## Naming & DNS



CS 475, Fall 2019 **Concurrent & Distributed Systems** 

#### **Review: Recurring Problem: Replication**



OK, we obviously need to actually do something here to replicate the data... but what?

Replication solves some problems, but creates a huge new one: consistency









## **Review: CAP Theorem**

- Pick two of three:
  - Consistency: All nodes see the same data at the same time  $\bullet$ (sequential consistency)
  - Availability: Individual node failures do not prevent survivors from continuing to operate
  - Partition tolerance: The system continues to operate despite message loss (from network and/or node failure)
- You can not have all three, ever





## **Review: CAP Theorem**

- C+A: Provide strong consistency and availability, assuming there are no network partitions
- C+P: Provide strong consistency in the presence of network partitions; minority partition is unavailable
- A+P: Provide availability even in presence of partitions; no sequential consistency guarantee, maybe can guarantee something else



## **Review: Relaxing Consistency**

- We can relax two design principles:
  - How stale reads can be  $\bullet$
  - The ordering of writes across the replicas



#### **Review: Choosing a consistency model**

- Sequential consistency
  - All over it's the most intuitive
- Causal consistency
  - "Increasingly useful" but not really widely used still pay coordination cost, unclear what the performance benefits are
- Eventual consistency
  - Very popular in industry and academia
  - File synchronizers, Amazon's Bayou and more

GMU CS 475 Fall 2019



## **Example: Facebook**

- Problem: >1 billion active users
- Solutions: Thousands of servers across the world
- Need 100% availability!



## Example: Facebook

- Problem: >1 billion active users
- Solutions: Thousands of servers across the world
- What kind of consistency guarantees are reasonable? Need 100% availability!
- If I post a story on my news feed, is it OK if it doesn't immediately show up on yours?
  - Two users might not see the same data at the same time
  - Now this is "solved" anyway because there is no "sort by most recent first" option anyway





## **Example: Airline Reservations**

- GDS needs to sell as many seats as possible within given constraints

 Reservations and flight inventory are managed by a GDS (Global Distribution) System), who acts as a middle broker between airlines, ticket agencies and consumers [Except for Southwest and Air New Zealand and other oddballs]





## **Example: Airline Reservations**

- If I have 100 seats for sale on a flight, does it matter if reservations for flights are reconciled immediately?
- If I have 5 seats for sale on a flight, does it matter if reservations are reconciled immediately?
- Result: Reservations can be made using either a strong consistency model or a weak, eventual one
- Most reservations are made under the normal strong model (reservation is confirmed immediately)
- GDS also supports "Long Sell" issue a reservation without confirmed availability, need to eventually reconcile it
- Long sells require the seller to make clear to the customer that even though there's a confirmation number it's not confirmed!





#### HW3 Grades, as of Sun Nov 03 2019



Grade (Out of 100 points)

#### HW3 Graded

		]
		-
	1	
	1	
	1	
	1	
0 0 0		
0 0	0	
0	õ	
0	0	
0		
	0	

# HW4 - Push-based replication

- Each KVStore client will have the entire dataset cached locally
- When updating values, the update will be propagated to each replica





# HW4 - Push-based replication

- Each KVStore client will have the entire dataset cached locally
- When updating values, the update will be propagated to each replica  $\bullet$





# IVY VS HW4

- Ivy never copies the actual values until a replica reads them (unlike HW4) Invalidate messages are probably smaller than the actual data! Ivy only sends update (invalidate) messages to replicas who have a copy of
- the data (unlike HW4)
  - Maybe most data is not actively shared
- Ivy requires the lock server to keep track of a few more bits of information (which replica has which data)
- With near certainty Ivy is a lot faster :)



#### Today

- This week case studies in replication
- Today: DNS and naming (partially explaining how the internet works)
- Reminder:
  - HW4 is due 11/18!



## How do we find data?

- readable domain names
- DNS is a distributed system
- measure of consistency?



DNS - Domain Name System - responsible for mapping IP address to human-

Not immediately obvious how to scale: how do we maintain replication, some









#### Partitioning + Replication



- If input is structured, can possibly leverage that structure to build these buckets (name spaces)
- Example: File system
  - Map from: /home/bellj/teaching/swe622 to file contents  $\bullet$
  - Could have different machines responsible for each tier?
  - We will look at file systems on Wednesday
- Example: DNS system
  - Maps from: www.jonbell.net TO: 104.24.122.171
  - Different machines for each tier?

### Partitioning + Replication



- Obvious solution: Local file
  - Keep local copy of mapping from all hosts to all IPs (e.g., /etc/hosts)  $\bullet$
  - Hosts change IPs regularly: Download file frequently ullet
  - IPv4 space is now full  $\bullet$ 
    - 32-bits: 4,294,967,296 addresses
    - At 1 byte per address, file would be 4GB
    - Not a lot of disk space (now, DNS introduced in the late 80s)



- Obvious solution: Local file
  - Keep local copy of mapping from all hosts to all IPs (e.g., /etc/hosts)
  - Hosts change IPs regularly: Download file frequently
  - IPv4 space is now full
    - 32-bits: 4,294,967,296 addresses
    - At 1 byte per address, file would be 4GB
    - Not a lot of disk space (now, DNS introduced in the late 80s)



atible on a user-installable card. Ilation. The easy way to get hard 1059
rd Disk Drive
99 <sup>00</sup> Low As \$40 Per Month *
ernative Expansion Option ows Greater Data Storage
able Kit and installation required ontroller Board (25-1007). 699.00
0 SX/SL and original Tandy 1000 or up to 40 million characters of 20-megabyte hard disks.
PRICES APPLY AT PARTICIPATING R

#### We need 200x of these to hold 4GB: \$270K+

	Inflat	ion Calculator	
lf in	1989	(enter year)	
I purchased	an item for \$	699.99	
then in	2018	(enter year)	
that same <u>i</u>	tem would cost:	\$1,391.65	
Cumulative	rate of inflation:	98.8%	
		CALCULATE	





- Obvious solution: Local file
  - Keep local copy of mapping from all hosts to all IPs (e.g., /etc/hosts)  $\bullet$
  - Hosts change IPs regularly: Download file frequently  $\bullet$
  - IPv4 space is now full  $\bullet$ 
    - 32-bits: 4,294,967,296 addresses
    - At 1 byte per address, file would be 4GB
    - Not a lot of disk space (now, DNS introduced in the late 80s)
    - But a lot of constant internet bandwidth
  - More names than IPs lacksquare
    - Aliases
  - Not scalable!



- Obvious solution: Local file
  - Keep local copy of mapping from all hosts to all IPs (e.g., /etc/hosts)
  - Hosts change IPs regularly: Download file frequently
  - IPv4 space is now full  $\bullet$ 
    - 32-bits: 4,294,967,296 addresses
    - At 1 byte per address, file would be 4GB
    - Not a lot of disk space (now, DNS introduced in the late 80s)
    - But a lot of constant internet bandwidth
  - More names than IPs
    - Aliases
  - Not scalable!
- Obvious solution: Well known centralized server  $\bullet$ 
  - Single point of failure lacksquare
  - Traffic volume
  - Access time
  - Not scalable!

Query Volume (Millions/Day)



http://a.root-servers.org/static/index.html





- Goals
  - Scalable ullet
  - Robust
    - High availability
  - Decentralized maintenance lacksquare
  - Global scope lacksquare
    - Names mean the same thing everywhere
- Non-goals
  - Atomicity
  - Strong consistency



#### DNS

#### Idea: break apart responsibility for each part of a domain name (**zone**) to a different group of servers



Each zone is a continuous section of name space Each zone has an associate set of name servers



#### DNS

- Can have more/less servers replicating each zone based on popularity DNS responses are cached at clients
- - Caches periodically time out; bigger zones tend to have longer timeouts  $\bullet$ Quick response for the same request, also for similar requests





#### DNS: Example



















- 13 root servers
  - [a-m].root-servers.org
  - E.g., d.root-servers.org
- Handled by 12 entities
- How many physical servers?
  - a) Less than 13
  - b) 13
  - c) Tens
  - d) Hundreds
  - e) Thousands
  - f) Millions

Verisign, Inc.	а
Information Sciences Institute	b
Cogent Communications	С
University of Maryland	d
NASA Ames Research Center	е
Internet Systems Consortium, Inc.	f
U.S. DOD Network Information Center	g
U.S. Army Research Lab	h
Netnod	i
Verisign, Inc.	j
RIPE NCC	k
ICANN	l
WIDE Project	m



- 13 root servers
  - [a-m].root-servers.org
  - E.g., d.root-servers.org
- Handled by 12 entities
- How many physical servers?
  - a) Less than 13
  - <del>b) 13</del>
  - <del>c) Tens</del>
  - d) Hundreds
    - 980
  - e) Thousands
  - <del>f) Millions</del>

Verisian, Inc.		
Information Sciences Institute	b	
Cogent Communications	С	
University of Maryland	d	
NASA Ames Research Center	e	
Internet Systems Consortium, Inc.	f	
U.S. DOD Network Information Center	g	
U.S. Army Research Lab	h	
Netnod	i	
Verisign, Inc.	j	
RIPE NCC	k	
ICANN	1	
WIDE Project	m	







www.root-servers.net







#### www.root-servers.net







#### www.root-servers.net



Operator:	U.S. Army Research Lab
Locations:	Sites: 2 Aberdeen Proving Ground, US San Diego, US





www.root-servers.net









www.roota-serviers.net



### **Domain Name System - Scale**





### **Domain Name System - Scale**







## **Domain Name System - Zones**



#### **Domain Name System - Questions/Answers**

- DNS message format is the same for questions and answers  $\bullet$ 
  - Header
    - ID
      - Question/Answer have the same ID
    - Flags
      - 1 bit for question/answer, etc.
    - Number of questions
    - Number of answers
    - Number of authority RRs (Resource Records)
    - Number of additional RRs
  - Questions  $\bullet$
  - Answers
  - Authority
  - Additional Info





#### **Domain Name System - Resource Records (RRs)**

- RR format: (class, name, value, type, ttl)  $\bullet$ 
  - Class: Internet (IN)  $\bullet$
  - Туре  $\bullet$ 
    - A (AAAA for IPv6)  $\bullet$ 
      - name is hostname  $\bullet$
      - value is IP address
    - NS  $\bullet$ 
      - name is domain  $\bullet$
      - value is authoritative server for domain  $\bullet$
    - CNAME  $\bullet$ 
      - name is alias for some canonical (real) name  $\bullet$
      - value is canonical name  $\bullet$
    - And more...



# **Domain Name System - Example Query** dig(1) - Linux man page

#### Name

dig - DNS lookup utility

#### **Synopsis**

**dig** [@server] [-**b** address] [-**c** class] [-**f** filename] [-**k** filename] [-**m**] [-**p** port#] [-**q** name] [-**t** type] [-**x** addr] [-**y**[hmac:]name:key] [-**4**] [-**6**] [name] [type] [class] [queryopt...]

**dig** [-h]

**dig** [global-queryopt...] [query...]

GMU CS 475 Fall 2019



#### **Domain Name System - Example Query**

> dig www.gmu.edu a						
;; Got answer:						
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 20159						
;; flags: qr rd ra; QUERY: 1, ANSWER: 2, AUTHORITY: 3, ADDITIONAL: 3						
;; OPT PSEUDOSECTION:						
; EDNS: version: 0, flags:; udp: 4096						
;; QUESTION SECTION:						
;www.gmu.edu.		IN	А			
;; ANSWER SECTION:						
www.gmu.edu.	77455	IN	CNAME	jiju3.gmu.edu.		
jiju3.gmu.edu.	46534	IN	А	129.174.1.59		
;; AUTHORITY SECTION:						
gmu.edu.	86013	IN	NS	eve.gmu.edu.		
gmu.edu.	86013	IN	NS	uvaarpa.virginia.edu.		
gmu.edu.	86013	IN	NS	magda.gmu.edu.		
;; ADDITIONAL SECTION:						
eve.gmu.edu.	3219	IN	А	129.174.253.66		
magda.gmu.edu.	1993	IN	А	129.174.18.18		
uvaarpa.virginia.edu.	84640	IN	А	128.143.2.7		
;; Query time: 2 msec						
;; SERVER: 192.168.1.1#53(192.168.1.1)						
;; WHEN: Thu Feb 15 10:19:24 EST 2018						
;; MSG SIZE rcvd: 212						

GMU CS 475 Fall 2019



#### **Domain Name System - Example Query**

> dig www.ic.ac.uk a

;; Got answer:

- ;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 28622</pre>
- ;; flags: qr rd ra; QUERY: 1, ANSWER: 2, AUTHORITY: 0, ADDITIONAL: 1
- ;; OPT PSEUDOSECTION:
- ; EDNS: version: 0, flags:; udp: 4096
- ;; QUESTION SECTION:
- ;www.ic.ac.uk. IN Α

;; ANSWER SECTION:

- www.ic.ac.uk. wrp.cc.gslb.ic.ac.uk. 271 CNAME ΙN
- wrp.cc.gslb.ic.ac.uk. 146.179.40.24 1 IN Α
- ;; Query time: 1 msec
- ;; SERVER: 192.168.1.1#53(192.168.1.1)
- ;; WHEN: Thu Feb 15 10:23:52 EST 2018
- ;; MSG SIZE rcvd: 83

GMU CS 475 Fall 2019



#### **Domain Name System - Resolution**





GMU CS 475 Fall 2019



#### **Domain Name System - Resolution**



GMU CS 475 Fall 2019



#### **Domain Name System - (Recursive) Resolution**



GMU CS 475 Fall 2019



#### **Domain Name System - Iterative Resolution**



![](_page_47_Picture_5.jpeg)

![](_page_47_Picture_13.jpeg)

#### **Domain Name System - Iterative Resolution**

![](_page_48_Figure_1.jpeg)

![](_page_48_Picture_5.jpeg)

#### **Domain Name System - Example Query**

> cat /etc/resolv.conf #which DNS server am I using

# Generated by resolvconf search fios-router.home nameserver 192.168.1.1

> dig @192.168.1.1 www.ic.ac.uk a #recursive query

- • •
- > dig +norecurse @192.168.1.1 www.ic.ac.uk a
- # same answer, why?
- # what if I had tried this first?

![](_page_49_Picture_10.jpeg)

https://xkcd.com/908/

![](_page_49_Picture_13.jpeg)

# **Domain Name System - Caching**

Some gamers steamed over alleged Valve anti-cheat DNS spying - CSO https://www.csoonline.com/.../some-gamers-steamed-over-alleged-valve-anti-cheat-dn... 💌 Feb 16, 2014 - Goes through all your DNS Cache entries (ipconfig /displaydns); Hashes each one with MD5; Reports back to VAC Servers. Valve is not the only company that uses an anti-cheat system, but it is perhaps one of the most highly regarded companies as countless millions of gamers have Steam. Various ...

Reddit user claims Valve Anti-Cheat scans your DNS cache - PC ... https://gamefags.gamespot.com/boards/916373-pc/68593356 -For PC on the PC, a GameFAQs message board topic titled "Reddit user claims Valve Anti-Cheat scans your DNS cache".

https://www.neogaf.com > Discussions > Gaming Discussion -

Feb 16, 2014 - Trust is a critical part of a multiplayer game community - trust in the developer, trust in the system, and trust in the other players. Cheats are a negative sum game, where a minority benefits less than the majority is harmed. There are a bunch of different ways to attack a trust-based system including writing a ...

Report: Valve anti-cheat scans your DNS history - Player Attack

Feb 17, 2014 - Even if you've never actively visited a cheat website, there may be traces of them in your DNS, and that's what VAC is reportedly now looking for. The news was first posted to the Counter-Strike: Global Offensive Reddit, explaining that VAC now: Goes through all your DNS Cache entries (ipconfig ...

#### IS IT OK THAT VAC SCANS YOUR DNS CACHE? :: VAC Discussion - Steam... steamcommunity.com > Steam Forums > VAC Discussion 👻

Feb 15, 2014 - 16 posts - 7 authors Valve ANSWER THIS! http://www.ghacks.net/2014/02/16/steams-vac-protection-now-scans-anstransfers-dns-cache/ What is going on with this DNS spying? We all know that various anti-cheat programs do check your DNS, some of them don't really collect data though. It has been proven that VAC collects ...

#### Valve Anti-Cheat seems to scan your DNS cache, but probably doesn ...

#### https://www.playerattack.com/news/.../report-valve-anti-cheat-scans-your-dns-history/ -

![](_page_50_Picture_17.jpeg)

# **Domain Name System - Caching**

There are a number of kernel-level paid cheats that relate to this Reddit thread. Cheat developers have a problem in getting cheaters to actually pay them for all the obvious reasons, so they start creating DRM and anti-cheat code for their cheats. These cheats phone home to a DRM server that confirms that a cheater has actually paid to use the cheat.

VAC checked for the presence of these cheats. If they were detected VAC then checked to see which cheat DRM server was being contacted. This second check was done by looking for a partial match to those (non-web) cheat DRM servers in the DNS cache. If found, then hashes of the matching DNS entries were sent to the VAC servers. The match was double checked on our servers and then that client was marked for a future ban. Less than a tenth of one percent of clients triggered the second check. 570 cheaters are being banned as a result.

Gabe Newell, Valve's CEO

https://www.reddit.com/r/gaming/comments/1y70ej/valve\_vac\_and\_trust/

![](_page_51_Picture_7.jpeg)

![](_page_51_Figure_9.jpeg)

![](_page_51_Picture_10.jpeg)

#### Conclusion

#### Resolving names requires a large scale distributed system $\bullet$

- Domain Name System (DNS) lacksquare
- Distributed between several entities and among all continents ullet
- World-wide scale ullet
- High availability

![](_page_52_Picture_10.jpeg)

#### This work is licensed under a Creative Commons Attribution-ShareAlike license

- This work is licensed under the Creative Commons Attribution-ShareAlike 4.0 International
- You are free to:
  - Share copy and redistribute the material in any medium or format
  - Adapt remix, transform, and build upon the material
  - for any purpose, even commercially.
- Under the following terms:
  - suggests the licensor endorses you or your use.
  - contributions under the same license as the original.
  - legally restrict others from doing anything the license permits.

License. To view a copy of this license, visit <u>http://creativecommons.org/licenses/by-sa/4.0/</u>

• Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that

• ShareAlike — If you remix, transform, or build upon the material, you must distribute your

No additional restrictions — You may not apply legal terms or technological measures that

![](_page_53_Picture_18.jpeg)