

# Space and Speed Tradeoffs in TCAM Hierarchical Packet Classification

Alexander Kesselman\*   Kirill Kogan<sup>†</sup>   Sergey Nemzer<sup>‡</sup>   Michael Segal<sup>§</sup>

August 10, 2010

## Abstract

Traffic classification in the Internet is a crucial mechanism necessary to support network services. Using Ternary Content Addressable Memories (TCAMs) to perform high-speed packet classification has become the de facto standard in industry. TCAMs concurrently match the packet headers against the rules in a classification database providing high throughput unparalleled by software-based solutions. The complexity of packet classification policies has been growing rapidly as the number of Internet services continues to increase. Many complex classification policies are naturally represented in a hierarchical fashion, where different layers perform classification based on the administrative domain and the traffic QoS parameters. However, multiple levels of classification hierarchy incur high lookup latency while high TCAM memory requirements of flattened classification policies is a major issue since TCAMs have very limited capacity. In this paper we focus on the fundamental trade-off between the TCAM space and the number of lookups in the TCAM classification policies. We consider two optimization problems of dual nature: the first problem is to minimize the number of TCAM entries subject to the constraint on the maximum number of levels in the policy hierarchy; the second problem is to minimize the number of levels in the policy hierarchy subject to the constraint on the maximum number of TCAM entries. We propose efficient algorithms for these problems, which do not require any hardware changes. To the best of our knowledge, this is the first work to study these problems. We also show experimental results that support our findings.<sup>1</sup>

---

\*Google Inc. Email: alx@google.com

<sup>†</sup>Cisco Systems, and Dept. of Com. Eng., Ben-Gurion Univ. E-mail: kkogan@cisco.com.

<sup>‡</sup>Compugen. E-mail: sergeyn@compugen.co.il.

<sup>§</sup>Dept. of Com. Eng., Ben-Gurion Univ. E-mail: segal@cse.bgu.ac.il.

<sup>1</sup>The preliminary version of this paper has been appeared in IEEE Sarnoff Symposium, 2008. The main changes versus the conference version are:

1. We revised the introduction to emphasize the importance of hierarchical packet classification and elaborated on the motivation of hierarchical policies.
2. We clarified the definition of classification architecture for multi-level hierarchies and the relationship between single-level and multi-level hierarchical algorithms.
3. We explicitly stated the assumptions related to the availability of header values at different stages of classification process.
4. We added a new section 7 with generalization of space optimization problem.
5. We added a new section 9 with experimental results.

The work on this paper has been partially supported by US Air Force European Office of Aerospace Research and Development, grant FA8655-09-1-3016, Deutsche Telecom, European project FLAVIA and Israeli Ministry of Industry, Trade and Labor (consortium CORNET)

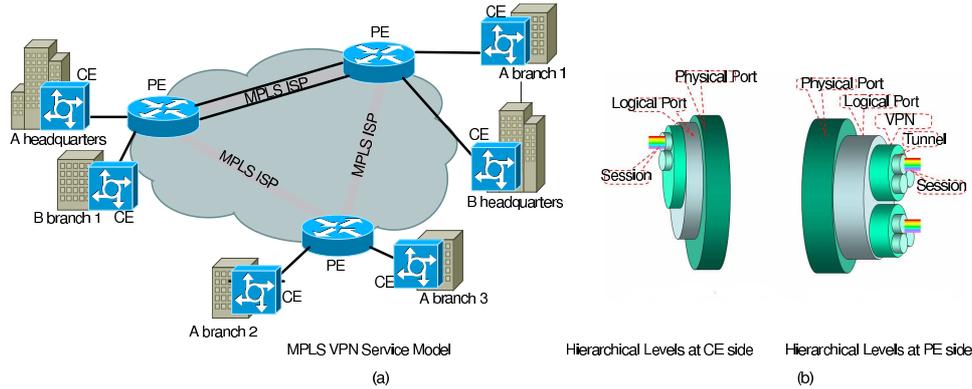


Figure 1: Application of Hierarchical QoS for MPLS VPN service.

## 1 Introduction

Packet classification is a crucial function for a variety of emerging network services such as Quality of Service (QoS), firewalling, and traffic monitoring and accounting, to name just a few.

An important characteristic influencing the network services is so called network hierarchy, that is the degree of concentration of traffic flows at interconnection points within the network. Hierarchies are important since they help to determine the sizes of networks, including the routing and addressing configurations as well as the scaling of network technologies, performance, and service levels. Furthermore, service-level agreements can include support for delay-constrained applications and may contain a variety of capacity and control mechanisms, such as traffic shaping and policing at multiple levels in the network.

Network devices maintain service policies under which incoming or outgoing packets are classified by matching against a hierarchical set of rules. For example, bandwidth can be shared between different users at the top level matching traffic according to the administrative affiliation while the second tier of the hierarchy can further classify the traffic according to the QoS class based on the Differentiated Services Code Point (DSCP), protocol type, or some other criteria. In addition, each rule can also specify a set of actions to be taken on packets matching this rule.

Several types of services such as *Multi Play*, *MPLS VPN*, *VPDN*, to name just a few, require implementation of hierarchical policies [31]. For instance, consider the following implementation of MPLS VPN service that is described in [31]. The *Customer Edge* (CE) side of the *Provider Edge* (PE) router links the branches of an enterprise, where each branch has a special bandwidth requirement. The enterprise has a limit on the outgoing traffic and the tunnels between PE routers bear traffic from different enterprises, which have different QoS requirements. The tunnels also have bandwidth limitations of their own. Figure 1 (a) depicts the MPLS VPN service model. Though for the CE side the traditional QoS can improve the service quality of each branch, it cannot manage the traffic of branches and that of the whole enterprise on PE. The problem is that for tunnels between PEs, the traditional QoS applies to either the tunnel or the internal traffic separately without being able relate between the two. On the other hand, the hierarchical QoS model can meet the QoS requirement of MPLS VPN service. Figure 1 (b) demonstrates a hierarchical structure of QoS policy that can be applied on CE and PE sides. Observe that on a PE side the hierarchy has five levels in depth. Therefore, supporting hierarchical service policies is a challenging task, which requires to perform hierarchical matching at the line rate.

Hierarchical classification polices are widely used in today's commercial routers [28, 29, 30, 31]. However,

most hierarchies have only a few levels due to the speed (line rate) constraints. In this work we study how to overcome this performance limitation. Specifically, we explore several interesting tradeoffs between the speed and the required space for wire-speed hierarchical packet classifiers that are implemented in *Ternary Content-Addressable Memory* (TCAM).

A TCAM is a memory device that stores data as a massive array of fixed-width ternary entries. A ternary entry is a string of bits where each bit is either 0, 1 or \* (“don’t care”). The TCAM searches the packet in parallel against all the ternary entries stored in the memory and produces the first rule that matches the packet. Remarkably, TCAM guarantees that each lookup is done in constant time. Usually each TCAM entry is wide enough to contain the concatenation of all the packet fields to be matched, possibly having room for some extra bits. If a matching rule consists solely of fields that specify exact or prefix matches, then it can be represented by a TCAM entry in a straightforward manner (a prefix match field is padded with the appropriate number of \*’s in the least significant bits). A range value may be converted to multiple prefixes or exact entries to fit the TCAM format.

## 2 Our Results

We study two fundamental problems dealing with hierarchical packet classification using TCAM. We are given a policy  $P$  with  $d$  levels of hierarchy and the goal is to convert  $P$  to an equivalent policy  $P'$  that needs to fulfill certain constraints minimizing the number of lookups or TCAM entries required. Our algorithms do not require any hardware modifications being very easy to deploy.

In the space optimization problem, we aim to minimize the number of TCAM entries required by  $P'$  subject to the limit on the maximum number of hierarchy levels. The motivation behind this problem is to reduce the required TCAM space allowing packet classification in wire-speed without incurring a possibly huge memory blowup if  $P$  is flattened to a single level. We propose two dynamic programming algorithms for this problem: the first algorithm is more efficient but is restricted only to policies that match disjoint fields in a packet header at different levels of the hierarchy; the second algorithm is slightly more complicated and can process general policies.

In the speed optimization problem, our goal is to minimize the number of hierarchy levels in  $P'$  subject to the limit on the maximum number of TCAM entries. The rationale behind this problem is to utilize the available TCAM capacity as efficiently as possible to reduce the number of lookups. We propose an algorithm that applies one of our speed optimization algorithms as a subroutine. This algorithm adds a factor of  $O(\log d)$  to the running time of the corresponding speed optimization algorithm. Please observe that the suggested algorithms are orthogonal to different approaches that minimize TCAM entries such as removing of redundant rules, conversion of ranges to prefixes, etc. (see SubSection 3) at a single hierarchical level. All our algorithms operate on already optimized classification rules at each hierarchical level and transparent to any specific TCAM HW architecture.

## 3 Related Work

Designing algorithms that scale to millions of rules and millions of searches per second has been and continues to be an important line of research. Many software-based sophisticated approaches have been proposed in the past few years including Recursive Flow Classification [10], Crossproducting [6, 22, 24], HyperCuts [21],

Extended Grid-of-Tries [2] and Aggregated Bit Vector [3], to name just a few. Comprehensive surveys on this subject can be found in [8, 12, 23, 25]. The complexity bounds derived by means of computational geometry imply that any software-based packet classifier with  $N$  rules and  $k > 2$  fields, uses either  $O(N^k)$  space and  $O(\log N)$  time or  $O(N)$  space and  $O(\log^{k-1} N)$  time [19]. Thus, many software-based approaches are either too slow or too memory intensive for  $k > 2$ . Though packet classification algorithms using decision trees achieve better time-space tradeoffs (see [11, 26]), they exploit statistical characteristics that are not reliable in general.

Due to the inherent limitations of software-based approaches, industry has increasingly employed hardware-based *Ternary Content Addressable Memory* (TCAM) for performing packet classification making it the dominant method [32, 33, 34]. A large class of packet classification systems that require up to a few hundred thousand rules have adopted TCAM for packet classification at multigigabit speeds [4, 9].

Some previous works consider TCAM space minimization for a packet classifier, for example [14, 18]. For the best of our knowledge all of them deal with a single level of hierarchy concentrating on elimination of overlapping rules and representation of filters ranges. Several schemes for converting ranges to TCAM rules have been proposed in [5, 13, 14, 15, 16, 17, 18]. Reducing of TCAM power consumption and increasing of throughput is explored in [1, 20, 27].

## 4 Paper Organization

The rest of the paper is organized as follows. The model description appears in Section 5. The algorithms for space and speed optimization are presented in Section 6 and Section 8, respectively. The generalization of space optimization problem is considered in Section 7. Finally, we conclude with Section 10.

## 5 Model Description

In this section we introduce the formal notation and define the Hierarchical Speed and Space Optimization problems.

### 5.1 Notation

A packet header contains  $k$  fields, where a field  $H_i$  ( $1 \leq i \leq k$ ) is a string of  $W_i$  bits. In an IPv4 packet, classifiers usually check the following six fields: the Type of Service (8 bits), the Destination Address (32 bits), the Source Address (32 bits), the Destination Port (16 bits), the Source Port (16 bits), and the Protocol Type (8 bits). Note that classifiers may access other fields besides TCP/IP header such as MAC or application headers. Packets are matched according to classification rules stored in a classification database.

The classification database of a router consists of a finite set of  $n$  rules,  $R_1 \dots R_n$ . Each rule  $R$  specifies matchings for one or more (up to  $k$ ) fields. For each header field  $H_i$ , a rule can specify a filter  $F_i$  of length  $|F_i|$  ( $|F_i| \leq |H_i|$ ), which can be any of two kinds of matches: exact match or prefix match.

1. A packet header field  $H_i$  exactly matches the filter  $F_i$  if and only if  $H_i = F_i$ .
2. A packet header field  $H_i$  is a prefix match for the filter  $F_i$  if and only if the  $|F_i|$  bits of  $H_i$  are equal to  $F_i$ .

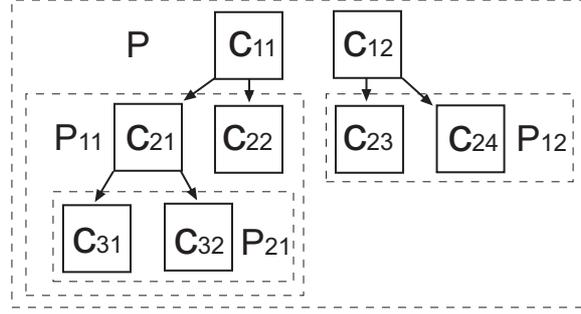


Figure 2: An example of policy and its graph  $G_p$ .

A packet  $p$  matches rule  $R$  if each of  $p$ 's header fields matches the corresponding filter of  $R$  if any. The header fields for which filters are not specified by the rule are matched in TCAM by a wildcard filter ("don't care"). Since a packet may match multiple rules, the classification problem is to determine the first matching rule in an ordered sequence of rules.

There is also a third type of matching, so-called range matching, where the header value should fall into a contiguous interval specified by the filter. In typical packet classifiers, such fields as the source and destination port numbers are represented as ranges rather than prefixes. Though range rules cannot be directly stored in TCAMs, they are usually converted to a corresponding set of prefixes each of which is stored in a separate TCAM entry. In this paper we deal with classification rules that reside in a TCAM device and assume that all filters are exact or prefix match (the classification database may undergo a conversion if necessary [5, 13, 14, 15, 16, 17, 18]). Such transformation of rules is orthogonal for our future discussion since our algorithms manipulate with already converted rules on each hierarchical level.

We define a *class*  $C$  to be an ordered set of rules and a set of *actions* to be taken on the packet. For instance, a QoS action may be packet marking with a pre-defined DSCP value while a security action may be packet accept or reject. We denote by  $|C|$  the number of rules or *cardinality* of  $C$ . We say that the class is matched if the first matching rule belongs to this class. A *policy*  $P$  is an ordered set of classes. The last class of a policy is usually so-called *default* class matching all packets that have not been matched by the other classes. Note that a global order of the rules is obtained by listing the rules in the corresponding classes.

The action of a class can also apply another policy in recursive manner, creating a *hierarchical policy*. Each recursive application creates a new *level of hierarchy*. A hierarchical policy can be viewed as a directed and acyclic graph  $G_P$  with classes acting as nodes and each edge representing a recursive policy application (see Figure 2). A directed path  $\mathcal{C} = C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_d$  from a *root* class with no incoming edges to a *terminal* class with no outgoing edges is called a *class chain*. The length of the longest class chain is defined as the *policy depth*. In order to match all rules on a class chain  $\mathcal{C}$ , a packet has to match rules  $R_1 \in C_1, R_2 \in C_2, \dots, R_d \in C_d$ . Observe that the total number of hierarchical class chains is equal to the number of terminal classes.

In the actual implementation of a hierarchical classifier, a special header field  $H_0$  is added to a packet at all levels of the hierarchy to identify the class sub-chain  $\mathcal{S} = C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_l$  ( $1 < l \leq d$ ) matched by the packet up to and including level  $l$ . In this way, the classifier can track the path of the packet in  $G_P$  and take the appropriate actions as matching proceeds. Furthermore, all the rules corresponding to the children of the class  $C_l$  in  $G_P$  have a common filter  $F_0$  that identifies the class sub-chain  $\mathcal{S}$ . As a result, different class chains in a hierarchical classifier are independent in the sense that distinct instances of the same class are created for

each node in  $G_P$ . Therefore, in order to minimize the number of entries across different hierarchical chains it is sufficient to consider separately each one of the chains. We assume that the size of  $F_0$  filter is at least as the cardinality of all hierarchical class chains in the policy graph. A value of  $H_0$  at each hierarchical level (from the second) is identified as result of the previous TCAM lookup. At the first level the value of  $H_0$  is known a priori and its value is identical for all hierarchical chains of the same policy.

Hierarchical policies allow a high degree of flexibility and modularity in policy definition. However, a separate TCAM lookup needs to be performed for each level of the hierarchy in the process of classification, which may incur large delays for policies with high depth. To speed up the classification process, a hierarchical policy can be converted to an equivalent policy of lower depth. At the same time, such a conversion can significantly increase the number of TCAM entries required to store the merged rules. We define the *TCAM space*  $M$  of a policy as the total number of rules in all levels of the classification hierarchy.

The operations defined below deal with policy flattening. We will describe how to *intersect* two filters  $F_i$  and  $F'_i$  specified for a common header field  $H_i$ . Suppose without loss of generality that  $|F_i| \leq |F'_i|$ . If  $F_i$  is a prefix of  $F'_i$ , then a resulted filter  $F_i \cap F'_i$  will consist of filter  $F'_i$ . Otherwise, a resulted filter  $F_i \cap F'_i$  is an *empty filter* that does not match any value of header field  $H_i$ . The merge of rules  $R \otimes R'$  is defined as a set of intersections of the corresponding filters in  $R$  and  $R'$ . If at least one empty filter is produced during this operation then the merged rule  $R \otimes R'$  is called an *empty rule* that does not match any packet. We define the *merge* of two classes  $C \otimes C'$  as a class consisting of all possible merges of rules  $R \otimes R'$  that are not empty, where  $R \in C$  and  $R' \in C'$  that sequentially applies actions of  $C$  and  $C'$ . For a class chain  $\mathcal{C} = C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_d$  we define a *virtual class*  $C_{1,d}$  as  $C_1 \otimes C_2 \otimes \dots \otimes C_d$  that represents the merge of all classes in  $\mathcal{C}$ .

Observe that  $|C_1 \otimes C_2 \dots C_n| \leq |C_1| |C_2| \dots |C_n|$ . Let  $C_1$  consist of three rules  $[ToS = 1]$ ,  $[ToS = 2]$  and  $[ToS = 3]$  and  $C_2$  consist of two rules  $[DstPort = 21]$  and  $[DstPort = 80]$ . In this case,  $|C_1 \otimes C_2| = 6$  as the filters are specified for disjoint header fields and  $C_1 \otimes C_2$  contains the following rules:  $[ToS = 1, DstPort = 21]$ ,  $[ToS = 2, DstPort = 21]$ ,  $[ToS = 3, DstPort = 21]$ ,  $[ToS = 1, DstPort = 80]$ ,  $[ToS = 2, DstPort = 80]$ ,  $[ToS = 3, DstPort = 80]$ . At the same time, the TCAM space required to represent the merge of classes may be smaller than the cardinality of the classes themselves if filters are specified for common header fields. For instance, suppose that  $C_1$  contains three rules  $[ToS = 1]$ ,  $[ToS = 2]$  and  $[ToS = 3]$  and  $C_2$  contains two rules  $[ToS = 3]$  and  $[ToS = 4]$ . We obtain that  $|C_1 \otimes C_2| = 1$  and  $C_1 \otimes C_2$  includes just one rule  $[ToS = 3]$ .

## 5.2 Problem Statement

We say that two packet classifiers are (semantically) equivalent if and only if they apply the same actions on each packet. Next, we define the optimization problems studied in this paper.

**Hierarchical Space Optimization Problem:** Given a hierarchical policy  $P$  with depth  $d > 1$ , the goal is to convert  $P$  to an equivalent policy  $P'$  with depth of at most  $l$  ( $l < d$ ) that minimizes the required TCAM space.

**Hierarchical Speed Optimization Problem:** Given a hierarchical policy  $P$  that requires TCAM space of  $M$ , the aim is to convert  $P$  to an equivalent policy  $P'$  that minimizes the policy depth subject to the constraint that the TCAM space cannot exceed the available TCAM space  $A$  ( $A > M$ ).

For now we assume that a set of classified packet headers is known a priori.

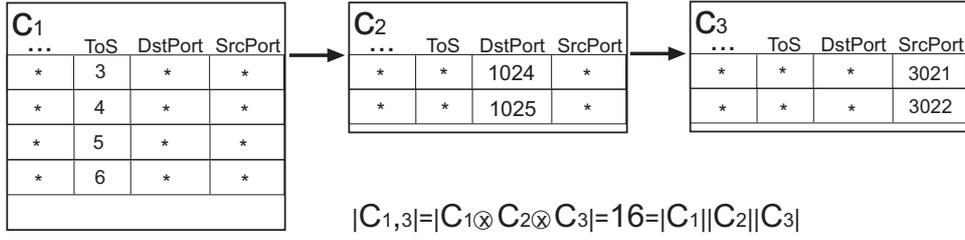


Figure 3: Example of well-structured hierarchy.

## 6 Space Optimization

In this section we consider the case where the values of classified packet headers are known a priori and are not changed during classification process. We consider the problem of minimizing the TCAM space subject to the constraint on the policy depth. First, we present an algorithm for the case where rules in any class chain apply only to disjoint header fields. Then we propose an algorithm for the general case where we impose no restrictions on the policy whatsoever.

### 6.1 Well-Structured Hierarchies

In this section we study *well-structured* hierarchical policies in which rules in any class chain apply only to disjoint header fields. The example of well-structured hierarchy is presented on Figure 3. We can learn from the same figure the importance of decision which sub-chains should be merged. At one hand, if we merge the classes  $C_2$  and  $C_3$  then the cardinality of flattened hierarchical chain is  $|C_1| + |C_2 \otimes C_3| = 4 + 4 = 8$ . From the other hand, if we merge the classes  $C_1$  and  $C_2$  then the cardinality of the flattened chain is  $|C_1 \otimes C_2| + |C_3| = 8 + 2 = 10$ .

Well-structured hierarchies have the following important property, which allows us to use a fast dynamic programming algorithm operating merely with cardinalities of merged classes without actually merging the rules themselves until the final stage.

*Input: policy  $P$  of depth  $d$ , integer  $l$  ( $l < d$ )*

*Output: policy  $P'$  equivalent to  $P$  of depth  $l$*

- **Step 1: Merging Long Chains.** For each class chain  $\mathcal{C} = C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_d$  in  $P$  such that  $d > l$  create a virtual class  $C_{1,d} = C_1 \otimes C_2 \otimes \dots \otimes C_d$  that represents the merge of all classes in  $\mathcal{C}$  (without merging the rules).
- **Step 2: Splitting Merged Chains.** Split each merged virtual class  $l - 1$  times using the Level Splitting Algorithm (see Figure 6).
- **Step 3: Merging Rules.** For each virtual class, merge the rules of the corresponding merged classes.
- **Step 4: Creating Output Policy.** Output  $P'$  as a union of all the original chains of length at most  $l$  and the converted long chains.

Figure 4: Space optimization algorithm (*SOAW*) for well-structured hierarchies.

**Observation 1** *In a well-structured hierarchy, for any class chain  $\mathcal{C} = C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_d$  the cardinality of the classes' merge in  $\mathcal{C}$ , that is  $|C_1 \otimes C_2 \otimes \dots \otimes C_d|$ , equals the product of the cardinalities of the individual*

classes  $|C_1||C_2| \dots |C_d|$ .

The space optimization algorithm for well-structured hierarchies (*SOAW*) is presented on Figure 4. *SOAW* proceeds by first merging all chains of length greater than  $l$  into a single *virtual* class without actually merging the rules. Then each virtual merged class is splitted  $l - 1$  times in an optimal way using a level splitting algorithm based on the dynamic programming technique. Finally, for each virtual class the rules of the corresponding merged classes are merged to produce the output policy. Observe that class chains of length smaller than  $l$  are left untouched by *SOAW*. The running time of *SOAW* is output sensitive. Basically, it build the optimized policy according to the output of level splitting algorithm, which is applied to all long chains in  $P$ . Recall that the total number of class chains equals the number of terminal classes.

Now we describe how to divide the level splitting problem into two sub-problems and combine solutions to these sub-problems into a solution to the original problem. For an merged class  $C_{1,d} = C_1 \otimes C_2 \otimes \dots \otimes C_d$ , we denote by  $V(i, j, n)$  ( $1 \leq i < d, i < j \leq d$ ) the cost of an optimal solution for the problem of  $n$  splittings in the merged sub-class  $C_{i,j} = C_i \otimes C_{i+1} \otimes \dots \otimes C_j$ . The cost is measured as the cumulative cardinality of the resulting  $n + 1$  sub-classes.

We define initial values for the special case of  $n = 0$  where nothing needs to be done and the special case of  $n > j - i$  where no feasible solution exists:

$$V(i, j, n) = \begin{cases} |C_{i,j}| & : n = 0, \\ \infty & : n > j - i. \end{cases} \quad (1)$$

The main recurrence relation is defined as follows for  $n > 0$  and  $n \leq j - i$ :

$$V(i, j, n) = \min_u (V(i, u, 0) + V(u + 1, j, n - 1)) \quad (2)$$

for  $i \leq u \leq j - n$ .

Basically, in order to make  $n$  level splittings we consider all possibilities for the first splitting and perform exhaustive search over all the remained at most  $n - 1$  splittings in  $C_{u+1,j}$  sub-class. The Figure 6.1 demonstrates this process.

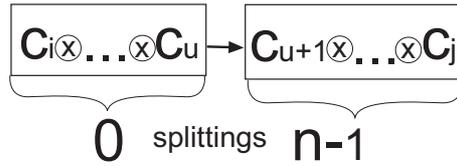


Figure 5: Dynamic programming for well-structured policy.

Our aim is to minimize the overall cost of the produced solution. The level splitting algorithm (*LSA*) appears on Figure 6.

It is easy to see that the space complexity of *LSA* is  $O(d^2l)$  and the running time is  $O(\max(dl^2, d^2l))$ , which is very reasonable since  $d$  and  $l$  are typically small numbers. The next theorem shows the correctness of *LSA*.

**Theorem 6.1** *The level splitting algorithm (LSA) finds an optimal solution of minimum cost for the problem of  $n$  splittings ( $n > 1$ ) in a class formed as the merge of  $d$  classes ( $d > n$ ).*

*Input: hierarchical chain  $\mathcal{C}$  consisting of  $d$  classes, integer  $l$  ( $l < d$ )*

*Output: chain  $\mathcal{C}'$  of length  $l$*

- **Step 1: Initialization.** Initialize  $V(i, j, n)$  for  $n = 0$  and  $n > j - i$  according to Equation 1.
- **Step 2: Calculation.** Calculate all values  $V(i, j, n)$  starting from  $n = 1$  up to  $n = l - 1$  for  $n \leq j - i$  using the recurrence Equation 2 and record the splitting of minimum cost.
- **Step 3: Reconstructing the optimal solution.** Construct the chain  $\mathcal{C}'$  by splitting the classes of  $\mathcal{C}$  with respect to the optimal solution of minimum cost  $V(1, d, l - 1)$ .

Figure 6: Level splitting algorithm (*LSA*) for well-structured hierarchies.

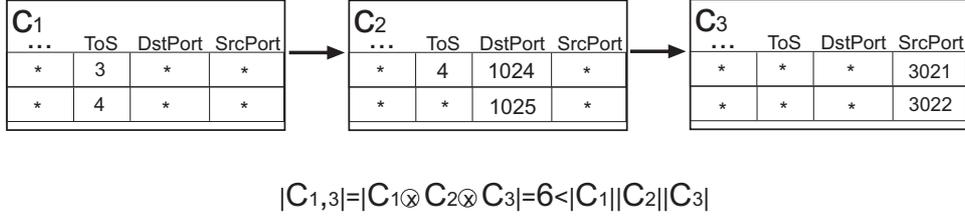


Figure 7: Example of arbitrary hierarchy.

**Proof:** The proof is by induction on the number of splittings  $n$ . Clearly, *LSA* finds an optimal solution for  $n = 0$  since it just returns the original input as no splittings are necessary.

Suppose that *LSA* finds an optimal solution for the number of splittings of at most  $n$  and let us show that it also finds an optimal solution for  $n + 1$  splittings. Note that *LSA* considers all options for making the first splitting, one of which necessarily corresponds to an optimal solution. Having done an optimal first splitting, it must be the case that the rest of this optimal solution consists of an independent optimal solution for the right sub-class with at most  $n$  splittings. By the induction hypothesis, *LSA* finds an optimal solution for the right sub-class with at most  $n$  splittings. Therefore, *LSA* finds an optimal solution for at most  $n + 1$  splittings. ■

## 6.2 Arbitrary Hierarchies

In this section we deal with the general case of arbitrary hierarchies. Unfortunately, such hierarchies require to calculate the actual merge of the rules for the classes' merge in order to obtain its cardinality as demonstrated by the next observation.

**Observation 2** *In an arbitrary hierarchy, for any class chain  $\mathcal{C} = \mathcal{C}_1 \rightarrow \mathcal{C}_2 \rightarrow \dots \rightarrow \mathcal{C}_d$  we have that the cardinality of the merge of all classes in  $\mathcal{C}$ , that is  $|\mathcal{C}_1 \otimes \mathcal{C}_2 \otimes \dots \otimes \mathcal{C}_d|$ , is bounded from above but may not be equal to the product of the cardinalities of the individual classes  $|\mathcal{C}_1| |\mathcal{C}_2| \dots |\mathcal{C}_d|$ .*

The example of arbitrary hierarchy is shown on Figure 7.

That necessitates a slightly more complicated dynamic programming algorithm for level merging compared to the level splitting algorithm used for well-structured hierarchies. The space optimization algorithm for general hierarchies (*SOAG*) appears on Figure 8. *SOAG* merges all chains of length greater than  $l$  by running a level merging algorithm based on the dynamic programming technique. The main component of *SOAG* is the level merging algorithm, which processes all long chains in  $P$ .

*Input: policy  $P$  of depth  $d$ , integer  $l$  ( $l < d$ )*

*Output: policy  $P'$  equivalent to  $P$  of depth  $l$*

- **Step 1: Merging Long Chains.** Merge all chains of length  $d'$  greater than  $l$  using the level merging algorithm with the number of merges  $m = d' - l$ .
- **Step 2: Creating Output Policy.** Output  $P'$  as a union of all the original chains of length at most  $l$  and the converted long chains.

Figure 8: Space optimization algorithm (*SOAG*) for general hierarchies.

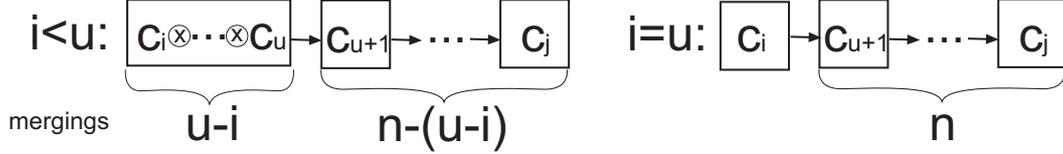


Figure 9: Dynamic programming for a chain with arbitrary structure.

In what follows we present a way of dividing the level merging problem into two sub-problems and combining solutions to these sub-problems into a solution to the original problem. For a class chain  $\mathcal{C} = C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_d$ , define  $V(i, j, n)$  ( $1 \leq i < d, i < j \leq d$ ) as the cost of an optimal solution for the problem of merging  $n$  levels in the class sub-chain  $\mathcal{C}_{i,j} = C_i \rightarrow C_{i+1} \rightarrow \dots \rightarrow C_j$ . We estimate the cost of a solution as the total cardinality of the produced  $j - i - n + 1$  sub-classes.

Initial values are set for the special case of  $n = 0$  where no merges have to be performed, the special case of  $n > j - i$  which permits no feasible solution and the case of  $n = j - i$  where all levels are completely merged. Note that there is no need to merge more than  $d - l$  levels to obtain the desired solution. We define initial values for the special case of  $n = 0$  where nothing needs to be done and the special case of  $n > j - i$  where no feasible solution exists:

$$V(i, j, n) = \begin{cases} \sum_{u=i}^j |C_u| & : n = 0, \infty : n > j - i, \\ |C_{i,j}| = |C_i \otimes \dots \otimes C_j| & : n = j - i, n \leq d - l. \end{cases} \quad (3)$$

We specify the main recurrence relation for  $n > 0$  and  $j - i > n$  in the following way:

$$V(i, j, n) = \min_u (V(i, u, u - i) + V(u + 1, j, n - (u - i))) \quad (4)$$

for  $i \leq u \leq i + n$ .

Essentially, we cover all options for a leftmost merged class sub-chain  $\mathcal{C}_{i,u}$  with  $u - i$  mergings. Then we consider all possible sub-divisions of the remaining  $n - (u - i)$  mergings of the class sub-chain  $\mathcal{C}_{u+1,j}$ .

The goal is to minimize the total cost of the resulting solution. The level merging algorithm (*LMA*) can be found on Figure 10.

We obtain that the space complexity of *LMA* is  $O(md^2)$  and the running time is  $O(\max(md^2, m^2d))$ . The subsequent theorem demonstrates the correctness of *LMA*.

**Theorem 6.2** *The level merging algorithm (LMA) finds an optimal solution of minimum cost for the problem of merging ( $n > 1$ ) levels in a chain of length  $d$  ( $d > n$ ).*

*Input: chain  $C$  (of length  $d$ ), integer  $m$  ( $m < d$ )*

*Output: chain  $C'$  of length  $d - m$*

- **Step 1: Initialization.** Initialize  $V(i, j, n)$  for  $n = 0$ ,  $n > j - i$ , and  $n = j - i$  ( $n \leq m$ ) according to Equation 3.
- **Step 2: Calculation.** Calculate all values  $V(i, j, n)$  starting from  $n = 1$  up to  $n = m$  for  $n > j - i$  using the recurrence Equation 4 and record the merging of minimum cost at each stage.
- **Step 3: Reconstructing the optimal solution.** Construct the chain  $C'$  by merging the classes of  $C'$  with respect to the optimal solution of minimum cost  $V(1, d, m)$ .

Figure 10: Level merging algorithm (*LMA*) for general hierarchies.

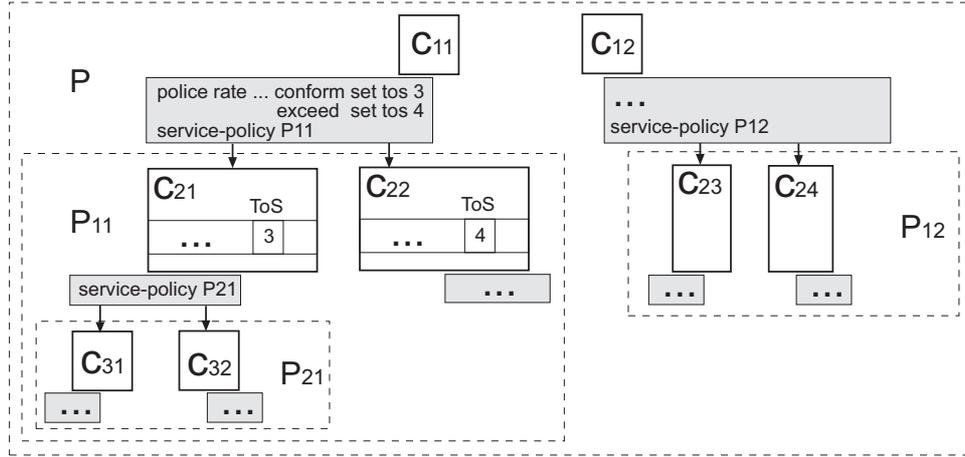


Figure 11: Motivation for generalized space optimization problem.

**Proof:** The proof is by induction on the number of mergings  $n$ . Obviously, *LMA* finds an optimal solution for  $n = 0$  since it just returns the original input as no mergings need to be done. Suppose that *LMA* finds an optimal solution for the number of mergings of at most  $n$  and let us prove that it finds an optimal solution for  $n + 1$  mergings as well. Observe that *LMA* examines all possibilities for making  $z$  ( $1 \leq z \leq n + 1$ ) mergings in the leftmost class sub-chain. It must be the case that one of these mergings is a merging in an optimal solution. Having done the first merging that is a part of optimal solution, the other mergings in this optimal solution is independent optimal solution for the right class sub-chain under at most  $n + 1 - z$  mergings. If the first merging is an empty merging (i.e.  $z = 0$ ), then the above argument is applied to the right class sub-chain. According to the induction hypothesis, *LMA* finds an optimal solution in the right class sub-chain with at most  $n + 1 - z \leq n$  mergings for the first non-empty merging (i.e.  $z > 0$ ). Hence, *LMA* finds an optimal solution for at most  $n + 1$  mergings. ■

## 7 Generalization of Space Optimization Problem

In this section we discuss a generalization of the space optimization problem in which we obtain an additional constraint that the inner protocol header fields cannot be extracted until the outer protocol header fields have been parsed. The rationale behind this problem is that certain levels of the hierarchy cannot be "merged" during optimization of a hierarchical chain as classifying complex protocols based on the analysis of encapsulated

content often requires sequential parsing of nested protocol headers. For instance, the application header cannot be extracted directly if the TCP header contains variable length options. Moreover, packet header fields can be dynamically changed during classification process as a result of class actions.

A *barrier* is the level of hierarchical chain where two neighboring classes belonging to it could not be “merged”. The classifier has to issue a separate TCAM lookup for each barrier in hierarchical chain. For instance, on Figure 11 we consider a packet that is matched to the class  $C_{11}$ . For correctness of classification process we have to update TOS value to 3 or 4 according to police action of  $C_{11}$  that defines a barrier in all hierarchical chains initiated from  $C_{11}$ . That means  $C_{11}$  could not be merged with  $C_{21}$  during optimization of hierarchical chain  $\mathcal{C} = C_{11} \rightarrow \dots \rightarrow C_{31}$ .

We will show that the generalized version of space optimization algorithms is simply extendable. The definition of *virtual sub-class*  $C_{i,j}$  previously being defined as  $C_i \otimes C_{i+1} \otimes \dots \otimes C_j$  now should take into account permanently located barriers. More precisely, if there is a barrier at the  $k$ -th position in an hierarchical sub-chain  $C_{i,j}$ ,  $i \leq k < j$  then the virtual sub-class  $C_{i,j}$  is defined as  $C_{i,k} \rightarrow C_{k+1,j}$  and its cardinality  $|C_{i,j}|$  is equal to  $|C_{i,k}| + |C_{k+1,j}|$ ,  $i < k < j$ . We assume that the number of barriers  $b$  in any class chain  $\mathcal{C}$  is less than  $|\mathcal{C}|$ ; otherwise, no feasible solution exists. The generalized version of *LSA* algorithm for the case of well-structured chain will find optimal placement of  $l - b - 1$  splittings. Moreover, the modified version of *LSA* will not try to place splittings at the predefined barrier positions. The generalized version of *LMA* algorithm for the case of general hierarchical chains can just skip during dynamical programming levels corresponding to barriers.

## 8 Speed Optimization

In this section we study the problem of minimizing the number of levels in the policy hierarchy subject to the constraint on the maximum TCAM space. We utilize the space optimization algorithms from the previous section. The speed optimization algorithm is presented on Figure 12.

*Input: policy  $P$  with TCAM space  $M$  of depth  $d$ , TCAM space  $A$  ( $A > M$ )*

*Output: policy  $P'$  equivalent to  $P$  with TCAM space of at most  $A$*

Perform *binary search* on the policy depth between 1 and  $d$  by applying either *SOAW* or *SOAG* on  $P$  depending on the structure of  $P$ 's hierarchy and find the minimum depth  $l$  for which the TCAM space of the produced policy  $P'$  does not exceed  $A$ .

Output  $P'$ .

Figure 12: Speed optimization algorithm.

Generally speaking, we need to find the optimal value  $l$  of the policy depth and then optimize the TCAM space of the transformed policy  $P'$  for depth  $l$ . The binary search is performed because this value is not known in advance. Once we have found the optimal depth, the space optimization algorithm guarantees minimization of the TCAM space. The running time of the speed optimization algorithm is  $O(\log d)$  times the running time of *SOAW* or *SOAG*, that is at most  $O(d^3)$  times the number of terminal classes in  $P$ , respectively. Remarkably, the running time of the speed optimization algorithm does not depend on  $A$ , which can be by orders of magnitude larger than  $d$ .

For the case of arbitrary structures the generalized version of *LMA* will not try to place mergings at the

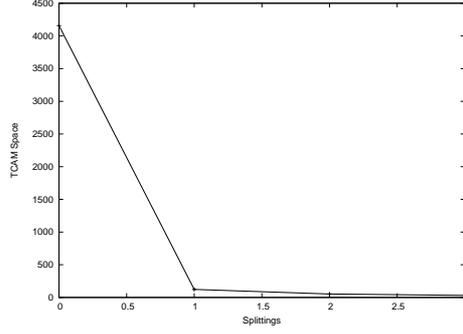


Figure 13: Required TCAM space as a function of a number of splittings.

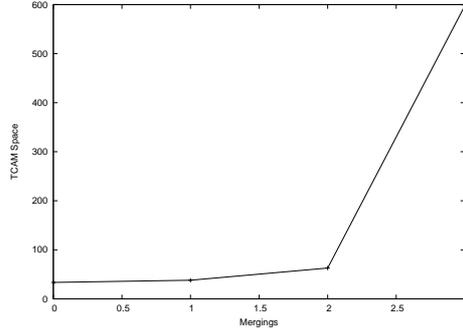


Figure 14: Required TCAM space as a function of a number of mergings.

predefined barriers positions. We obtain that the time and space complexities of generalized versions of *LSA* and *LMA* are the same as original.

## 9 Experimental Results

In this section, we evaluate the efficiency of the proposed algorithms, which use additional TCAM lookups to minimize the TCAM space required to represent a hierarchical policy. Since hierarchical chains are independent, we evaluate the algorithms for merging/splitting of a single hierarchical chain rather than the whole classification policy. Unfortunately, there is no available for evaluation real-world data for classifiers with more than two levels of hierarchy. Hence, we synthesize data for our experimental study.

In our simulations we generate hierarchical chains with a length of 4. The cardinality of any class belonging to such a chain is at most 16 being chosen uniformly at random. In addition, for the case of general hierarchical chains that are used for simulation of *LMA* algorithm, the intersection between any two neighboring classes is also chosen uniformly at random. For each simulated algorithm we perform over 1000 trials.

Firstly, we consider *LSA* algorithm. Figure 13 presents the dependency of TCAM space on the number of splittings for the case of well-structured class chains. The simulation results demonstrate very fast exponential growth of the required TCAM space when the number of splittings in a hierarchical chain decreases.

Next, we evaluate the *LMA* algorithm for the case of general class chains. Figure 14 shows the dependency of TCAM space on the number of mergings. As in the previous case, the simulation demonstrates an exponential

growth of the required TCAM space with increasing of the number of mergings. Note that as the cardinality of the intersections between classes of the general hierarchical chain is growing, the TCAM space requirement drops down.

As evidenced by the performed experiments for synthetically created classifiers, there exists a fundamental tradeoff between the TCAM classifier space requirements and the lookup speed, which is crucial for efficient implementation of hierarchical classification. Specifically, we need to find the smallest memory footprint that still satisfies the latency requirement so as to optimize the memory cost.

## 10 Conclusion

Hierarchical packet classification is a key operation needed in provisioning of many crucial network services. One of the major challenges in design of the next generation high-speed switches is to deliver wire-speed packet classification. TCAMs are the dominant industry standard used for multi-gigabit classifiers. However, as packet classification policies grow in depth and complexity, there arises a fundamental tradeoff between the TCAM space and the number of lookups for hierarchical policies.

In this paper we propose novel algorithms based on dynamic programming for solving two important problems concerned with hierarchical packet classification. The algorithms for the first problem minimize the TCAM space given a constraint on the policy depth while the algorithm for the second problem minimizes the policy depth subject to the constraint on the maximum TCAM space. Also we study extensions of space optimization problem. Our algorithms do not require any modification to existing packet classification systems and can be easily deployed. Exploring tradeoff between required TCAM space and performance states a good balance for the future efficient implementation of hierarchical classifications that do not require any hardware changes. As far as we aware, this is the first work to study TCAM speed and space optimization for hierarchical packet classification. We believe that studying of the proposed tradeoffs and interconnection between them is interesting and provides additional insight on generalized classification problem with several hierarchical levels.

## References

- [1] B. Agrawal, T. Sherwood, "Ternary CAM power and delay model: extensions and uses", *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, v.16 n.5, pp. 554–564, May 2008
- [2] F. Baboescu, S. Singh, and G. Varghese, "Packet Classification for Core Routers: Is there an Alternative to CAMs?", *Proc. of IEEE INFOCOM*, 2003.
- [3] F. Baboescu and G. Varghese, "Scalable Packet Classification", *Proc. of ACM SIGCOMM*, 2001.
- [4] J. Bolaria and L. Gwennap, "A guide to search engines and networking memory", [http://www.linleygroup.com/reports/memory\\_guide.html](http://www.linleygroup.com/reports/memory_guide.html).
- [5] R. Cohen and D. Raz, "Simple Efficient TCAM Based Range Classification" *Proc. of IEEE INFOCOM*, 2010.

- [6] Sarang Dharmapurikar, Haoyu Song, Jonathan S. Turner, John W. Lockwood, "Fast packet classification using bloom filters", *ANCS 2006*, pp. 61–70.
- S. Dharmapurikar, H. Song, J. Turner, J. Lockwood, "Fast packet classification using bloom filters", *Proceedings of the 2006 ACM/IEEE symposium on Architecture for networking and communications systems*
- [7] Yeim-Kuan Chang, "A 2-Level TCAM Architecture for Ranges", *IEEE Trans. Computers* 55(12), pp. 1614–1629, 2006
- [8] H. J. Chao, "Next Generation Routers", *Proceedings of the IEEE*, Vol. 90, No. 9, pp. 1518–1558, 2002.
- [9] P. Gupta, K. Etzel and J. Bolaria, "A Scalable and Cost-Optimized Search Subsystem for IPv4 and IPv6", <http://www.eetimes.com/netseminar.html>.
- [10] P. Gupta and N. McKeown, "Packet Classification on Multiple Fields", *Proc. of ACM SIGCOMM*, pp. 147–160, 1999.
- [11] P. Gupta and N. McKeown, "Packet classification using hierarchical intelligent cuttings", *Proc. of Hot Interconnects VII*, Aug. 1999.
- [12] P. Gupta and N. McKeown, "Algorithms for Packet Classification", *IEEE Network*, pp. 24–32, March/April 2001.
- [13] A. Bremler-Barr, D. Hay, D. Hendler, B. Farber, "Layered interval codes for tcam-based classification", *ACM SIGMETRICS Performance Evaluation Review*, v.36 n.1, June 2008
- [14] A. Bremler-Barr, D. Hendler, "Space-Efficient TCAM-based Classification Using Gray Coding", *INFOCOM 2007*, pp. 1388–1396.
- [15] K. Lakshminarayanan, A. Rangarajan, and S. Venkatachary, "Algorithms for advanced packet classification with ternary cams", *Proc. of ACM SIGCOMM*, pp. 193–204, Aug. 2005.
- [16] H. Liu, "Efficient mapping of range classifier into ternary-cam". *Proc. of the Hot Interconnects*, pp. 95–100, 2002.
- [17] J. van Lunteren and T. Engbersen, "Fast and scalable packet classification", *IEEE Journals on Selected Areas in Communications*, Vol. 21, No. 4, pp. 560–571, 2003.
- [18] Chad R. Meiners, Alex X. Liu, and Eric Torng, "TCAM Razor: A Systematic Approach Towards Minimizing Packet Classifiers in TCAMs", *Proc. of the 15th IEEE International Conference on Network Protocols (ICNP)*, pp 226–275, 2007.
- [19] M. H. Overmars and A. F. van der Stappen, "Range searching and point location among fat objects", *Journal of Algorithms*, Vol. 21, No. 3, pp. 629–656, 1996.
- [20] R. Panigrahy and S. Sharma, "Reducing TCAM Power Consumption and Increasing Throughput", *Proc. of HOT Interconnects*, pp. 107–112, 2002.
- [21] S. Singh, F. Baboescu, G. Varghese, and J. Wang, "Classification using Multidimensional Cutting", *Proc. of ACM SIGCOMM*, 2003.

- [22] V. Srinivasan, G. Varghese, S. Suri, and M. Waldvogel, "Fast and Scalable Layer Four Switching", *Proc. of ACM SIGCOMM*, 1998.
- [23] D.E. Taylor, "Survey and taxonomy of packet classification techniques", *ACM Computer Surverys*, pp. 238-275, 2005.
- [24] D. E. Taylor and J. Turner, "Scalable Packet Classification using Distributed Crossproducting of Field Labels", *Technical Report WUCSE-2004-38 Washington Univ., St. Louis*, 2004.
- [25] G. Varghese, "Network Algorithmics: An Interdisciplinary Approach to Designing Fast Networked Devices", *Morgan Kaufmann Publishers Inc.*, 2004.
- [26] T. Y. C. Woo, "A modular approach to packet classification: Algorithms and results", *Proc. of IEEE INFOCOM*, pp. 1213-1222, 2000.
- [27] W. Wu , Jian Shi , Ling Zuo , Bingxin Shi, "Power-Efficient TCAMS for Bursty Access Patterns", *IEEE Micro*, v.25 n.4, pp.64-72, July 2005.
- [28] "Cisco ASR 1000 Series Aggregation Services Routers",  
<http://www.cisco.com/en/US/products/ps9343/index.html>.
- [29] "Modular QoS CLI Three-Level Hierarchical Policer",  
<http://www.cisco.com/en/US/docs/ios/12.2t/12.2t13/feature/guide/ft3level.html>.
- [30] "Juniper Networks E320 Broadband Services Router",  
[https://www.apac-juniper.net/juniper\\_public/campaign/pdf/iptv/en/100103.pdf](https://www.apac-juniper.net/juniper_public/campaign/pdf/iptv/en/100103.pdf).
- [31] "Technical White Paper for MSCG Hierarchical - QoS",  
[www.huawei.com/products/datacomm/pdf/view.do?f=916](http://www.huawei.com/products/datacomm/pdf/view.do?f=916).
- [32] "Cypress Semiconductor Corp. Content addressable memory",  
<http://www.cypress.com/>
- [33] "Integrated Device Technology, Inc. Content addressable memory",  
<http://www.idt.com/>
- [34] "Netlogic Microsystems. Content addressable memory",  
<http://www.netlogicmicro.com>