ADVANCED DATA STRUCTURES IN TELECOMMUNICATION

János Tapolcai, Rétvári Gábor, Körösi Attila

Budapest University of Technology and Economics (BME) Faculty of Electrical Engineering and Informatics (VIK) Department of Telecommunications and Media Informatics (TMIT) High-Speed Networks Laboratory (HSNLab) MTA-BME Future Internet Research Group



BME Modellalkotás szeminárium16/9/2014

MTA-BME Future Internet Research Group

Gábor Rétvári, András Gulyás, József Bíró, Zalán Heszberger, Péter Babarczi, Balázs Sonkoly, Felicián Németh, Attila Körösi, Biczók Gergely, Toka László, Csikor Levente

+ 8 PhD and >30 MSc students



IEEE Infocom papers 2009-2013



Weighted by the population (1999-2011)

Countries	#Citizents	#Papers	#Pa	pers/million people	e o	.0 5	.0	10.0	15.0
Israel	7521400	C	82	10,	9				
Hong Kong	702640	О.	53	7,	5				
Switzerland	778290	О.	56	7,	2				
Cyprus	80185	1	3	3,	7				
Singapore	498760	C	16	3,	2				
Iceland	31763	C	1	3,	1				
Hungary	1001362	8	21	2,	1				
Taiwan	23119772	2	41	1,	8				
Finland	535680	C	9	1,	7				
United States	30896400) 4	96	1,	6				
France	65447374	4	88	1,	3				
Greece	11306183	3	15	1,	3				
United Kingdom	n 6204170	8	77	1,	2				
Italy	6027696	9	72	1,	2				
Sweden	934513	5	11	1	2				

Applied computer science

Scientific value



Main Research Topics

- Network design
 - Backbone optical, Metro IP, OpenFlow
- Reliability
 - "Your computer can go wrong, but the internet never."
- Routing and addressing
 - Software Defined Networking
- Packet processing
 - Router design, Firewall rules

Packet processing

 G. Rétvári, J. Tapolcai, A. Kőrösi, A. Majdán Heszberger, "Compressing IP Forwarding Tables: Towards Entropy Bounds and Beyond", In ACM SigCOMM, 2013.

Accepted to IEEE/ACM Transactions on Networking

 A. Kőrösi, J. Tapolcai, B. Mihálka, G. Mészáros, G. Rétvári, "Compressing IP Forwarding Tables: Realizing Information-theoretical Space Bounds and Fast Lookups Simultaneously", In Proc. IEEE ICNP, 2014.

A Core Internet Router

- Holds info on the whereabouts of every single IP host
- That ought to be a huge amount of information



Cisco CRS-3 line card up to 8 Gbyte memory 533 MHz DDR2 >300 Watt http://www.cisco.com/en/US/docs/routers/

Internet Routing

Packets are marked by 32 bit IP addresses weber.tmit.bme.hu=152.66.130.2 =10011000010000101111010001101111

Hierarchy of subnetworks by increasing prefix length



IP Forwarding Information Base

- The key data structure in Internet routing: tells a router where to forward a packet
- FIB: a database of prefix-to-next-hop associations
 lookup a 32 bit long key: longest prefix match

update the association for some prefix

Address/prefix length	Label
-/0	2
0/1	3
00/2	3
001/3	2
01/2	2
011/3	1



IP Forwarding Information Base

- Stores more than 440K IP-prefix-to-nexthop mappings as of January, 2013
 - consulted on a packet-by-packet basis at line speed
 - queries are complex: longest prefix match
 - updated couple of hundred times per second
 - takes several MBytes of fast line card memory and counting
- May or may not become an Internet scalability barrier

http://lendulet.tmit.bme.hu/fib_comp/

Measurements



rtr.bme.hbone.hu

FIB Representations

- Prefix tree: search tree over the address space
- Essentially a labeled ordinal binary tree (binary trie)
- \Box Lookup/ update are O(W) for W bit address size

Address/prefix length	Label
-/0	2
0/1	3
00/2	3
001/3	2
01/2	2
011/3	1

Compressed data structures

- Compression not necessarily sacrifices fast access!
- Store information in entropy-bounded space and provide fast in-place access to it
 - take advantage of regularity, if any, to compress data drifts closer to the CPU in the cache hierarchy operations are even faster than on the original uncompressed form
- No space-time trade-off!
- Goal: advocate compressed data structures to the networking community
 - IP forwarding table compression as a use case

prefix	neighbour	
010/3	1	
11/2	2	$\begin{array}{c} \begin{array}{c} \begin{array}{c} \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} $
0/1	3	
1111/4	1	
100/3	3	
-/0	1	

 \sim

prefix	neighbour	
010/3	1	X X \overline{X} \overline{X} X X \overline{X} \overline{X}
11/2	2	$\begin{array}{c} \dot{\Lambda} & $
0/1	3	
1111/4	1	
100/3	3	
-/0	1	

prefix 010/3 11/2 0/1 1111/4 100/3 -/0	neighbour 1 2 3 1 3 1	
-/0	1	

prefix	neighbour	
010/3	1	
11/2	2	$ \begin{array}{c} $
0/1	3	
1111/4	1	
100/3	3	
-/0	1	



prefix	neighbour	
010/3	1	A A A A A A A A
11/2	2	
0/1	3	
1111/4	1	
100/3	3	
-/0	1	





prefix	neighbour
010/3	1
11/2	2
0/1	3
1111/4	1
100/3	3
-/0	1

					2
profix	naimhhaur	OM COM	~~~	8	
prenx	neignbour	_	*	* *	\rightarrow
010/3	1				Ŭ Ă
11/2	2				Ó 🌢
0/1	3				
1111/4	1				
100/3	3				
-/0	1				

Example – Leaf Pushed

		•	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		2
prenx	neignbour		🖌 🍾	🏅 🍾	*
010/3	1				
11/2	2				é è
0/1	3	prefix	neighbour		
1111/4	1	010/2	1		
100/3	3	010/5	1		
100/0	1	11/2	2		
-/0		1111/4	1		
		101/3	1		
		-/0	3		

FIB Aggregation

- remove redundancy from the binary trie without changing forwarding behavior
 - Label-optimized trie
 - Leaf-pushed trie: proper leaf-labeled trie
 - unique (normalized) form
 - Level-compressed trie: remove excess levels



Optimal Routing Table Constructor (ORTC)

- Minimum label
 - Dynamic programming, INFOCOM'99

Assign the set of labels for each branch



Merge children





Merging finished



□ From root of the tree

Skip if child next hop is the same as its parent



- □ From root of the tree
 - Skip if child next hop is the same as its parent



- □ From root of the tree
 - Skip if child next hop is the same as its parent



- □ From root of the tree
 - Skip if child next hop is the same as its parent



Optimal solution



Dynamic FIBs: Trie-folding

Practical FIB compression, a good old pointer machine
 Fold the trie into a prefix DAG (DAFSA, DAWG, BDD)



- For good compression, we need the tree to be in a prefix-free form
 - But prefix-free forms are expensive to update
 - \square Balance by a parameter λ , called the leaf-push barrier
- Deployable with minimal modification to router ASICs
 ORTC (Min Label) for DAG is NP hard

Prefix DAG Size (Pointers)

View the problem as string compression: encode a string S into a prefix DAG D(S)



Entropy bound

For Shannon Entropy - We need a random model

- String compression
- Next hops are generated with a random process
- $\square p_s$ is the probability of a next hop



Coupon collector's problem with arbitrary coupon probabilities

\Box Set of coupon C, we draw *m* coupons

\Box At each draw coupon o has a probability p_o

Let V denote the set of coupons we have after m draw. The probability of having coupon o in V is

$$P(o \in V) = 1 - (1 - p_o)^m \tag{1}$$

Thus the expected cardinality of V is

$$E(|V|) = \sum_{o \in C} E(I(o \in V)) = \sum_{o \in C} P(o \in V) =$$
$$= \sum_{o \in C} (1 - (1 - p_o)^m) \quad (2)$$

Let H_C denote the entropy of the coupon distribution

$$H_C = \sum_{o \in C} p_o \log_2 \frac{1}{p_o} \tag{3}$$

Coupon collector Lemma

Lemma 1.

n number of coupons

$$E(V) \le \min\left\{\frac{m}{\log_2(m)} \cdot H_C + 3, m, n\right\}$$

for $m \geq 3$.

Proof: Trivially holds that $E \leq m$ and $E \leq n$. Next, let us expand

$$\sum_{o \in C} (1 - (1 - p_o)^m) \le \frac{m}{\log_2(m)} \sum_{o \in C} p_o \log_2 \frac{1}{p_o} + 3 \quad (4)$$

$$3 > \frac{1}{e}$$

The above inequality holds if the inequality holds for each $o \in C$. Thus next we prove that

$$1 - (1 - p_o)^m \le \frac{m}{\log_2(m)} p_o \log_2 \frac{1}{p_o} \quad \text{if } p_o < \frac{1}{e} \tag{5}$$

Case 1

$$1 - (1 - p_o)^m \le \frac{m}{\log_2(m)} p_o \log_2 \frac{1}{p_o} \quad \text{if } p_o < \frac{1}{e} \quad (5)$$

Let us assume $m \ge \frac{1}{p}$. Note that the right hand size is a monotone increasing function of m, when m > e. Thus we can substitute $m = \frac{1}{p}$ if $\frac{1}{p} > e$ in the right hand side and we get

$$1 - (1 - p_o)^m \le \frac{1/p_o}{\log_2(1/p_o)} p_o \log_2 \frac{1}{p_o} = 1.$$
(6)

Case 2

In the rest of the proof we focus on the other case, which is $m < \frac{1}{p_0}$. Let us define 1 > x > 0 as

$$x := \log_{\frac{1}{p_o}} m \quad . \tag{7}$$

After substituting $m = \frac{1}{p_o^x}$ we have

$$1 - (1 - p_o)^{\frac{1}{p_o^x}} \le \frac{\frac{1}{p_o^x}}{\log_2\left(\frac{1}{p_o^x}\right)} p_o \log_2\left(\frac{1}{p_o}\right) = \\ = \frac{\frac{1}{p_o^x}}{x \log_2\left(\frac{1}{p_o}\right)} p_o \log_2\left(\frac{1}{p_o}\right) = \frac{\frac{1}{p_o^x}}{x} p_o = \frac{1}{x p_o^{x-1}} , \quad (8)$$

which can be reordered as

$$(1-p_o)^{\frac{1}{p_o^x}} \ge 1 - \frac{1}{xp_o^{x-1}}$$
 (9)

Case 2 Cont'

$$(1-p_o)^{\frac{1}{p_o^x}} \ge 1 - \frac{1}{xp_o^{x-1}}$$
.

(9)

Taking the $p_o^{x-1} > 0$ exponent of both sides we get

$$(1-p_o)^{\frac{1}{p_o}} \ge \left(1-\frac{1}{xp_o^{x-1}}\right)^{p_o^{x-1}}.$$
 (10)

Note that x < 1, thus $\frac{1}{x} > 1$, and we can prove

$$(1-p_o)^{\frac{1}{p_o}} \ge \left(1-\frac{1}{p_o^{x-1}}\right)^{p_o^{x-1}}.$$
 (11)

Bernoulli discovered that $(1 - p_o)^{\frac{1}{p_o}}$ is monotone decreasing function and equals to $\frac{1}{e}$ for $p_o \rightarrow 0$ []. Thus the inequality holds if

$$p_o \le \frac{1}{p_o^{x-1}}.\tag{12}$$

which holds because

$$p_o^x \le 1. \tag{13}$$

Number of nodes in the DAG

Coupon collector problem

h height
$$E(|V_D^j|) \le \frac{2^{h-j}}{\log_2(2^{h-j})} H_C + 3 = \frac{2^{h-j}}{\log_2(2^{h-j})} H_O 2^j + 3$$

 $= \frac{2^h}{h-j} H_O + 3,$
where $H_C = H_O 2^j$ is the entropy of a 2^j long string
 V_D^j
 $E(|V_D^j|) \le \min\left\{\frac{H_O}{h-j}2^h + 3, 2^{h-j}, \delta^{2^j}\right\}$

 δ number of next hops



Evaluation

Entropy bound (E) is way smaller than information theoretic limit (I)

100-400

KBytes

	FIB	N	δ	H_0	Τ	\tilde{E}	YRW-b	DAG
	1112	410,510		1.00			ADW-0	170
SS	taz	410,513	4	1.00	94	56	63	1/8
	hbone	410,454	195	2.00	356	142	149	396
See	access(d)	444,513	28	1.06	206	90	100	370
a	access(v)	2,986	3	1.22	2.8	2.2	2.5	7.5
	mobile	21,783	16	1.08	0.8	0.4	1.1	3.6
	as1221	440,060	3	1.54	130	115	111	331
e le	as4637	219,581	3	1.12	52	41	44	129
S	as6447	445,016	36	3.91	375	277	277	748
	as6730	437,378	186	2.98	421	209	213	545
'n.	fib_600k	600,000	5	1.06	257	157	179	462
S	fib_1m	1,000,000	5	1.06	427	261	297	782

- Several million lookups per sec both in HW and SW
 - faster than the uncompressed form
- pDAG tolerates more than 100, 000 updates per sec

Level Compressed prefix DAG

Min Size DAG NP hard

Lower bound by a Linear Program relaxation



Evaluation

		-					fib_	trie	lcTri	e	bDA	G	lcDAG	
		N	σ	H_0	Ι	E	M	ν	M	ν	M	ν	M	ν
ш]	bme	499,211	89	1.20	196	71	32198	451.05	635.95	8.91	203.72	2.85	162.05	2.27
Z	szeged	499,236	87	1.20	196	71	32198	450.75	636.06	8.90	203.76	2.85	162.09	2.27
Ĕ I	vh1	499,143	207	2.34	407	185	32197	173.60	1157.48	6.24	503.22	2.71	393.34 2.12	
ш т	vh2	499,302	101	1.20	196	72	32202	450.31	636.49	8.90	204.13	2.85	162.36	2.27
et2 2.	atlanta	14,312	93	1.91	38	14	1017	71.82	110.75 7.82		48.65	3.44	41.11	2.90
n I	houston	14,305	101	1.12	38	14	1016	71.93	109.47	7.75	49.00	3.47	41.49	2.94
Inte	kansas	14,335	101	1.06	38	14	1019	73.31	110.43	7.95	47.38	3.41	39.99	2.88
1	taz	410,513	4	0.97	94	56	26698	474.62	519.23	9.23	172.95	3.07	138.75	2.47
Sess .	access(d)	444,513	28	1.61	206	95	28713	302.07	869.80	9.15	252.88	2.66	201.23	2.12
acc	access(v)	2,986	3	0.99	3	2	192	88.99	15.23	7.06	8.26	3.83	7.16	3.32
I	mobile	21,783	10	1.62	1	0	290	655.67	2.47	5.59	1.27	2.88	1.19	2.69

		fib_trie					lcTrie					bDAG					lcDAG				
		\overline{h}	$h_{ m max}$	mpps	#CPU	cm	\overline{h}	h_{\max}	mpps	#CPU	cm	\overline{h}	h_{\max}	mpps	#CPU	cm	\overline{h}	h_{\max}	mpps	#CPU	cm
HBONE	bme	2.43	8	4.07	758	23900	6.62	13	10.39	297	2145	25.85	31	10.69	289	1926	9.61	15	10.92	283	1903
	szeged	2.43	8	4.09	756	23583	6.62	13	10.40	297	2135	25.57	31	10.69	289	1960	9.49	15	10.92	283	1902
	vh1	2.43	8	4.09	754	23407	5.56	13	10.11	306	2533	25.52	31	9.92	311	2359	7.56	15	10.24	301	2247
	vh2	2.43	8	4.05	763	24463	6.62	13	10.40	297	2137	25.86	31	10.69	289	1965	9.55	15	10.88	284	1902
Internet2	atlanta	3.37	10	5.87	526	2442	6.31	14	9.94	311	1803	29.45	31	10.65	290	1842	11.04	17	11.15	277	1834
	houston	3.37	10	5.86	527	2476	6.15	14	9.99	309	1837	28.72	31	10.64	290	1839	11.12	16	11.11	278	1839
	kansas	3.37	10	5.87	526	2496	6.04	15	10.21	303	1796	29.44	31	10.65	290	1836	10.89	16	11.25	274	1823
access	taz	2.42	6	4.45	694	20735	6.53	13	10.57	292	2081	21.00	31	10.73	288	1882	7.75	16	10.96	282	1874
	access(d)	2.44	8	4.25	726	21941	6.06	14	9.25	334	2417	28.85	31	10.81	286	2105	10.96	16	11.00	281	1982
	access(v)	5.53	9	12.45	248	1798	10.31	15	12.94	239	1817	20.67	31	13.40	230	1828	10.67	18	14.83	208	1846
	mobile	6.68	9	12.33	250	1813	9.70	14	14.13	218	1776	26.44	31	17.67	174	1828	11.00	16	22.56	137	1832



- Applied mathematics
 - Combinatorial optimization, graph theory, algebra

Compressed data structures
 Packet Processing

tapolcai@tmit.bme.hu