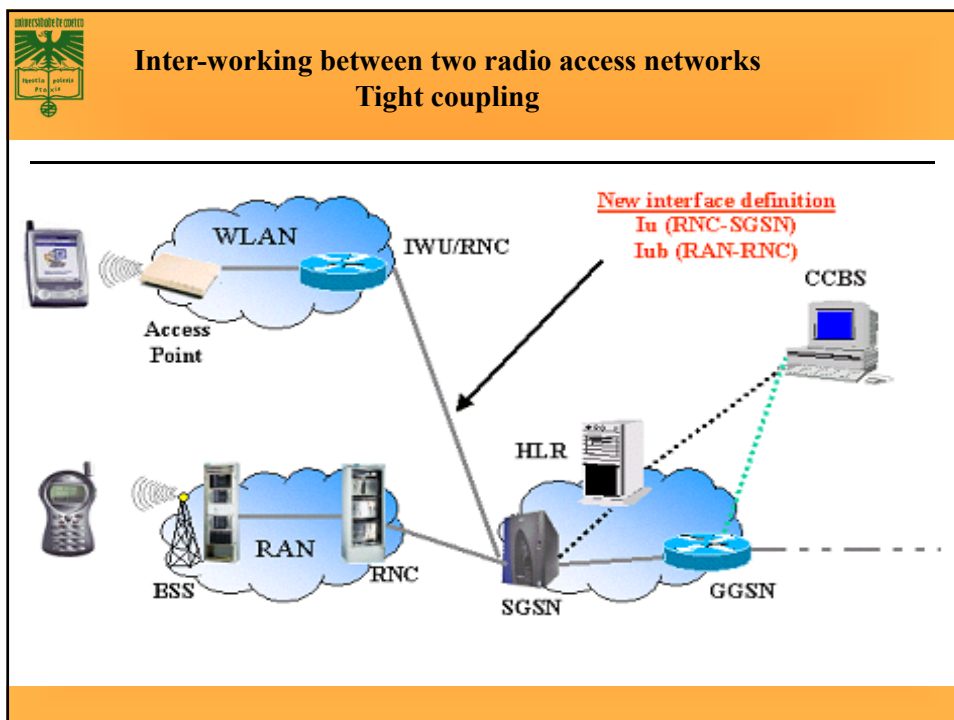


Fixed Mobile Broadcast Convergence (FMC)

- One customer service
 - Handles mobile and fixed calls
 - Any network — mobile, WiFi, Broadband Cable...
 - Avoid mobile charging when in-building
- Single (customer) number with common suite of services
 - One voice mailbox, one phone directory...
 - Mobile, fixed, conference room
- New services? Irrespective of location, access technology or terminal device
 - Potentially gradative provision

Slide 72







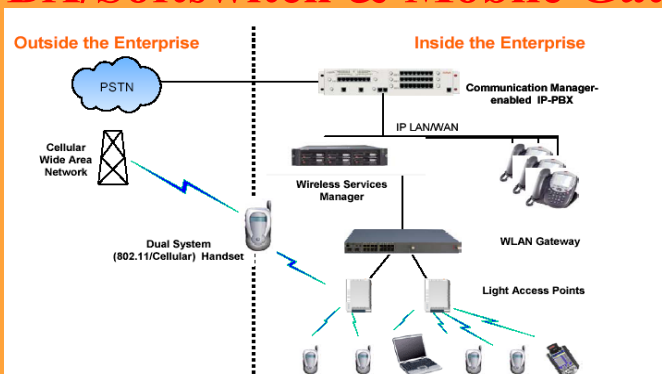
Implementing the FMBC concepts

- **Wireless “fixed” line services**
 - New (not FMBC) in developing nations, mobile, no handoffs
- **IP-PBX or softswitch with mobile network interface**
 - Centered in company internal communications
- **Unlicensed Mobile Access (UMA)**
 - GSM & GPRS services over WiFi or Bluetooth
- **Mobile VoIP technology (pre-IMS)**
 - emulating mobile network entities
- **IP Multimedia Subsystem (IMS)**
 - 3G vision of future IP-based mobile communications

Slide 77



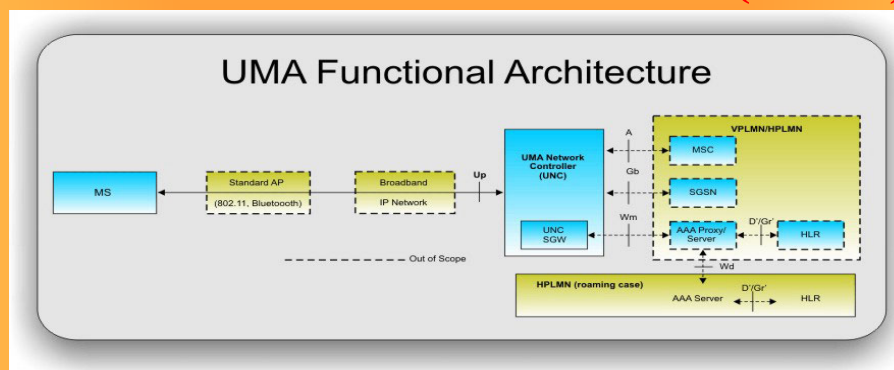
IP-PBX/Softswitch & Mobile Gateway



- IP-PBX or softswitch is in charge -> directed to company contacts
- Service hands off to mobile network when out of WLAN range
- Mobile network only used when necessary
 - “Tunnel” through the mobile network
 - Voice or text messages generated when necessary to connect enterprise calls to remote users



Unlicensed Mobile Access (UMA)



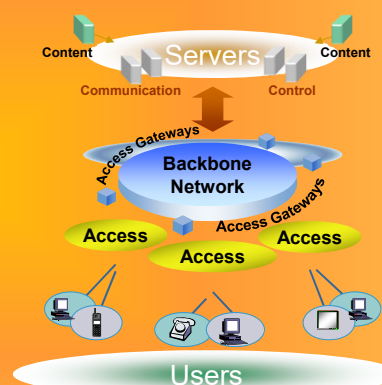
- Seamless delivery (roaming and handover) of voice (GSM) and data (GPRS) over wireless IP networks
- Works with 2.5G+ mobile networks
 - Tunnels GSM & GPRS over IP to mobile core network
 - No impact to operations of cellular RAN



80

UMTS from release 5 on: IMS: IP Multimedia Subsystem

- Same Core network
- Same User on different accesses
- Same Services
- Can use WLAN, ADSL, LAN, UTRAN (GPRS) etc. as accesses in ONE system
- Can have several devices and move between them





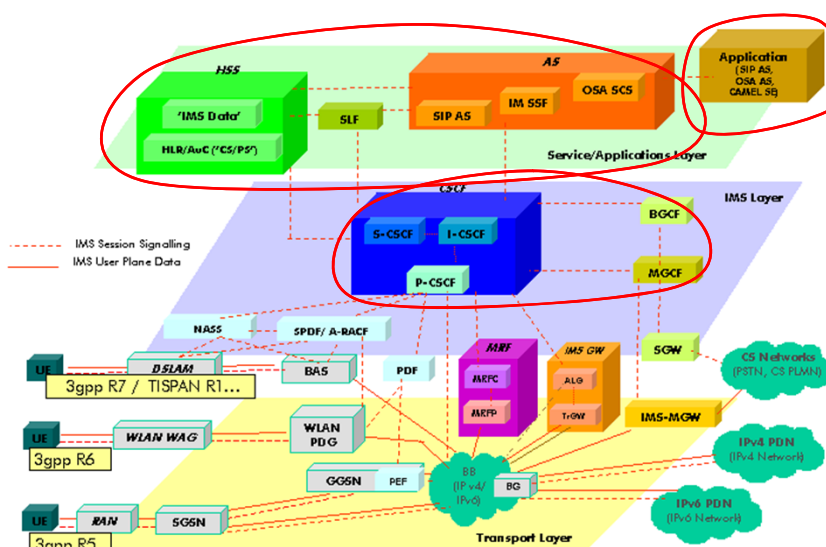
IP Multimedia Subsystem (IMS)

- New IP-based mobile core network for 3G evolution
- Uses 3GPP variant of SIP & other IP protocols
- “Intelligent Network” over IP?
- New services drive IMS deployment
 - Push-to-Talk, FMC, IP Centrex
- PTT (PoC) & UMA FMC specs already turned over to 3GPP
- Developed by 3GPP for GSM-to-3G evolution
 - Defined in release 5; fully specified in release 6



83

Where is IMS ?

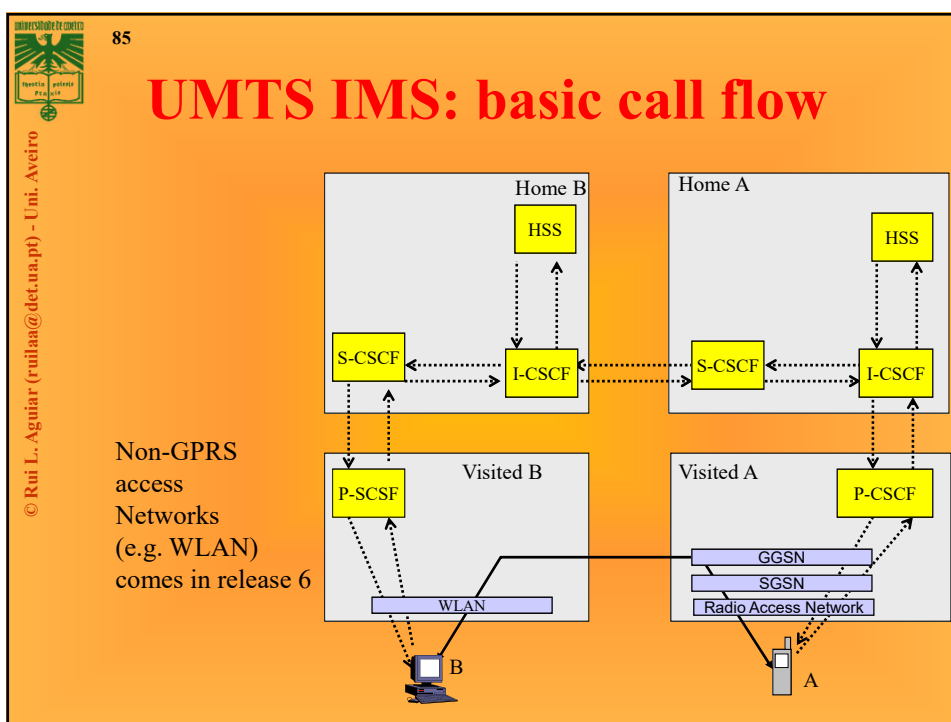


3GPP IMS - Eng. de Serviços - Rui Aguiar/Diogo Gomes <dgomes@ua.pt>

84

x-CSCF

- **Serving - CSCF**
 - Controls the user's SIP Session
 - very few per domain
 - Located in the home domain
 - Is a SIP registrar (and proxy)
- **Proxy-CSCF**
 - IMS contact point for the user's SIP signaling
 - Several in a domain
 - Located in the visited domain
 - Terminals must know this proxy (e.g. DHCP used)
 - Compresses and decompresses SIP messages
 - Secures SIP messages
 - Assures correctness of SIP messages
- **Interrogating – CSCF**
 - domain's contact point for inter-domain SIP signaling
 - one or more per domain
 - In case there are more than one S-CSCFs in the domain, locates which S-CSCF is serving a user





86

SIP Protocol

- **Defined in IETF RFC 3261**
 - “... an application-layer control (signaling) protocol for creating, modifying, and terminating sessions with one or more participants. These sessions include Internet telephone calls, multimedia distribution, and multimedia conferences.”
- **In IMS, SIP is modified to include extra functionality and support a specific set of functions only**
 - SIP is to the Internet what SS#7 is to telephony
- **At the core of IMS there are several SIP proxies:**
 - I-CSCF, S-CSCF, P-CSCF
 - The Call Session Control function (CSCF) is the heart of the IMS architecture
 - The main functions of the CSCF:
 - provide session control for terminals and applications using the IMS network
 - secure routing of the SIP messages,
 - subsequent monitoring of the SIP sessions and communicating with the policy architecture to support media authorization.
 - responsibility for interacting with the HSS.



87

IMS Identity and User Profiles

- **IMS uses SIP identity: SIP URIs**
 - e.g. sip:ruilaa@ua.pt
 - Opposed to phone numbers
 - A user is uniquely identified in the HSS by his IMPI (Private User Identity).
 - IMPI is a unique global identity defined by the Home operator
 - used only in the process of registration
- **to establish communication with a user IMPU (Public User Identity) is necessary.**
 - Every user has one or more IMPUs.
 - Each IMPI can have several IMPUs
 - Users can classify their public identities: business, family, friends, ...
 - E.g. sip:ruilaa@ua.pt, sip:steve.jobs@left.apple.com

