

کارشناسی ارشد شبکه های کامپیوتری

شبکه های پهن باند



استاد:

دکتر جبار پور

نویسنده:

مهندس علیرضا روحانی نژاد

۱۳۹۷

Trie-based Algorithms

- Binary Trie,
- Path Compressed Trie,
- Multi-Bit Trie
- Level Compression,
- Tree Bitmap Algorithm
- Tree-Based Pipelined Search,
- Binary Search on Prefix Lengths,
- Binary Search on Prefix Range.

دو روش جهت افزایش سرعت مسیریاب ها ذکر شد: روش سخت افزاری - روش نرم افزاری

روش های نرم افزاری بر مبنای درخت هستند و به آنها Tree Based Algorithm گفته می شود. در ادامه مباحث، روش های بر این مبنا را بررسی می کنیم.

Binary Tree: درخت دودویی: سرعت جستجو بالاست ولی حافظه مصرفی کاملاً به نوع درختی که داریم استفاده می کنیم، بستگی دارد. می تواند کم یا زیاد باشد. سرعت Update هم بستگی به درخت دارد.

Binary Trie

- Represents prefixes of different lengths

- 1-bit trie

- Left link: 0

- Right link: 1

- Search

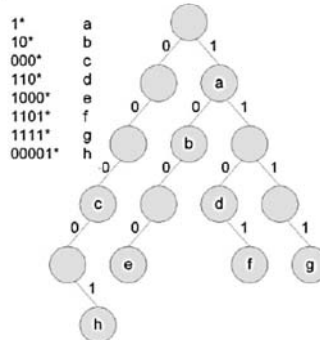
- Start from root, move to left or right if the current bit of the address is 0 or 1 respectively
- If a node containing a prefix mark (*) is seen, store it somewhere as the longest match up to now

- Addition

- Follow the path and create new nodes if needed and finally mark the last node as a prefix

- Deletion

- Follow the path and delete the last node and its parents until a marked node or a node with another child is seen



یکی از ساختارهای داده برای جستجو (Longest Match) حالت درختی است. در ساده ترین

حالت درختی دودویی داریم که اطلاعات Prefix ها را در یال ها ذخیره می کند. هر عدد روی یال

مشخص کننده عددهایی است که در جدول مسیریابی هستند. اعداد در این درخت معادل IP های

Destination هایی هستند که ما در Routing Table داریم و a, b, c, ..., h پورت نامبرهای

خروجی هستند که مشخص می کند ما از کدام پورت می توانیم بسته ها را به سمت مقصد هدایت کنیم. با

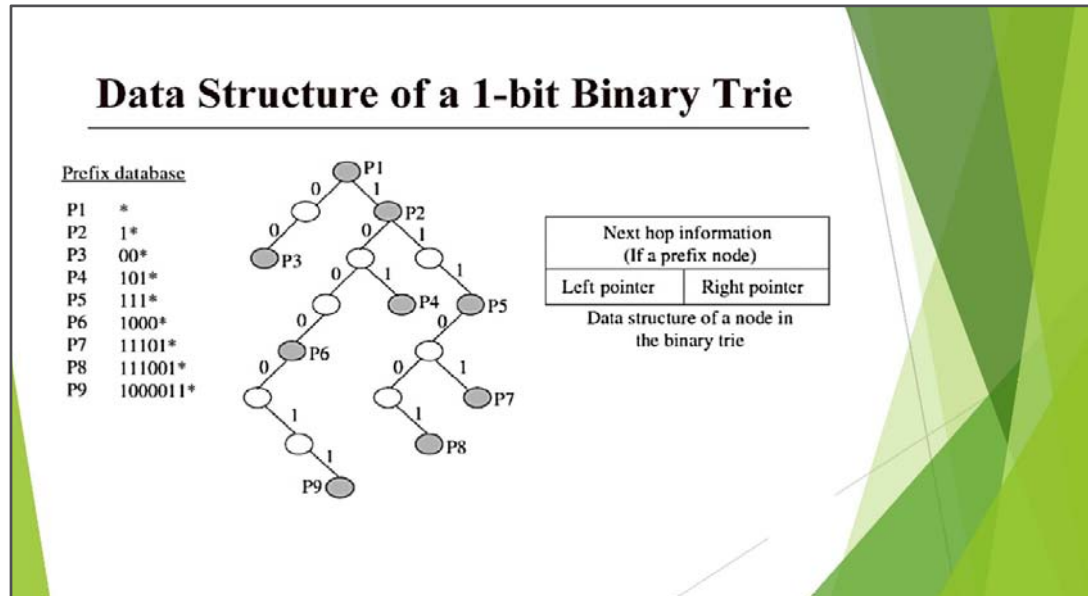
شروع از بیت پر ارزش که می تواند ۰ یا ۱ باشد، یکی از فرزندان ریشه درخت را برای پیش رفتن انتخاب می کنیم. بعد از رسیدن به آخرین بیت، اسم آن پیشوندی که در Node فعلی است را به آن اضافه می کنیم. مطابق با شکل، نام پیشوندها در Node های غیر برگ هم قرار گیرند. الزامی نیست که همیشه در برگ ها بخواهیم Prefix ها را قرار دهیم. مثلاً پیشوند b. ۱۰ است. خود آن زیر مجموعه ای از حالت e است به همین دلیل امکان تطبیق آدرس پیش رو با چند پیشوند وجود دارد و تاکید می شود در این مورد بایستی طولانی ترین پیشوند را به عنوان نتیجه در نظر بگیریم.

جستجو برای این درخت بسیار ساده است: با شروع از ریشه، با توجه به صفر و یک بودن، بیت های آدرس را مورد جستجو قرار می دهیم پس از چک کردن هر بیت به فرزند چپ یا راست آن منتقل می شویم. حین انتقال هر جا Node ی جلوی یک Prefix رسیدیم آن را به عنوان طولانی ترین تطبیق به حافظه می سپاریم تا به Node بزرگ برسیم یا امکان ادامه جستجو مقدور نباشد. آن جاست که جستجوی ما به پایان رسیده است. (مثال جلسه قبل)

وقتی می خواهیم Entry جدید اضافه کنیم، باید قبل آن پیشوند جدید اضافه شود. ابتدا باید شبیه جستجوی درخت جلو برویم و مسیر موجود را تا جایی که وجود دارد، ادامه دهیم، اگر این مسیر به طور کامل وجود داشت، کافی است نام پیشوند را در آن Node ای که در انتهاست ثبت کنیم در غیر این صورت ادامه مسیر مربوط به آن Node را در درخت ایجاد می کنیم و بعد از ایجاد کردن، Prefix ی که مدنظرمان است را در Node پایانی ذخیره می کنیم.

برای حذف (Delete) اول مسیر آن پیشوند را تا انتهایش می رویم. اگر آخرین Node مربوط به آن مسیر دارای فرزندی بود مثل پیشوند C کافی است نام C را در درخت حذف کنیم. کار دیگری نیاز نیست چون اگر بخواهیم از C به بعد را پاک کنیم، Prefix h را ازدست خواهیم داد. اگر بخواهیم Node ی را حذف کنیم که فرزندی ندارد باید خود آن Node را حذف کنیم و والد آن را هم باید حذف کنیم.

برای مثال: اگر بخواهیم h را حذف کنیم، h و Node والد را حذف کنیم به نحوی که C به برگ تبدیل شود. (چند مثال دقیقه ۱۱ کپچر)

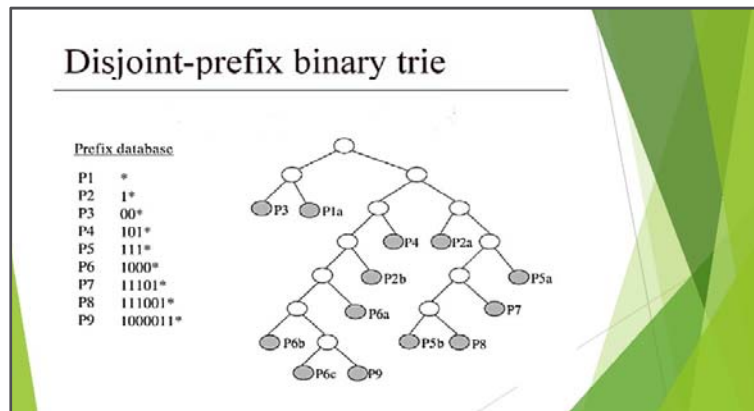


درخت باینری ۱بیتی را توضیح می دهد. به این صورت است که هر گره، دارای دو اشاره گر چپ و راست است و اطلاعات Next Hop داخل گره ذخیره می شود. این ها اطلاعاتی است که برای حالت ۱بیتی باید داشته باشیم.

Performance of Binary Trie

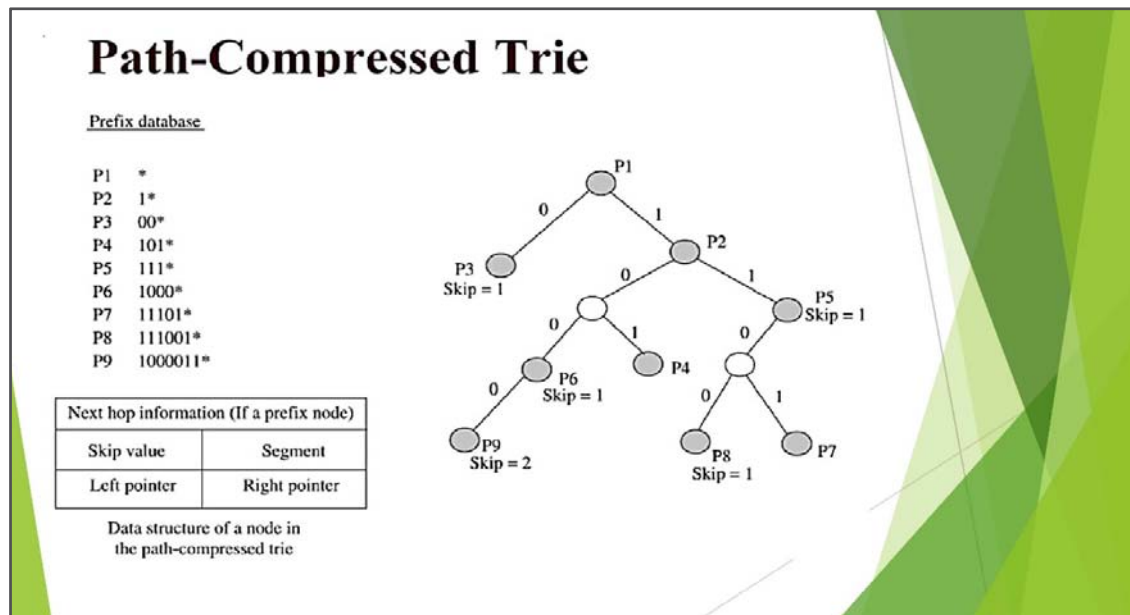
- ❑ The Number of Memory Accesses in The Worst Case is 32 for IPv4.
- ❑ To Add a Prefix to The Trie, In The Worst Case It Needs to Add 32 Nodes. In This Case, The Storing Complexity is $32N \cdot S$.
- ❑ The Lookup Complexity is $O(W)$,
- ❑ The Storage complexity is $O(NW)$

اگر بخواهیم کارایی این درخت را بررسی کنیم، گفتیم که در بدترین حالت تعداد دسترسی ها به حافظه برای جستجو حالتی است که عمق درخت ۳۲بیتی باشد چون IPv4، ۳۲بیتی است و Prefix های ما باید ۳۲بیتی باشد چون در این حالت باید ۳۲ بار به حافظه رجوع کنیم. این مشخص می کند کدام Entry مدنظر است. برای اضافه کردن در بدترین حالت، بایستی ۳۲ Node اضافه شود. هیچ مسیری وجود نداشته باشد و مجبور شویم همه Node ها را اضافه کنیم. اگر عمق درخت را W فرض کنیم، در این اسلاید که می شود طول بزرگترین Prefix مان، پیچیدگی جستجو برابر با $O(W)$ خواهد بود و حجم مصرفی ها $(N \times W)$ خواهد شد به ازای Node ی که وجود دارد ضربدر عمقی که خواهیم داشت. حال اگر IPv6 مطرح شود حداکثر عمق ما به جای ۳۲، ۱۲۸ خواهد بود.



می‌خواهیم درخت باینری که جلسه قبل درباره‌اش حرف زدیم را تغییر دهیم تا حافظه مصرفی‌اش را کم کنیم و تا حد امکان سرعت جستجو را افزایش دهیم. هرچقدر حافظه مصرفی کم شود، در حقیقت زمان نیاز به دسترسی حافظه هم کمتر خواهد بود. درخت اسلاید همان درخت دودویی است که «Disjoint» شده‌است. در حالت Disjoint، Node های میانی دیگر اطلاعات Next Hop را در خودشان نگهداری نمی‌کنند فقط گره‌های برگ هستند که اطلاعات Next Hop را ذخیره می‌کند. پس در گره‌های میانی مقداری نداریم و گره‌های برگ هستند که مقدار Next Hop را به‌خودشان گرفته‌اند. برای ایجاد این درخت، اول به Node هایی که یک فرزند دارند، یک برگ اضافه می‌کنیم، این برگ‌ها به‌عنوان Prefix مورد استفاده قرار می‌گیرند که بتواند اطلاعات Next Hop را در نزدیک‌ترین پدر خودشان که Prefix است را به ارث ببرد. حال Node های میانی که به عنوان Prefix هستند را به Node معمولی تبدیل می‌کنیم. مزیت این کار این است که در حالت قبلی هر Node ی ممکن است لازم باشد اطلاعات Next Hop را در خودش ذخیره کند. بنابراین برای همه Node ها بایستی 1 byte فضای Next Hop هم در نظر می‌گرفتیم. (در اسلاید قبلی دیدیم که اشاره‌گر سمت چپ و راست را خواهیم داشت و علاوه بر آن فضایی هم برای ذخیره Next Hop نیاز داشتیم). در صورتی که در حالت جدید Node های میانی، اطلاعات Next Hop را در خودشان ذخیره نمی‌کنند و نیازی به این فیلد ندارند. اگر درخت‌مان بزرگتر شود، شما می‌توانید (طبق محاسبات) یک سوم از حجم حافظه برای Node های میانی را حذف کنید. دلیلش این است که Node های میانی فقط باید اطلاعات فرزند چپ و راست خودشان را نگه‌دارند و برای گره‌های برگ نیاز به اشاره‌گر به فرزندها نخواهیم داشت. بنابراین می‌توانیم فضای یکی از آن اشاره‌گر برگ را برای ذخیره کردن اطلاعات Prefix استفاده کنیم.

مشکل این روش: در این روش تعداد برگ‌های ما زیاد خواهد بود بنابراین تعداد Node های ما بیشتر خواهد بود ولی Node ها ساده‌تر شده و فضای کمتری را اشغال خواهد کرد. در مجموع به کمتر شدن فضای خالی حافظه‌مان می‌تواند کمک کند. پس برای تبدیل کردن درخت به Disjoint باید درخت را کامل کنیم. همه Node هایی که فرزند چپ یا راست ندارند را به آن‌ها باید فرزند چپ و راست بدهیم. باید گره‌ای به آن اضافه کنیم که بشود گره بزرگ، حالا مقدارش را چه بدهیم؟ از والد یا پدر خودش به ارث می‌برد. (مانند P1a در اسلاید). (از نزدیک‌ترین والد خودشان به ارث می‌برند).



می‌توانیم درخت را کمی فشرده‌تر کنیم تا باز هم حجم را کاهش دهیم. از Path

Compressed Tree استفاده می‌کنیم. فشرده‌سازی مسیرها در درخت اولیه، در برخی نقاط

مسیرهایی‌ست که فقط به یک Prefix ختم می‌شوند. مثلاً در درخت اولیه مسیر 00 را با هم چک کنیم:

فقط به P3 ختم می‌شود بنابراین اگر این دو بیت در آدرس مقصد، 00 باشد را می‌توانیم پیمایش کنیم و

اگر 01 باشد، از اول کار معلوم است که نمی‌توانیم پیمایش کنیم و بایستی شاخه دیگری پیدا کنیم. یا

مثلاً اگر در درخت اولیه از P6 تا F6 داریم، بنابراین یا ترتیب 011 را داریم که می‌توانیم پیمایش کنیم

یا نمی‌توانیم. برای مواردی که فقط یک مسیر وجود دارد، می‌توانیم Nodeهای اضافی را حذف کنیم و

به این ترتیب درخت‌مان را فشرده کنیم. از روی چند بیت پریده‌ایم و این بیت‌ها چه مقداری دارند برای

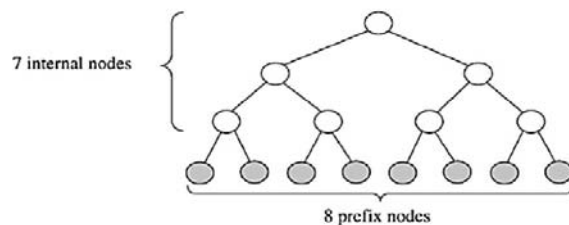
ما مهم است. در ساختمان داده‌مان، دو فیلد جدید مطابق با اسلاید ۱ اضافه می‌کنیم. این دو فیلد برای

درخت‌هایی که خلوت هستند و شاخه‌های طولانی دارند، می‌توانند به کم کردن حافظه کمک کند.

در اسلاید به P3 رسیدیم، Nodeی بوده که پریده‌ایم، Skip مساوی ۱ شده است یا برای P6 همینطور.

با اضافه کردن Skip Value نشان می‌دهیم چه اتفاقی رخ داده و چگونه درخت کوچک شده‌است.

Example of Path-Compressed Trie with N Leaves



نشان می‌دهد که اگر ۸ نود Prefix داشته باشیم، ۷ نود میانی خواهیم داشت که می‌توانیم آن‌ها را فشرده‌تر کنیم اگر آن Prefix مقدار یکسانی داشته باشد.

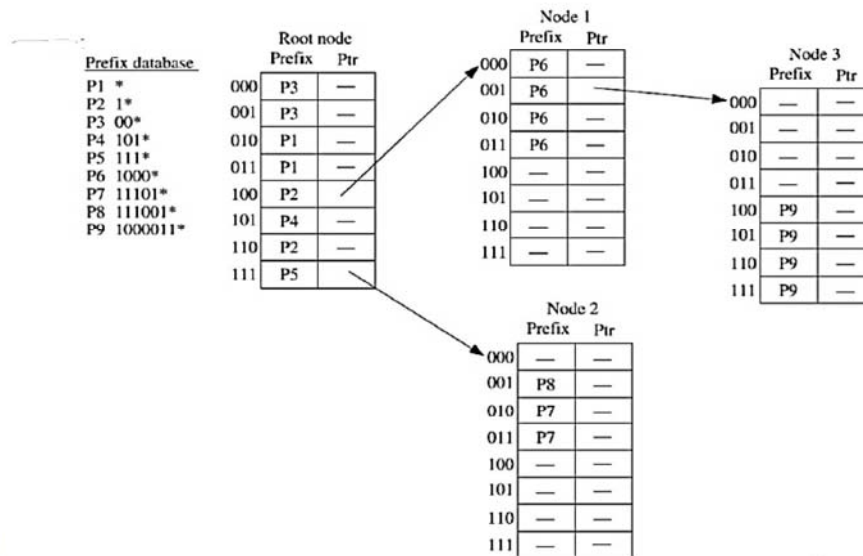
Performance of Path-compressed Tire

Path compression reduces the height of a sparse binary trie. When the tree is full and there is no compression possible, a path-compressed trie looks the same as a regular binary trie. Thus, its lookup and update complexity (the worst case) is the same as a binary trie, $O(W)$. Considering a path-compressed trie as a full binary trie with N leaves, there can be $N - 1$ internal nodes between the root and each leaf node (including the root node). Since the path can be significantly compressed to reduce the internal nodes, the space complexity becomes $O(N)$, independent of W .

کارایی روش Path Compress نشان داده شده‌است. برای درخت‌هایی که Sparse هستند و

خلوت هستند، کارایی زیادی دارد ولی برای درخت‌های کامل کارایی مناسبی ندارد چون هیچ فشرده‌سازی ایجاد نمی‌کند پس کارایی این روش در بدترین حالت با روش باینری ساده‌ای که صحبت شد، هیچ تفاوتی نخواهد داشت. چون گره‌های میانی را فشرده می‌کنیم. $O(N)$ برای مصرف حافظه می‌تواند باشد. (۷ نود را می‌توانیم در مثال Compress کنیم).

Multi-Bit Trie



به جای ۱ بیت می توانیم از چند بیت استفاده کنیم. در درخت ۱ بیتی می شود در هر بار فقط یک بیت را چک کرد. در این روش به جای ۱ بیت می توان چند بیت را همزمان چک کرد. در این شکل یک حالت ۳ بیتی را ایجاد کرده. اصطلاحاً با هین تعداد بیت هایی که چک می کند شاخه شاخه شدن Stride گفته می شود. اینجا ۳ است.

با ۳ بیت $2^3 = 8$ حالت می توان ساخت.

۸ حالت را در اسلاید برای ما مشخص کرده است. هر Node آن Tree که شامل ۸ عدد Entry است را ایجاد می کنیم Prefix هایی که طول کمتر از Stride را دارند. به عنوان مثال همه این Prefix ها طولشان از ۳ بیشتر نیست. $(P1, P2, P3)$ به حالت ۳ بیتی باید آن را توسعه دهیم. مثال: اگر $P2$ را در نظر بگیریم، اولش ۱ است، بعد دو بیت بعدی ۴ حالت را می تواند داشته باشد، باید به سه بیت تبدیل شود. می شود: 111 110 101 100 تا سه بیت شود.

حالت هایی که با Prefix طولانی تری Match شود، مثلاً با $P4$ و $P5$ آن ها را به همان اسم می گذاریم ولی مابقی را به $P2$ تبدیل می کنیم. سه صفر بیشترین Match را با $P3$ دارند. آن را می نویسیم. یا

100 بیشترین تطبیق را با P2 دارد چون فرمان دوصفر اضافه کرده ایم. وقتی این ها را پر کردیم، حالت هایی که مشترک هستند مثل 100، هم P2 هم P6 و هم P9 را نشان می دهد.

100 می تواند P2، P6 و P9 باشد. با استفاده از Node های بعدی مشخص می کنیم کدام یک از این Prefix ها مدنظر ماست. مثلا اگر 100 بود \leftarrow P2

اگر 000 شد $\xleftarrow{\text{Longest Match}}$ می شود حالت P6

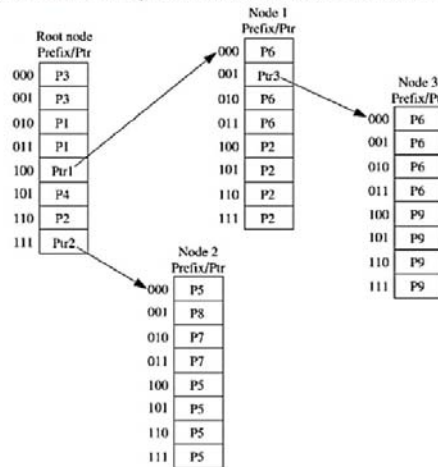
در درخت ایجاد شده در این اسلاید هر سه بیت را در یک گام می توانیم بررسی کنیم. چون هر Node درخت را یکبار از حافظه می خوانیم، اگر اندازه ی ماکزیمم Prefix ها برابر با W بیت باشد و Stride ها K باشد، تعداد دفعات مراجعه به حافظه حداکثر برابر W/K خواهد بود. در هر بار می توانیم K بیت را بررسی می کنیم. پس پیچیدگی جستجو می شود $O(W/K)$. سرعت می تواند نسبت به حالت تک بیتی K برابر باشد.

مشکل این است که ما مجبوریم یکسری Prefix ها را تکرار کنیم، مثلا دیدیم که P2 و P6 را چندین مرتبه تکرار کردیم دلیلش این است که می خواهیم Prefix ها را مضربی را از K بکنیم که ما را دچار مشکل می کند. K هر چه شود براساس آن همه Prefix ها آن را شامل می شود. این کار حافظه مصرفی ما را افزایش می دهد. در بدترین حالت حجم می تواند به 2^k افزایش یابد.

اگر مثل Disjoint Prefix ها را در K، Push کنیم، درخت ما تبدیل به چنین حالتی خواهد شد، حالت Pointer ها را نیاز نخواهیم داشت. P2 نوشته نمی شود، Pointer می کنیم به خود شماره ۱. براساس سه Node بعدی تصمیم می گیریم اگر 000 بود P6 می شود، اگر سه Node بعدی 001 بود مجددا Pointer دارد و همینطور تا آخر. (حالت قبلی گفتیم ۱ اضافه می کنیم. 10 و 11 مشخص کننده این است که کدام حالت را پیش روی خواهیم داشت)

Multi-Bit Trie

Example With Each Entry a Prefix or a Pointer to Save Memory Space



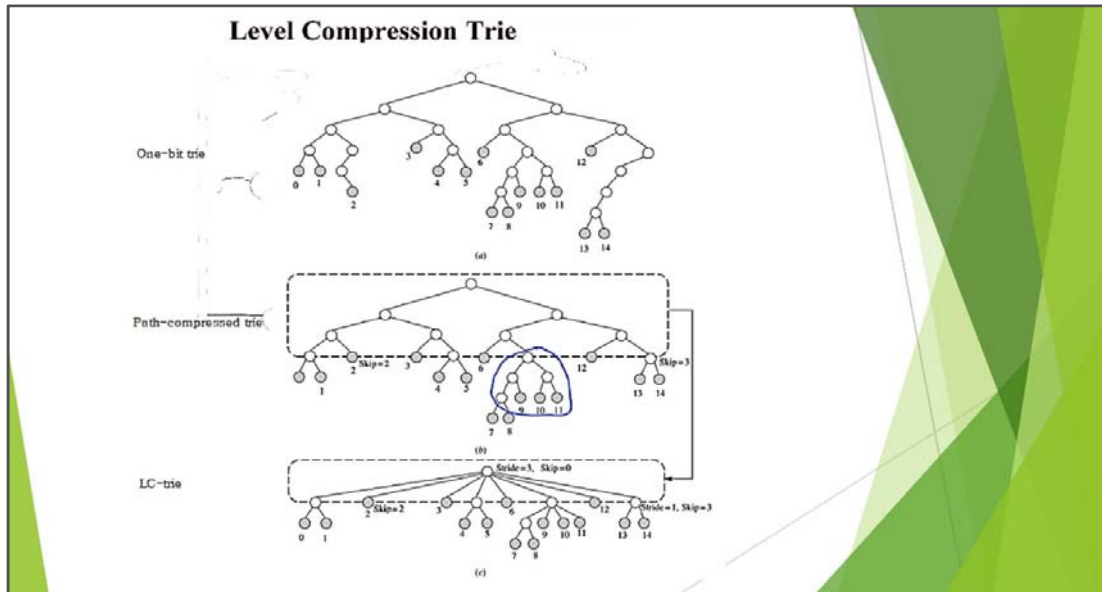
Performance of Multi-Bit Trie

The advantage of the k -bit trie structure is that it improves the lookup by k times. The disadvantage is that a large memory space is required. One way to reduce the memory space is to use a scheme called 'leaf pushing'.

The lookup is performed in strides of k bits. The lookup complexity is the number of bits in the prefix divided by k bits, $O(W/k)$. For example, if W is 32 and k is 4, then 8 lookups in the worst case are required to access that node. An update requires a search through W/k lookup iterations plus access to each child node (2^k). The update complexity is $O(W/k + 2^k)$. In the worst case, each prefix would need an entire path of length (W/k) and each node would have 2^k entries. The space complexity would then be $O((2^k * N * W)/k)$.

با حذف کردن Pointerها حجم حافظه را نصف کردیم، پیچیدگی این روش $O(W/K+2^k)$:

k ای که انتخاب می شود برای ما مهم است. پیچیدگی جستجو W/K خواهد بود. بروزرسانی هم همینطور و 2^k برای بدترین حالت باید بررسی شود. برای همین در بدترین حالت، پیچیدگی $W/K+2^k$ خواهد شد. در انتخاب K بایستی Trade Of ایجاد کنیم، اگر K کوچک باشد، تعداد مراجعه به حافظه زیاد می شود ولی حجم حافظه کم. چون سه بیت را چک می کنیم و اگر K بزرگ شود تعداد دفعات مراجعه شده ولی حجم حافظه بیشتر می شود.



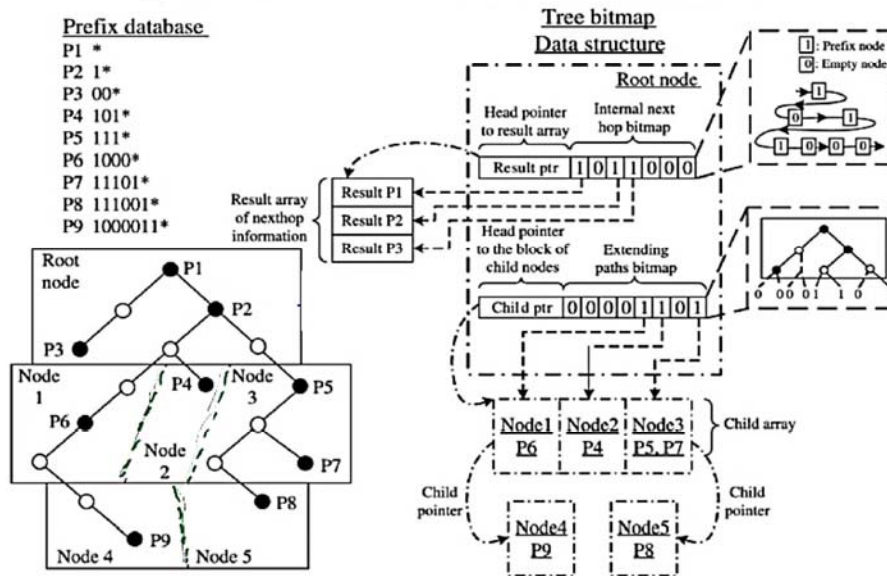
روشی داریم به نام Level Compress Tree که بر اساس 1bit tree و Path Compression هم اضافه کرده است. در این روش، روش چند بیتی و فشرده سازی را به صورت همزمان استفاده کرده می کند. هر جایی که شاخه Sparse دارد شاخه هایش پر نیست، فشرده سازی کرده و هر جا که شاخه ها کامل هستند به صورت چند بیتی کار می کند. مطابق با اسلاید ۳ سطح اول را که کامل هستند به روش ۳ بیتی عمل می کند و تعداد زیادی از Node ها می توانیم حذف کنیم. در سطح بعدی شاخه ها را به روش فشرده سازی انجام داده. برای شاخه هایی که تنها هستند. (در انتها ۳ Node را Skip کرده تا به اینجا رسیده است). Skip: 1 K=1 (Stride)

Performance of Level Compression Trie

An LC-trie searches in strides of k bits, and thus the lookup complexity of a k -stride LC-trie is $O(W/k)$. To update a particular node, we would have to go through W/k lookups and then access each child of the node (2^k). Thus, the update complexity is $O(W/k + 2^k)$. The memory consumption increases exponentially as the stride size (k) increases. In the worst case, each prefix would need an entire path of length (W/k) and each node has 2^k entries. The space complexity would then be $O(2^k * N * W/k)$.

فشرده سازی در بدترین حالت شبیه حالت چند بیتی می شود. اشکالی که وجود دارد این است که ساختار پیچیده تر می شود چون دو روش را داریم با هم ترکیب می کنیم. همان تعاریف اینجا هم پا بر جا هستند.

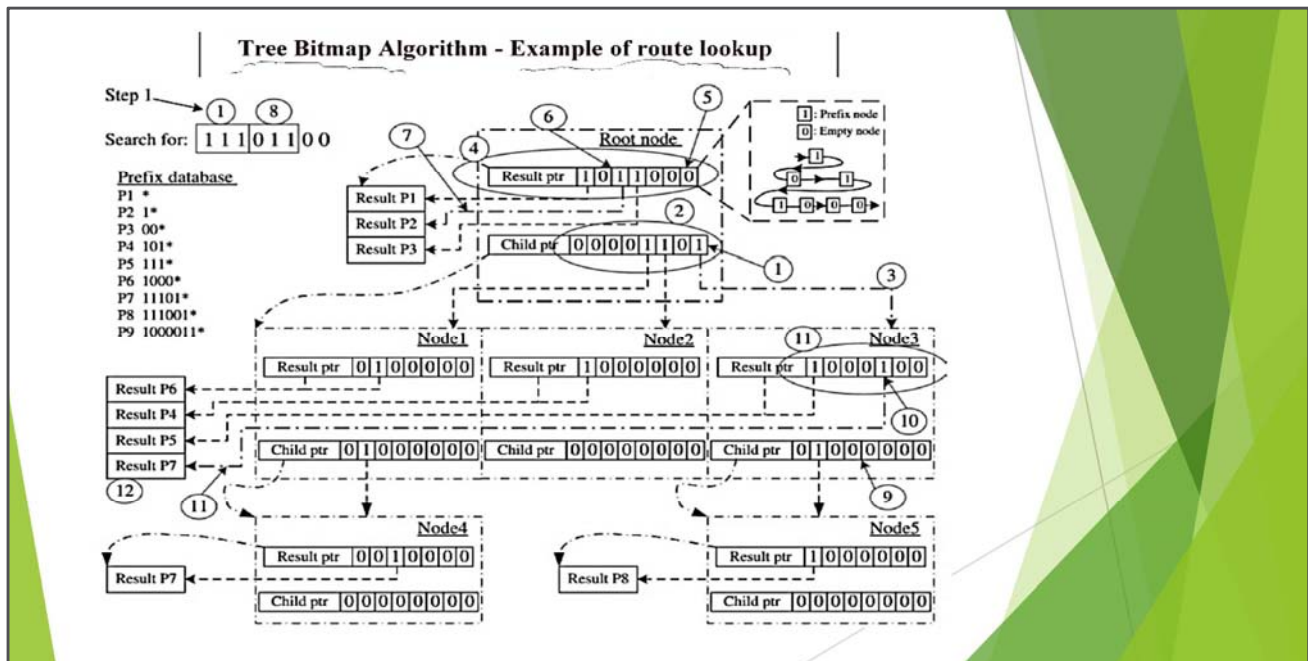
Tree Bitmap Algorithm



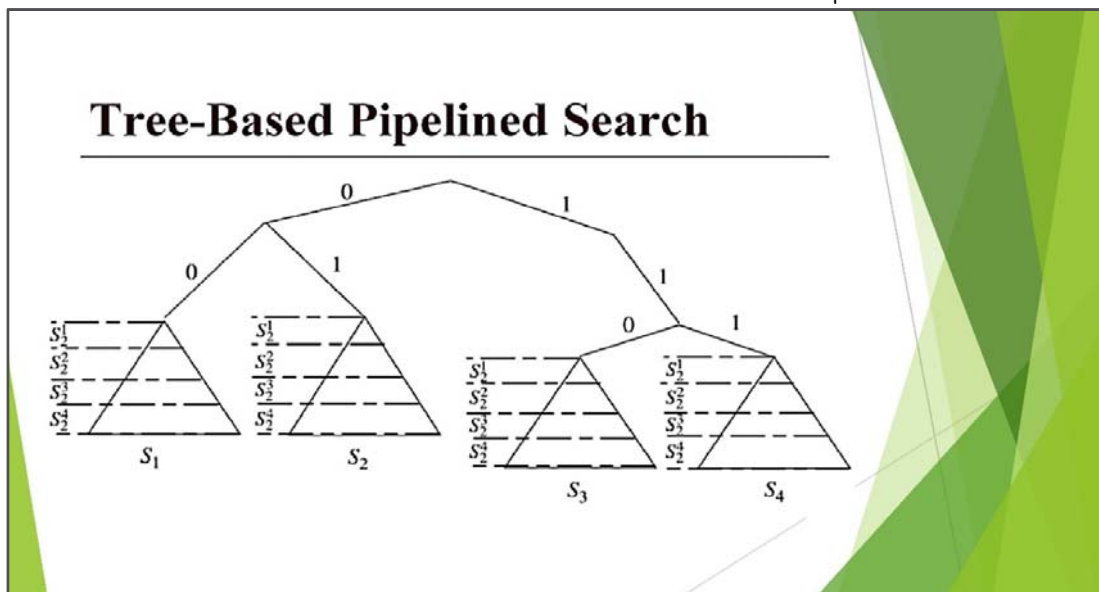
روش Tree Bitmap: از روی درخت باینری یک Bit Map تهیه می کنیم. بجای این که درخت را نگهداری کنیم Bit Map را نگهداری می کنیم. به عنوان مثال اگر St Ride برابر ۳ بگیریم، سه سطح اول از ریشه به عنوان یک Node از درخت جدا می شود. دوباره از گره بعدی سه سطح جدا می کنیم. برای گروه های بعدی هم همین اتفاق می افتد. اگر بخواهیم سه بیت را چک کنیم، سطح اول می شود یک Node، سه تایی بعدی می شود یک Node و سه تایی بعدی می شود یک Node. با این روش درخت به یک Node ریشه می آید و به پنج Node دیگر تقسیم می شود. یک گره ریشه دارد و پنج Node را شامل خواهد بود. برای هر Node دو عدد Bit Map نگهداری می کنیم. مثل مثال سمت چپ اسلاید. (Result & Child Pointer) را نشان می دهد.

دو بیت مپ نگهداری می کنیم یعنی: در درخت باینری چه Node هایی Prefix جدول هستند و چه Node هایی Prefix نیستند. P1, P2, P3 نیستند. شماره گذاری شان هم به این شکل انجام دهیم که هر سوپر Node درخت، حداکثر می توانیم ۷ گره از درخت اصلی را داشته باشیم. این درخت را کامل کنیم. (در هر سطح در Node ریشه می توانیم ۷ Node را داشته باشیم)

فرض کنید همه ۷ گره وجود داشته باشند، برای Prefix ها عدد ۱ قرار می دهیم، برای گره های معمولی و خانه های خالی صفر را قرار می دهیم. (همان اتفاقی که در شکل نشان داده شد). (دقیقه ۴۰ کیچر)

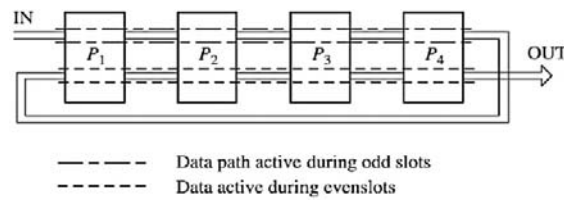


می‌توانیم ساختار قبلی را به صورت کاهش حجم و سادگی به این صورت نمایش داد. هر ۷ گره معمولی یک سوپر نود ایجاد می‌کنیم که فضای ۲ بیت Map ۷ بیتی و ۸ بیتی را خواهد داشت، به این صورت دو Pointer هم به آن اضافه خواهد شد.



وقتی حافظه ای داریم و روی آن عمل Look Up انجام دهیم در هر لحظه، فقط می‌توانیم یک جستجو انجام دهیم. ولی اگر حافظه یک درخت را به چندین بخش اختصاص دهیم و کپی کنیم، می‌توانیم از تکنیک Pipelined استفاده کرده و جستجوی متفاوتی انجام دهیم.

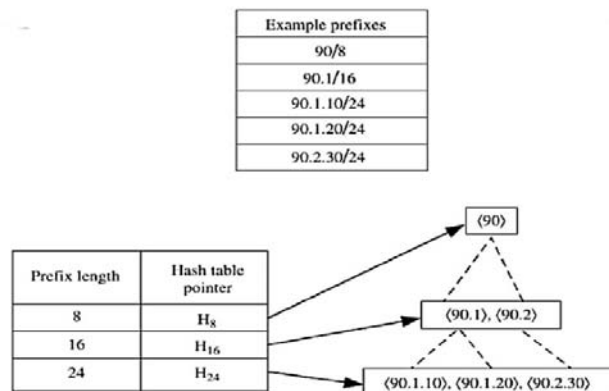
Tree-Based Pipelined Search



Random ring pipeline architecture with two data paths

یک درخت جستجو را به ۴ قسمت تقسیم کرده، یک بخشی از جستجو در حافظه اول انجام می شود، نتیجه در بخش دیگر. حال می توانیم به صورت همزمان یک جستجوی جدیدی در بخش اول انجام دهیم به همین صورت می توانیم چهار جستجو را به صورت Pipeline انجام دهیم و سرعت را ۴ برابر کنیم. فرض کنید دنبال Prefix ای هستیم. ۱،۲،۳،۴ منظور این است که اگر بتوانیم حافظه را به ۴ بخش تقسیم کنیم می توانیم به صورت متوالی این عملیات را انجام دهیم به شرط این که این سطوح مختلف درخت را بتوانیم انجام دهیم. جستجو در یک سطح تمام شد. ادامه جستجو در سطح بعدی امکان پذیر است.

Binary Search on Prefix Lengths



مثالی نمایش یافته است. سطح بندی براساس طول Prefix انجام شده است. طول ۲۴:۱۶:۸ و یک Pointer هم به سه تا جدول ایجاد کرده است. اگر Prefix، ۸ باشد، ۱ عدد Entry بیشتر نداریم. در حافظه آمده است که به صورت مستقیم آن را سرچ می کند و دیگر نیازی نیست دوتای بعدی را سرچ کند. اگر ۱۶ است. ۹۰.۲ و ۹۰.۱ را قرار داده اگر ۲۴ باشد ۱۰۲۰۳۰ را مشخص کرده به این صورت براساس Prefix معلوم می شود کدام جدول بررسی شود. یک سری Pointer هم اضافه شده. مثلا اگر دنبال ۹۰.۲. ۳۰. ۲۰ IP هستیم به صورت سلسله مراتب چک می کند. ۹۰ چک می شود ← ۹۰.۲ و بعد ۹۰.۲. ۳۰ می شود Longest Match و ادامه کار.

Binary Search on Prefix Lengths

Performance. The algorithm requires $O(\log_2 W)$ hashed memory accesses for one lookup operation, taking no account of the hash collision. So does the update complexity. This data structure has storage complexity of $O(NW)$ since there could be up to W markers for a prefix-each internal node in the trie on the path from the root node to the prefix. However, not all the markers need to be kept. Only the $\log_2 W$ markers that would be probed by the binary search algorithm need be stored in the corresponding hash tables. For instance, an IPv4 prefix of length 22 needs markers only for prefix lengths 16 and 20. This decreases the storage complexity to $O(N \log_2 W)$.

می دانیم که روش Binary Search، سریعترین روش جستجو است. (نصف کردن داده ها به شرط Sort بودن، پیدا کردن میانه، حذف یک طرف و...) بنابراین $O(\log(n))$ بود پیچیدگی الگوریتم. در این روش یک شبیه سازی جستجوی دودویی مدنظرمان است. در جستجوی دودویی باید اطلاعات به صورت Sort شده باشند. حال زمانی که می خواهیم عددی پیدا کنیم فضای جستجو را به دو قسمت تقسیم می کنیم. برای P1.3 جدول Range و Prefix را تشکیل داده. Range هایی را درست کرده و در اشکال B و C نشان داده است.

از حالت 100000 شروع کرده و این Range بندی ها را ادامه داده در هر بار اندازه فضای Search نصف می گردد. از این ایده در Routing Table می توانیم استفاده کنیم. اشکال این است که اطلاعات جدول عدد نیستند، بلکه Range ی را برای ما مشخص می کنند. می توانیم شروع و پایان رنج ها را مشخص کنیم. هر IP یک رنج برای ما مشخص می کند. آدرس مقصد که در یکی از این رنج های مرتب شده قرار گرفت، بسته متعلق به آن زیر شبکه خواهد بود. اگر حالت های منطقی برای رنج هایش بیابد مثلاً بزرگتر مساوی بایستی متناسب آن پورت را تعیین کنیم. رنج های شروع و پایان را که می نویسیم در اصل داریم اندازه جدول مسیریابی را دو برابر می کنیم. بعلاوه یکسری اطلاعات را هم به Entry ها اضافه می کنیم. در این روش داریم اندازه حافظه را بزرگتر می کنیم اما Search ها سریعتر خواهد شد.

معماری مسیریابها و سوئیچ های با سرعت بالا

IP Address Lookup
(Hardware-based Schemes)

Hardware-based Algorithms

- ❑ DIR-24-8-BASIC Scheme
- ❑ DIR-Based Scheme with Bitmap Compression (BC-16-16)
- ❑ Ternary CAM for Route Lookup
- ❑ The Algorithms for Reducing TCAM Entries
- ❑ Reducing TCAM Power – CoolCAMs
- ❑ TCAM-Based Distributed Parallel Lookup

روش های سخت افزاری جستجوی اطلاعات مسیریابی:

روش های سخت افزاری، سرعت بسیار بالاتری دارند. اما هزینه پیاده سازی بالاست و انعطاف پذیری (Flexibility) کمتر است.

به طور کلی روش های سخت افزاری دو دسته اند:

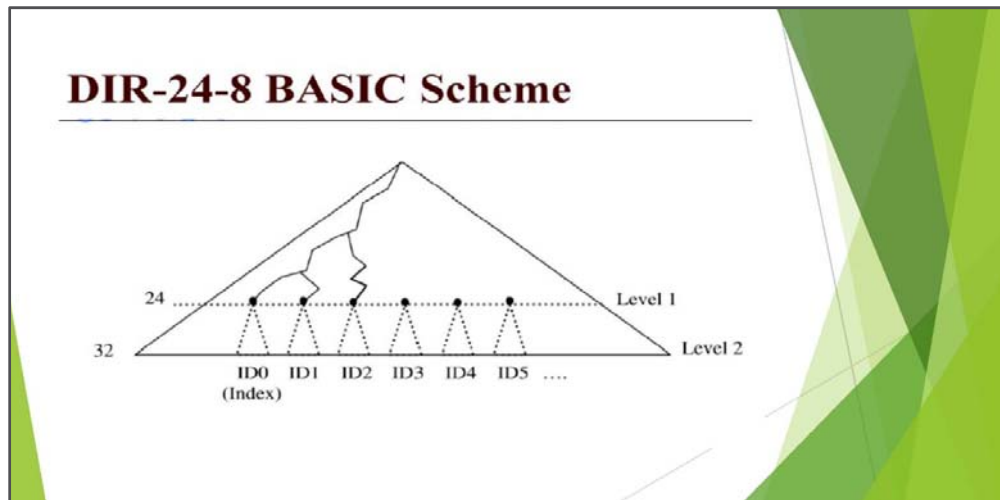
۱. DIR

۲. Ternary

روش DIR (Direct): مستقیم: با یکبار مراجعه به حافظه، جستجو را انجام می دهد. حجم حافظه مصرفی بالاست. هزینه Update کردن جداول بالاست.

روش Ternary: مبتنی بر Content Addressable Memory است.

این ها براساس محتوا، آدرس دهی را انجام می دهند.



روش های Direct دو شماره دارند. مثلاً DIR-24-8 یعنی یک جستجوی دو سطحی خواهیم داشت درخت جستجو، درختی با عمق ۳۲ است. اگر همه برگ ها را در یک حافظه ذخیره کنیم، نیازمند $2^{32}=4GB$ هستیم. می توانیم آدرس مقصد بسته را به عنوان آدرس حافظه استفاده کنیم و به خانه حافظه اشاره کنیم در محتوای حافظه هم Next Hop را ذخیره کنیم. به این ترتیب با یکبار مراجعه به حافظه می توان عمل Look Up را انجام داد. اما از نظر هزینه پیاده سازی، حدود 4GB حافظه نیاز خواهیم داشت.

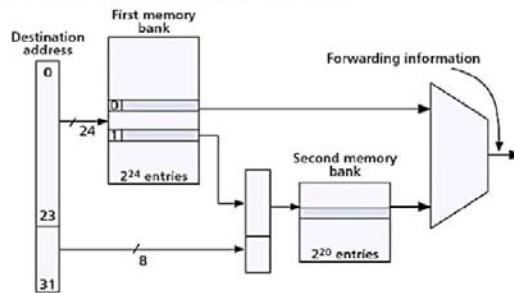
اگر نیاز به Update جداول باشد، هزینه بیشتر خواهد شد چون هر تغییر در جدول مسیریابی، می توان سبب هزاران تغییر درجه اول Direct باشد. (چون Pointerها هم مورد اشاره اند).

طبق اطلاعات آماری که از جداول مسیریابی به دست آمده است مشاهده شده که اکثر Entryها Prefix کمتر از ۲۴ دارند و تعداد کمی از Prefixها طولشان از ۲۴ بیشتر است. به همین دلیل دو سطح ۲۴ و ۸ تعریف شده است. پس یک راه حل این است که سطح ۲۴ اطلاعاتش ذخیره شود. تعداد Nodeها در این حالت 2^{24} یعنی ۱۶ مگابایت که در برابر ۴ گیگ بسیار کمتر است. حال اگر Prefix کمتر از طول ۲۴ داشته باشیم، در عمق حداکثر ۲۴ می شود Next Hop را پیدا کرد. اگر Prefixهایی با طول بیشتر از ۲۴ داشته باشیم، بایستی عمل جستجو را بیش از ۲۴، ادامه دهیم و تا عمق ۳۲ این امکان است که این اتفاق بیفتد.

برای Prefixهایی که طول بزرگتر از ۲۴ دارند، یک حافظه جداگانه در نظر بگیریم. از آنهایی که تعداد اینها کم است، نیازی به حافظه بزرگی نداریم پس ساختار و معماری مان شبیه شکل اسلاید ۴ می شود.

DIR-24-8 implementation

- Gupta et al.
- Two levels
 - First memory bank: 24 bits of address
 - Second memory bank: 8 bits of address



پیاده سازی سخت افزاری یک Tree چند سطحی نشان داده شده است. دو سطح دارد: یک سطح همان ریشه می توان فرض کرد با ۲۴ بیت و یک سطح ۸ بیت باقی مانده را تصمیم گیری می کند.

یک بانک حافظه با 2^{24} Entry را در نود ریشه پیاده سازی کرده و براساس ۱۴ بیت اول به حافظه مراجعه می کند با یافتن طولانی ترین تطبیق Prefix (Longest Match) مورد جستجو را پیدا می کند یا به اشاره گر نود بعدی، (بانک دوم) اشاره می کند و به آن ارجاع می دهد. ۸ بیت باقی مانده سراغ بانک دوم می رود، آنجا هم Longest Match را بایستی انجام دهیم تا مشخص شود که چه Prefix مناسب است. برای این که Throughput مان بهتر شود می توان دو بانک حافظه ای را به صورت Pipe (لوله ای) کنار هم قرار داد و برای استفاده از حافظه های پویا که زمان دسترسی می تواند 50ns باشد، می توان در حدود ۲۰ میلیون جستجو در ثانیه انجام دهد.

مشکل پیاده سازی این روش این است که برای بسط دادن آن ممکن است به تعداد زیادی خانه حافظه، نیاز پیدا کنیم و مجبور شویم خانه های حافظه را زیاد تغییر دهیم. ضمناً به علت حجم حافظه بالای مورد نیاز، امکان پیاده سازی آن در حافظه های ایستا (که از پویا سریع تر است) وجود ندارد. می توانی بانک ها را سه سطحی کنیم تا حجم حافظه هر بانک، کمتر شود. هر چند به بدترین حالت جستجو یک دسترسی به حافظه افزوده خواهد شد. به جای دوبار، سه بار مراجعه به حافظه. آن وقت خواهیم داشت: DIR-21-3-8

۲۱ بیت اول بعد ۳ بیت و در نهایت ۸ بیت.

DIR-24-8 implementation

TBL24 entry format.

If longest route with this 24-bit prefix is < 25 bits long :

0	Next Hop
1 bit	15 bits

If longest route with this 24-bit prefix is > 24 bits long :

1	Index into 2nd table
1 bit	15 bits

در این مثال محتوای حافظه سطح اول را نشان می‌دهد که اگر بیت اول صفر باشد یعنی Next Hop در آن ذخیره شده و اگر بیت اول ۱ باشد محتوای آن Index ی است که به جدول دوم اشاره می‌کند. می‌توان برای خانه‌های حافظه از آن استفاده کرد.

DIR-24-8 implementation

Example of two tables containing three routes.

Key to table entries
A = 10.54/16
B = 10.54.34/24
C = 10.54.34.192/26

TBL24			TBLlong		
Entry Number :	Contents		Entry Number :	Contents	
10.53.255			123*256		B
10.54.0	0	A	123*256+1		B
10.54.1	0	A	123*256+2		B
10.54.33	0	A	123*256+191		B
10.54.34	1	123	123*256+192		C
10.54.35	0	A	123*256+193		C
10.54.255	0	A	123*256+255		C
10.55.0			124*256		C

256 entries allocated to 10.54.34 prefix

همان مسئله تکرار شده است اگر خانه اول صفر باشد یعنی A، Next Hop ها است ولی اگر ۱ است مشخص کننده آدرس یا Index در جدول حافظه ها است. حال در این مثال فرض کرده ۳ عدد Entry در جدول ها باشد، اول باید رکوردهایی که طول کوچکتر از ۲۴ بیت دارند را ۲۴ بیتی کنیم. مثلا 10.54 باید به 10.54.0 تا 10.54.255 یعنی به آن Range دهیم چون کمتر از ۲۴ است و ما داریم درباره DIR248 حرف می زنیم باید حداقل به ۲۴ افزایش دهیم. در همه حالات به غیر از حالات B و C، Next Hop ها A خواهد بود. برای حالت B و C صرفا با استفاده از 10.54.34 نمی توانیم تصمیم بگیریم Next Hop چیست. بایستی آدرسی به Index حافظه سطح دوم ایجاد کنیم و با بیت های باقی مانده (۸ بیت) تصمیم بگیریم، Next Hop چه خواهد شد. با ۸ بیت باقی مانده شروع کرده و تا ۲۵۶ جلو رفته است. تمامی حالات را در حافظه سطح دوم گنجانده است.

بسته ای گرفتیم با آدرس 10.54.34.185 ابتدا ۲۴ بیت اول را چک می کند حافظه سطح اول Index به ما می دهد که باید سراغ جدول سطح دو برویم. ← سراغ ۸ بیت باقی مانده ← ۸ بیت باقی مانده 185 را نشان می دهد. 185 مشخص کننده B خواهد بود. می بینیم که یک Entry داخل جدول روی تعداد زیادی از خانه ها در سطح اول و دوم تکرار می شود و این ایرادی است که در این روش وجود دارد و باعث افزایش حجم حافظه می شود.

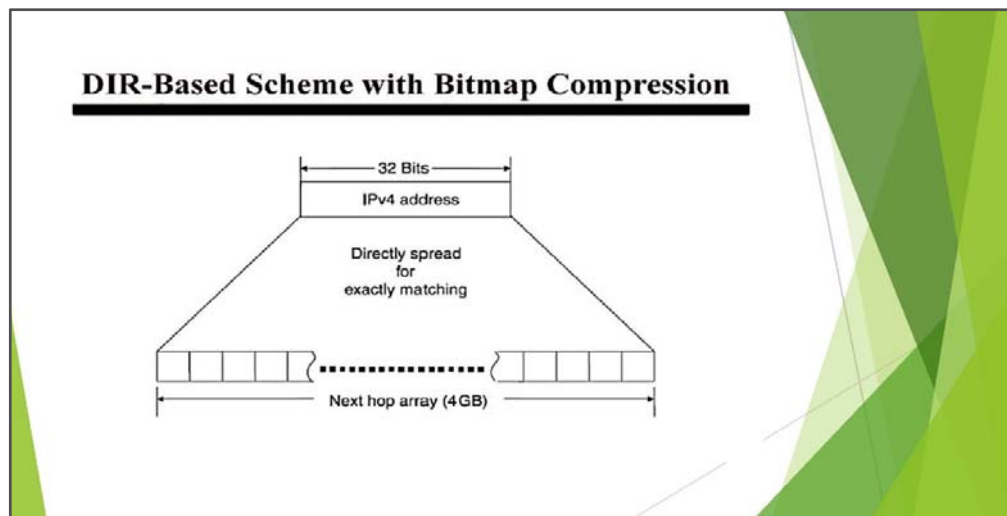
سوال: چرا بعد از ۲۵۶ هنوز ادامه دارد؟ چون داریم درباره Longest Match صحبت می کنیم. بعد از آن هر چه باشد، Longest Match می شود به همین دلیل تا C ادامه دارد. همه حالت ها پورت A هستند. فقط 10.54.34 مشکل ساز خواهد شد که برای آن Index تعریف کرده تا بتواند آن را در جدول مشخص کند. Index را ۱۲۳ در نظر گرفته است. 156.1.2 تا برسد به ← ۱۹۱ تا ۱۹۱، B پورت خروجی خواهد بود و از آن جا به بعد پورت C برای خروجی در نظر گرفته می شود.

DIR-24-8 implementation

- Performance
 - Two pipelined memory accesses per lookup
 - DRAM delay of 50ns => 20 mps
 - 33 Mbytes of DRAM
- Drawbacks
 - High memory usage
 - Many memory places may need to change for an update

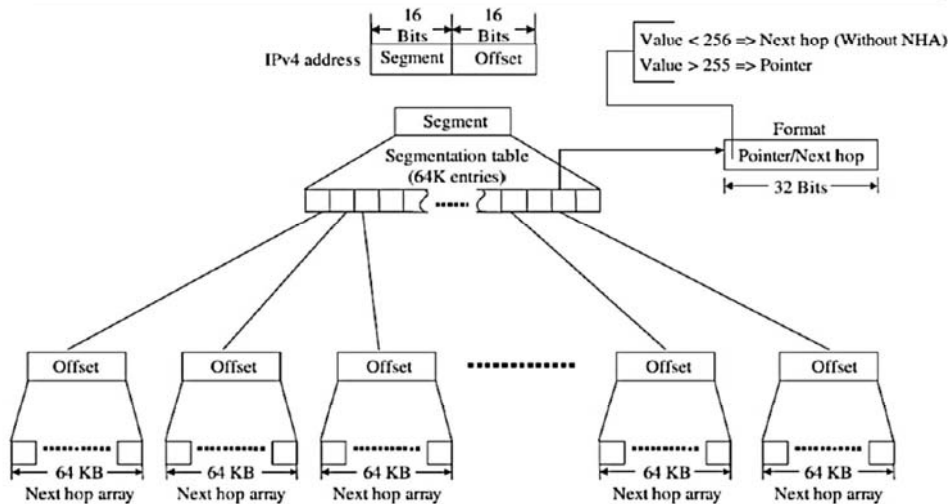
در این روش جستجو با یک یا دوبار به حافظه انجام می شود. ولی اگر بخواهیم تغییراتی را در جدول بدهیم مجبور خواهیم بود تعداد زیادی از خانه های حافظه را تغییر دهیم.

اشکال اصلی: حافظه مصرفی زیاد و تغییرات زیادی برای Update نیاز است.

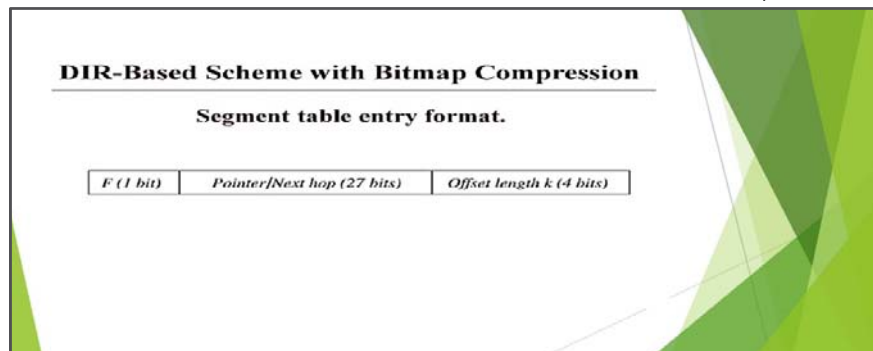


روش دیگری داریم به نام ذکر شده در اسلاید که به آن DIR-16-16 هم گفته می شود.

DIR-Based Scheme with Bitmap Compression



از ایده Direct دو سطحی در آن استفاده شده است. اما سعی شده اندازه عمق حافظه سطح ۲ دیگر ثابت نباشد. در این روش تا ۱۶ ادامه می دهد. بعد درخت شکسته می شود و همه Prefix ها را با ۱۶ مشخص می کنیم. بعضی از Prefix ها که طول بزرگتر از ۱۶ بیت داشته باشند، عمق بزرگترین شاخه را در نظر می گیریم، حداکثر تا ۳۲ می تواند برسد. (کمتر بود تا همان جا را نگه می داریم) در حافظه دوم به جای این که به ازای هر شاخه 2^{16} تا حافظه را در سطح ۲ در نظر بگیریم، به اندازه ی Max عمق در نظر می گیریم. (عمق ثابت نیست). عمق K بیت باشد، 2^k خانه حافظه در نظر گرفته می شود. به این ترتیب در حافظه سطح اول یا Next Hop پیدا خواهد شد یا Pointer به حافظه سطح دوم خواهیم داشت. اگر Pointer به حافظه دوم داشته باشیم یک حافظه ۴ بیتی را در نظر می گیریم که عمق سطح دوم مان را برآیمان مشخص کند.



یک بیت مشخص شده اگر Next Hop است، که هیچ، اگر Pointer است Offset ی که در ۴ بیت بعدی است مشخص کننده خواهد بود.