Towards TCAM-based scalable virtual routers

Layong Luo, Gaogang Xie, Yingke Xie (ICT/CAS) Kavé Salamatian (U. of Savoie) Steve Uhlig (QMUL) Laurent Mathy (U. of Liège/Lancaster U)







Motivation

Virtual routers (VRs)

- key building blocks for enabling network virtualization
- VPN, network testbeds, Future Internet ...

Memory issue

- The number of FIBs, and the size of each FIB, are expected to increase continuously
- FIBs are preferably stored in high-speed memory (SRAMs or TCAMs) but limited amount of it
- Scalability challenge: support as many FIBs as possible in the limited high-speed memory?

Related work

- SRAM-based scalable virtual routers (How to merge multiple tries to achieve good scalability?)
 Trie overlap, CoNEXT 2008
 Trie braiding, INFOCOM 2010
 - Multiroot, ICC 2011
 -
- None of previous work has exploited the possibility of using TCAMs to build scalable virtual routers

Background





Non-shared approach for TCAM-based virtual



Merged data structure



Lookup issue Example 1: IP 100, VID 0 Example 2: IP 000, VID 1 FIB 0 FIB 1 next hop prefix prefix next hop next hop correct NH! prefix A4 **B4** 100* 0* A1 **B1 B**5 0 101* A5 1* A2 **B2** 1* **B6** 111* A6 00* **B3 Incorrect NH!** A3 11* 00* A3 0 100 100* A4 **B4** 11* 0 **B**3 101* **B5** A1 101[°] A5 0* **B1** ??? 111* 111* **B6** A6 A2 **B2** (a) (b) TCAM SRAM

TCAM FIB merging principle

 Incorrect matching, resulting from the masking of a shorter prefix (e.g., <0*, B1>) in an original FIB by a longer prefix (e.g., <00*, 0>) in the merged FIB, must be avoided.

Two approaches

Two FIB merging approaches that respect the principle

- FIB Completion
- FIB Splitting

FIB completion

Basic idea

 Whenever a prefix from the merged FIB doesn't appear in a given individual FIB, we simply associate it with a valid NH in this FIB according to the LPM rule.



Fig. 1. (a) The basic merged FIB, and (b) its completed version

Completion process

Auxiliary tries help the completion process



Fig. 1. two tries built from the two sample FIBs

Fig. 2. (a) a merged trie using trie overlap^[1], and (b) its completed version

[1] J. Fu and J. Rexford, Efficient IP-address lookup with a shared forwarding table for multiple virtual routers, CoNEXT 2008

Update scenarios



FIB splitting

- The main drawback of FIB completion
 - High update overhead in the worst case due to the masking prefix correction process
- Another way to address the lookup issue is
 - To ensure that when a masking prefix is a hit for a lookup in the merged FIB, the corresponding masked prefix is also made available to the lookup process, rather than correct the NHs of masking prefixes, thus reducing the worst-case update overhead

Basic idea

- The naturally disjoint leaf prefixes, which are about 90% of the total prefixes, are merged in one TCAM
- The remaining small overlapping prefix set is stored in another TCAM using the non-shared approach





Update scenarios



Performance evaluation

- Routing tables
 - 14 full routing tables from core routers

Evaluation

- TCAM size
- SRAM size
- Total cost of the system
- Lookup and update performance

TCAM size



SRAM size



Total cost of the system

Table 2. Reference prices of TCAMs and SRAMs

Memory	Part No.	Capacity	Speed	Price
TCAM	NL9512	512K×40bit	250MHz	\$387.2
SRAM	CY7C1525	8M×9bit	250MHz	\$89.7

Table 3. Cost of the three approaches for IPv4 FIBs

	# of	# of	Total
	TCAMs	SRAMs	Cost
Non-shared	11	1	\$4348.9
FIB completion	1	1	\$476.9
FIB splitting	3	2	\$1341

Lookup & update performance

Table 5. Worst-case lookup and update overhead

	Lookup	Update
Non-shared	O(1)	W/2
FIB completion	O(1)	$2^{W+1}-1$
FIB splitting	O(1)	NW/2

W: the length of the IP address

N: the number of virtual routers

Conclusions

- Main contributions
 - The first work to exploit the possibility of using TCAMs to build scalable virtual routers
 - Merged data structure and merging principle for TCAMs
 - Two different approaches
 - FIB completion: best scalability but high worst-case update overhead
 - FIB splitting: good scalability with a more reasonable upper bound on the worst-case update overhead

Future work

 Quantify the actual update overhead of our approaches based on realistic update workloads

Thank you!