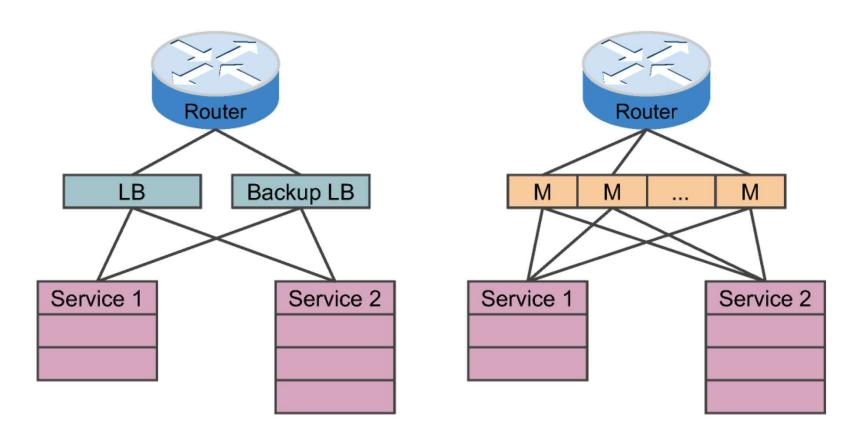


# HY-559 Infrastructure Technologies for Large-Scale Service-Oriented Systems

Kostas Magoutis magoutis@csd.uoc.gr http://www.csd.uoc.gr/~hy559

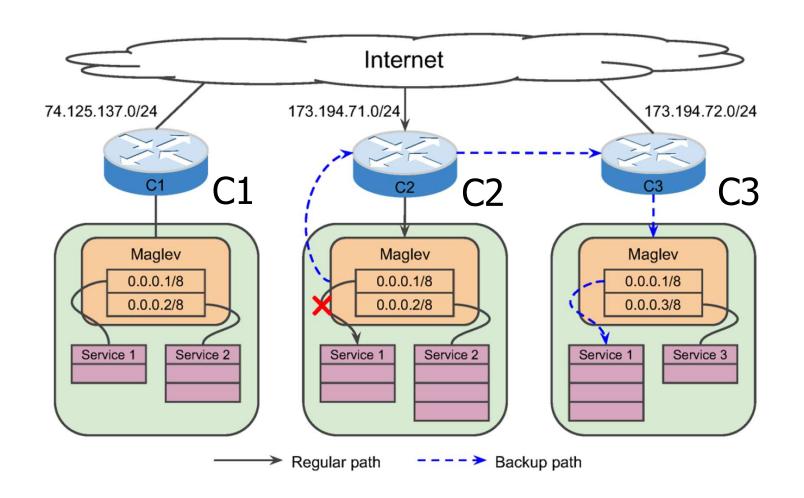
## Network Load Balancing: Google Maglev

#### Hardware load balancer and Maglev



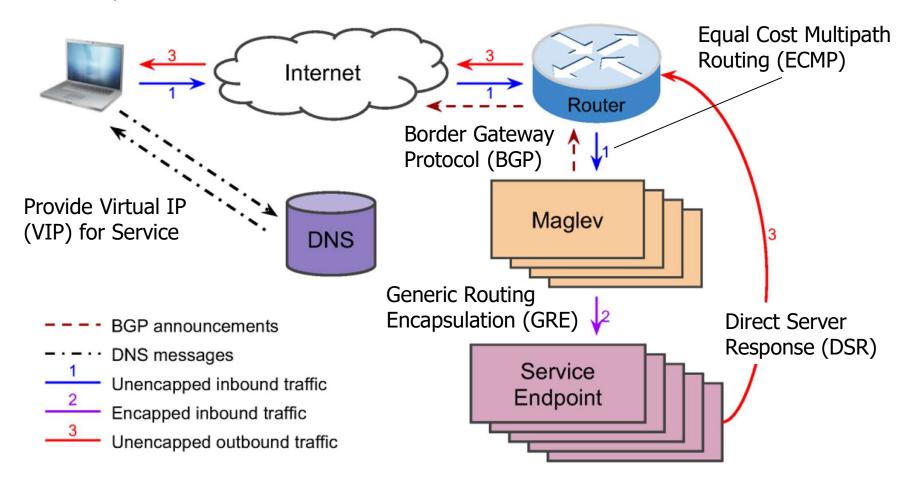
# Virtual IP (VIP) assignment

- The same service is configured as different VIPs in different clusters
- The user is directed to one of them by DNS
- Service1 is configured as 74.125.137.1 in C1 and 173.194.71.1 in C2



#### Maglev packet flow

BGP advertises prefix of each cluster e.g., 74.125.137.0/24

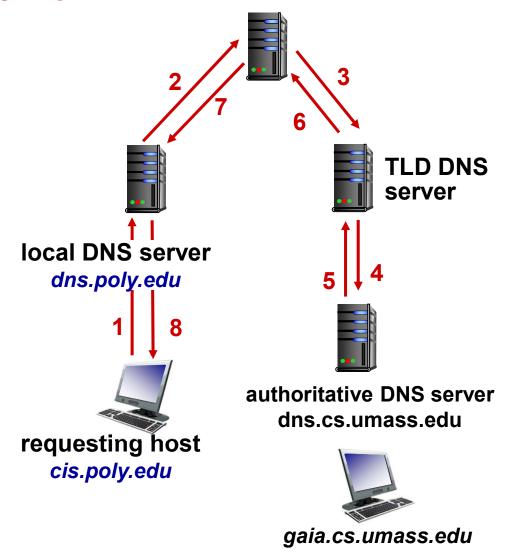


#### DNS recursive resolver

#### root DNS server

#### EDNS0 extension

- EDNS Client Subnet
- Allows a recursive resolver to include the client's IP address or subnet in the query
- Provide more relevant and geographically closer responses, improving performance and lowering latency for end user



# DNS: caching and updating records

- Once any name server learns mapping, it caches it
  - Cache entries timeout after some time (TTL)
  - TLD servers cached in local name servers
    - Thus root name servers are not visited often
- update/notify mechanisms under design by IETF
  - RFC 2136
  - http://www.ietf.org/html.charters/dnsind-charter.html

#### **DNS** records

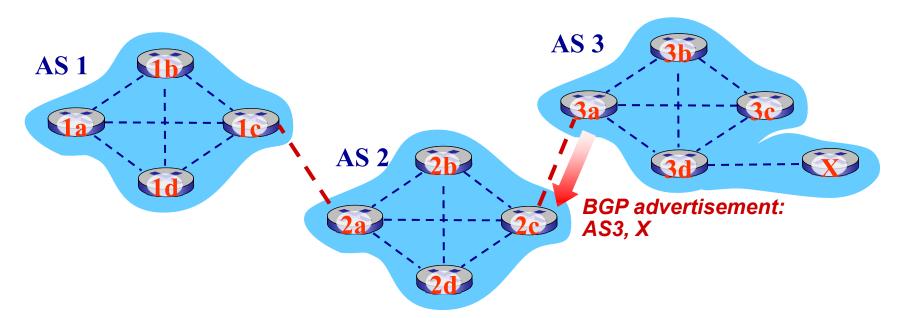
RR format: (name, value, type, TTL)

- ☐ Type=A
  - name is hostname
  - value is IP address
- Type=NS
  - name is domain (e.g. foo.com)
  - value is hostname of authoritative name server for this domain

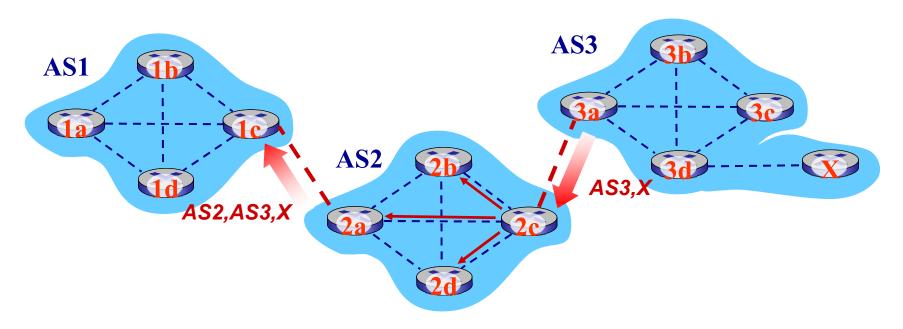
- □ Type=CNAME
  - name is alias for some "canonical" (real) name www.ibm.com is really servereast.backup2.ibm.com
  - value is canonical name
- ☐ Type=MX
  - value is name of mail server associated with name

#### Border Gateway Protocol (BGP) basics

- BGP session: two BGP routers ("peers") exchange BGP messages over semi-permanent TCP connection:
  - advertising paths to different destination network prefixes (BGP is a "path vector" protocol)
- when AS3 gateway router 3a advertises path AS3,X to AS2 gateway router 2c:
  - AS3 promises to AS2 it will forward datagrams towards X



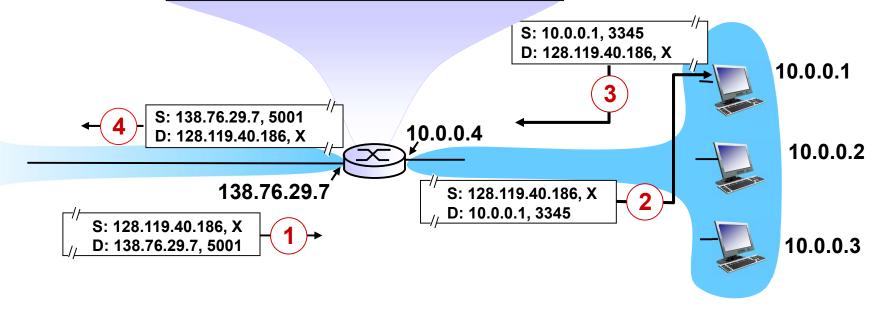
#### BGP path advertisement



- AS2 router 2c receives path advertisement AS3,X (via eBGP) from AS3 router 3a
- Based on AS2 policy, AS2 router 2c accepts path AS3,X, propagates (via iBGP) to all AS2 routers
- Based on AS2 policy, AS2 router 2a advertises (via eBGP) path AS2, AS3, X to AS1 router 1c

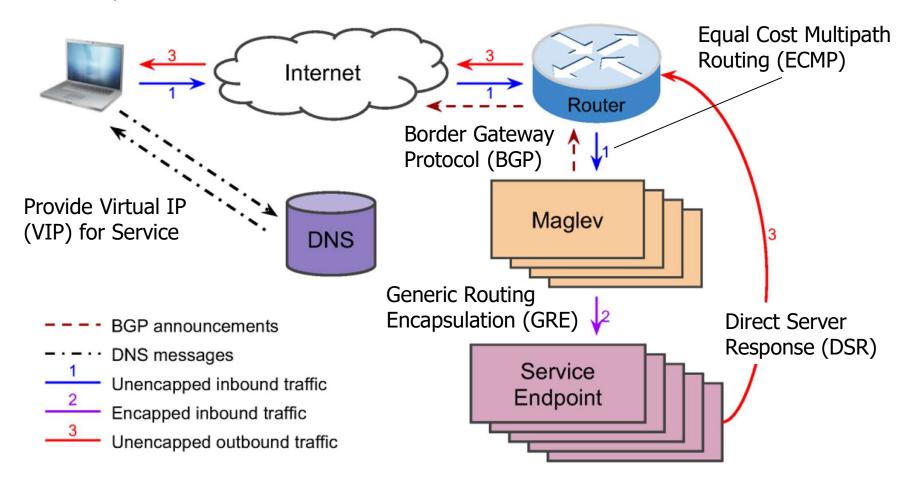
#### NAT / Reverse NAT

NAT translation table			
WAN side addr	LAN side addr		
138.76.29.7, 5001	10.0.0.1, 3345		

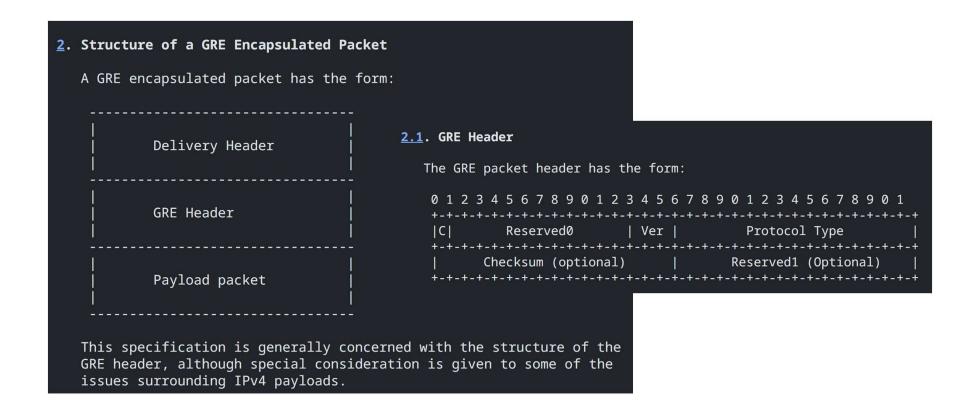


#### Maglev packet flow

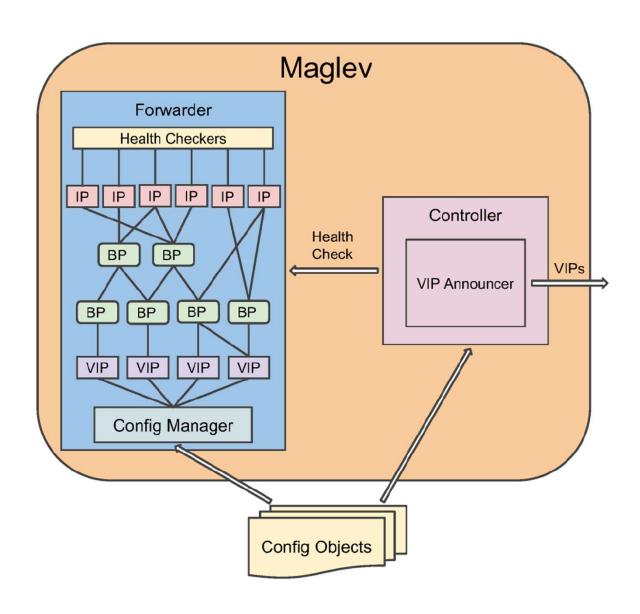
BGP advertises prefix of each cluster e.g., 74.125.137.0/24



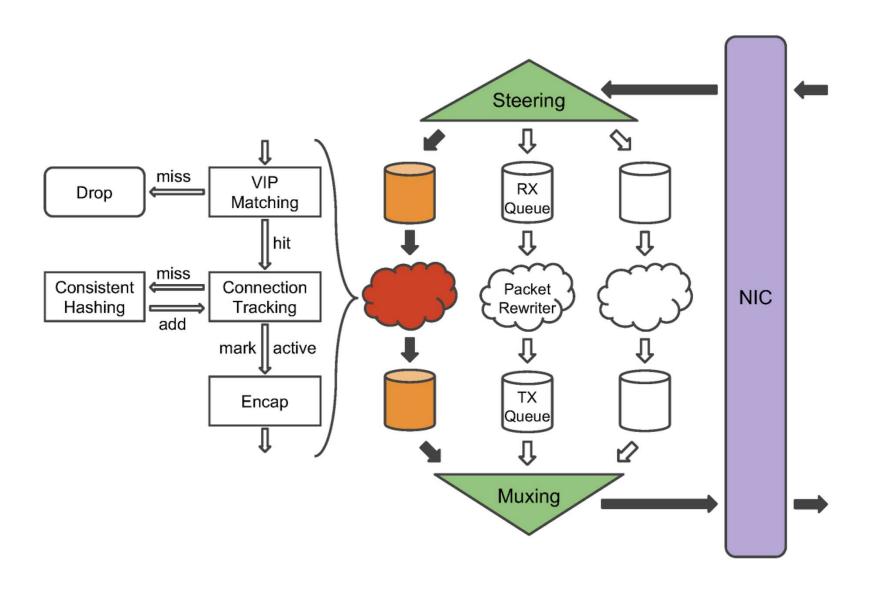
## Generic Routing Encapsulation (GRE)



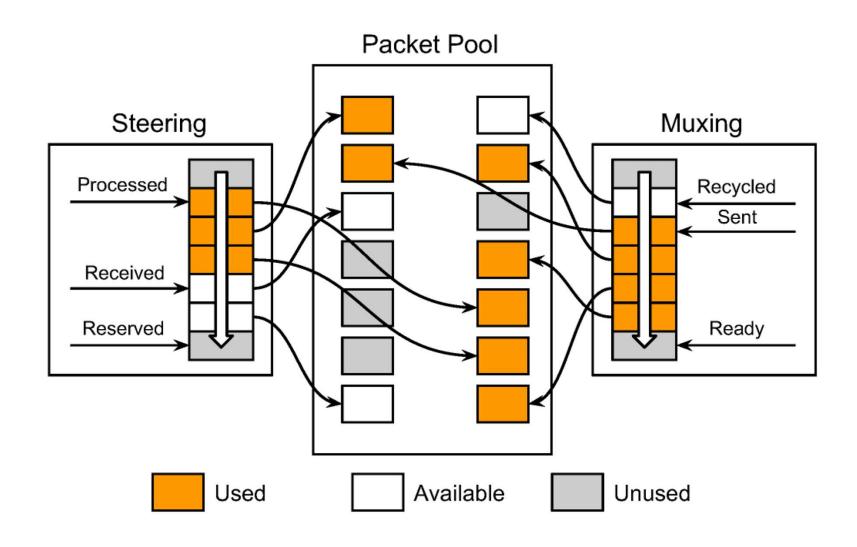
# Maglev config



# Maglev forwarder structure



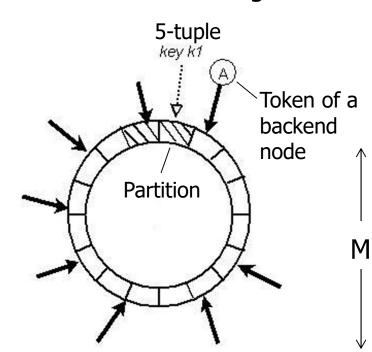
#### Packet movement in and out of forwarder



## Maglev consistent hashing

#### Consistent hashing

 $offset \leftarrow h_1(name[i]) \mod M$  $skip \leftarrow h_2(name[i]) \mod (M-1)+1$  $permutation[i][j] \leftarrow (offset + j \times skip) \mod M$ 



Per VIP: (3,4), (0,2) and (3,1)

Table 1: A sample consistent hash lookup table.

		B0	B1	B2
	0	3	0	3
	1	0	2	4
	2	4	4	5
	3	1	6	6
	4	5	1	0
	5	2	3	1
·	6	6	5	2

	Before	After
0	<i>B</i> 1	<i>B</i> 0
1	<i>B</i> 0	B0
2	<i>B</i> 1	<i>B</i> 0
3	<i>B</i> 0	<i>B</i> 0
4	<i>B</i> 2	<i>B</i> 2
5	<i>B</i> 2	<i>B</i> 2
6	<i>B</i> 0	<i>B</i> 2

N backend nodes T random tokens per backend node M equal-sized partitions M >> N\*T

backends.

Permutation tables for the Lookup table before and after B1 is removed.

#### Certain VIPs may have hundreds of backends

M > 100\*N to ensure at most a 1% difference in hash space assigned to backends

#### Maglev consistent hashing

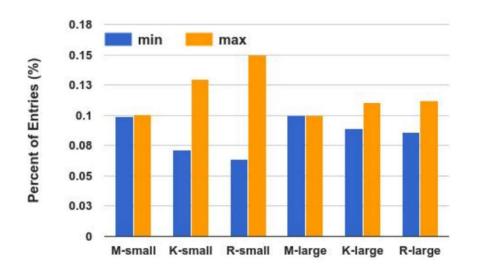


Figure 11: Load balancing efficiency of different hashing methods. M, K and R stand for Maglev, Karger and Rendezvous, respectively. Lookup table size is 65537 for *small* and 655373 for *large*.

- Total number of backends: 1000
- Lookup table size is 65537 (Small) and 655373 (Large)
- Max, min % of table entries per backend for each table size

# Summary: Traditional datacenter load balancing architecture

