

IBM Research Report

A Programmable Message Classification Engine for Session Initiation Protocol (SIP)

Arup Acharya, Xiping Wang, Charles P. Wright
IBM Research Division
Thomas J. Watson Research Center
P.O. Box 704
Yorktown Heights, NY 10598



Research Division
Almaden - Austin - Beijing - Haifa - India - T. J. Watson - Tokyo - Zurich

A Programmable Message Classification Engine for Session Initiation Protocol (SIP)

Arup Acharya, Xiping Wang, and Charles P. Wright

IBM TJ Watson Research Center, Hawthorne, NY

{arup, xiping, cpwright}@us.ibm.com

Abstract—Session Initiation Protocol (SIP) has begun to be widely deployed for multiple services such as VoIP, Instant Messaging and Presence. Each of these services uses different subsets of SIP messages, and depending on the value of a service, e.g. revenue, the associated messages may need to be prioritized accordingly. Even within the same service, different messages may be assigned different priorities. In this paper, we present the design and implementation of a *programmable* classification engine for SIP messages in the Linux kernel. This design uses a novel algorithm that in addition to classifying messages can extract and maintain state information across multiple messages. We apply the classifier for overload control using operator-specified rules for categorizing messages and associated actions, augmented with a protocol-level understanding of SIP message structure. When faced with loads beyond their capacity (e.g., during catastrophic situations and major network outages), SIP servers must drop messages from the input stream. It is therefore desirable that the server process high-value messages in preference to dropping lower-value messages. We evaluated our in-kernel classifier implementation with a commonly-used open source SIP server (SER) for such an overload scenario. The workload consists of a mix of call setup and call handoff SIP messages and the classifier is programmed with rules that prioritize handoffs over call setups (reflecting typical message prioritization used by mobile service providers). We show that, while SER can process about 40K messages/sec (in a FIFO manner), our classifier can examine and prioritize about 105K messages/sec during overload. With the classifier operating at peak throughput, SER's processing rate drops to about 31.6 K messages/sec, but it should be noted that the processed messages reflect as much of the high-value messages as available in the input stream.

I. INTRODUCTION

Session Initiation Protocol (SIP) is a control plane for establishing, manipulating, and terminating multimedia sessions with one or more participants. SIP is media agnostic and can establish voice, text, video, and other types of sessions. SIP has gained widespread acceptance and deployment already among wire line service providers for introducing new services such as VoIP, within the enterprise for Instant Messaging and collaboration, and for push-to-talk service amongst mobile carriers. Service

providers ranging from cable companies to mobile providers are looking to deploy IP Multimedia Subsystem (IMS) as a common platform for deploying new services and applications [8]. SIP is an extensible protocol and extensions have been proposed for new functionality such as presence [4].

In this study, we study the problem of how to classify SIP messages before they are processed by a SIP server. Today, SIP servers process messages in a first-in-first-out manner. This does not lend itself to prioritizing messages before they are processed at the server. We design an efficient algorithm that takes as input a set of user-defined rules, and morphs them into suitable data structures that enable fast matching of rules against the input message stream. The rules specify both how to identify specific subsets of messages, based on a combination of message header values including complex functions such as set membership as well operations on the state amassed at the classifier from previous messages, and the actions to be executed on the matching packets. Since the classifier is driven by these user-specified rules, it can be programmed to suit specific goals ranging from overload control to denial-of-service prevention for SIP servers.

We showcase overload control as a defining example for our classifier in this paper. Given the variety of usage contexts for a SIP server (e.g. Voice-over-IP, Instant Messaging, Presence,..), it is not surprising that each service provides a different value to the operator (e.g., revenue or customer satisfaction). Moreover, different types of messages within a service can also provide different amounts of value (e.g., “411” information calls cost \$1 on some mobile services). Thus, our classifier-based solution for overload control will aim to maximize the value of the messages processed by the server. SIP servers can become overloaded despite being provisioned correctly. During overload, only some requests can be handled and the rest are dropped to decrease the server load and bring the load down to maximum server capacity. Rather than dropping requests randomly or in a FIFO fashion, our goal is to prioritize requests in order to maximize value for an operator. Additionally, each server operator may have a different notion of value attached to a specific type of request. We demonstrate how our classifier prioritizes messages according to operator-specified metrics, so that under overload conditions, revenue is maximized by servicing the higher-value requests first. Additionally, the acceptable service delay could vary with the type of service request. For example, instant messages may tolerate more delay than voice calls. In that case, it is not sufficient to handle the requests in terms of highest value first, but rather a trade-off between delay and value is required. For now, our prototype assumes that the service delay is the same for all request types.

II. SIP BACKGROUND

As shown in Figure 1: SIP architecture, a SIP infrastructure consists of *user agents* and a number of *SIP servers*, such as registration servers, location servers and SIP proxies deployed across a network. A *user agent* is a SIP endpoint that controls

session setup and media transfer. [6] describes the SIP protocol in detail.

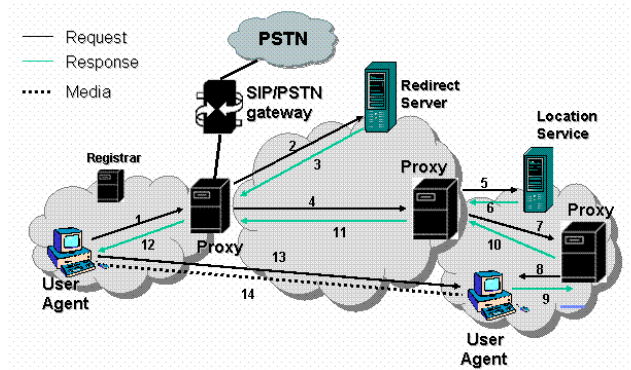


Figure 1: SIP architecture

All SIP messages are *requests* or *responses*. For example, INVITE is a request while “180 Ringing” or “200 OK” are responses. A SIP message consists of a set of headers and values, all specified as strings, with a syntax similar to HTTP but much richer in variety, usage and semantics. For example, a header may occur multiple times, have list of strings as its value, and a number of sub-headers, called parameters each with an associated value. In the following example from [6] Alice invites Bob to begin a dialog:

```
INVITE sip:bob@biloxi.com SIP/2.0
Via: SIP/2.0/UDP pc33.atlanta.com;branch=z9h
Max-Forwards: 70
To: Bob <sip:bob@biloxi.com>
From: Alice <sip:alice@atlanta.com>;tag=192
Call-ID: a84b4c76e66710@pc33.atlanta.com
CSeq: 314159 INVITE
Contact: <sip:alice@pc33.atlanta.com>
Content-Type: application/sdp
Content-Length: 142
(Alice's SDP not shown)
```

SIP messages are routed through SIP proxies to setup sessions between user agents. All requests (such as an INVITE) are routed by the proxy to the appropriate destination user agent based on the destination SIP URI included in the message. A session is setup between two user agents through an INVITE request, an OK response and an ACK to the response. This is shown in Figure 2 where the call setup is followed by media exchange using RTP. The session is torn down through an exchange of BYE and OK messages.

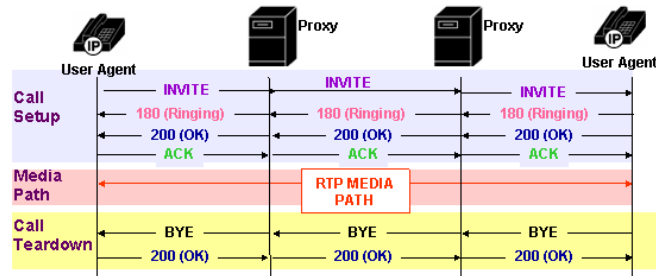


Figure 2: SIP Call setup and Media Path

SIP separates signaling from the media – signaling messages are carried via SIP, whereas media is typically carried as RTP over UDP [3]. Signaling messages are routed through the different SIP servers while the media path is end-to-end. The body of a session setup message (e.g., INVITE) describes the session using Session Description Protocol (SDP) [3]. The IP address and port numbers exchanged through SDP are used for the actual data transmission (*media path*) for the session. Any of these parameters can be changed during an ongoing session through a re-INVITE message, which is identical to the INVITE message except that it occurs within an existing session. The re-INVITE message is used most often in mobile networks to support handoff of an existing VoIP call due to user mobility (and subsequent change of endpoint addresses).

SIP messages primarily belong to three functional classes: (a) session setup/modification/teardown, (b) instant messaging and (c) event subscription/notification. RFC 3261 [6] defines the basic set of messages and interactions that define sessions, such as REGISTER (for registering a user agent), INVITE, and ACK for session setup, BYE for session teardown, and a variety of other control messages such as OPTIONS. Additional messages have also been defined (e.g., INFO[9], UPDATE [10], etc.). The MESSAGE request is an extension for ‘paging-mode’ instant messaging [1] and the more recent, Message Session Relay Protocol (MSRP) [1] defines methods for session-mode instant messaging. Another set of extensions enable presence applications with the PUBLISH, SUBSCRIBE, and NOTIFY primitives for event notification [1].

SIP can operate over multiple transport protocols such as UDP, TCP or SCTP. Use of UDP is probably more prevalent today especially for proxy-to-proxy connections, but TCP usage is expected to grow down the road. Additionally, when using TCP, SIP can use SSL (secure sockets layer) for security and encryption. It may also use IPSec underneath any of the transport protocols as well.

III. MOTIVATION

Overload is an inevitable condition for servers. Flash crowds, emergencies, and denial-of-service attacks can all initiate loads that exceed a server’s resources. Therefore, servers must be designed with overload in mind. Given that a server can not handle all of the requests it receives, it would be desirable for it to handle those requests which produce the most value for its operator.

For example, “911” emergency calls should take precedence over other calls; text or picture messages may generate more revenue than local calls; and dropped calls are more frustrating for users than “system busy” messages. Furthermore, each operator may have different policies and values associated with each type of message.

Our solution is to leverage the rich header information contained within SIP messages to classify the incoming stream of messages according to operator-defined rules; and then based on the classification deliver the highest priority messages to the server first. We achieve this with a novel SIP message classification algorithm described in detail in Section IV.

A key motivation behind the design of our classification engine, which is an implementation of the above algorithm, is that it must be programmable with a set of user-specified rules. In this paper, we use the classification engine to provide overload control. The need for programmability for overload control arises from our earlier observation that SIP can be used to support multiple services / applications and each provider may offer a different subset of these services and assign a different set of values (e.g., revenue metrics) to their mix of offered services. Clearly then, the rules for overload control that aims to maximize value for an operator under overload, must be different, and rather than create a specialized engine for each operator, it is eminently better to program a common classification engine with different sets of rules.

Because our classification engine is programmable, it can be used in multiple contexts besides overload control. For example, the classifier can be used as a SIP-aware load balancer in front of a SIP server farm to provide either transaction affinity or session affinity. It could potentially be also used to prevent denial-of-service attacks by programming it with rules that drop undesirable messages. The multiple different scenarios and their corresponding architectures where a fast, efficient classification engine could be useful, is currently under study. In this paper, we will provide a brief sketch of one additional scenario, namely for use as a session-aware dispatcher of requests for a SIP server farm.

Another key point of our classifier design is that it is independent of the underlying transport protocol. As mentioned earlier, SIP can operate over multiple transport protocols: our classification engine operates on SIP messages and assumes that the transport layer connections have been terminated and provide a single FIFO stream of messages as input to the classification engine. However, depending on the transport protocol used, additional transport-layer mechanisms will be needed. We discuss some of these issues briefly in Section III.B. In this paper, we assume that UDP is transport protocol used which requires no additional (transport-layer) mechanisms in order to apply our classification engine for overload control.

We would also like to point out that the classification algorithm is independent of the queuing policy used. The end-result of the classification process is to place an incoming message in one of multiple categories. In case of overload control, the categories will be realized as queues and they may be serviced using one of many possible queuing schemes such as weighted round-robin or priority schemes. Furthermore, usage of classification in contexts besides overload control may benefit from

schemes better suited for the specific context. Thus our design decouples the classification algorithm from the queuing mechanism/policy used.

It is worthwhile also mentioning here that a SIP server would normally parse a SIP message before processing it. Thus, a key goal for our classifier design is *not* to duplicate server functionality, and instead extract only the needed information (to enable rule matching) from a subset of the message headers, as will be described in detail in Section IV. This is especially relevant for use in overload control for two reasons: (a) the classifier should take away as few processing cycles from the SIP server and yet provide a disproportionately higher return, (b) when messages need to be dropped, it is better to drop them earlier in the processing path. The second point (b) ties in with the positioning of the classification engine within the overall server system: as will be discussed in Section III.C, the classifier engine is best positioned as an in-kernel module from a performance standpoint.

We expect the usage of the classifier not *require* any modification to the SIP server application (proxy, redirect server, presence server etc), i.e. the classifier is self-contained with its own rules. However, we also expect configurations where a SIP server would cooperate with the classification engine by programming rules and/or input to the classifier, such as routing policies. An example of such cooperation is in using the classifier for VoIP denial-of-service (DoS) prevention, where the SIP server may use its own methods to detect onset of DoS attacks, and thereupon, insert rules in the classifier to detect and drop undesirable messages before they are processed by the server.

Lastly, it should be noted that the classification-based overload control scheme is triggered only at the onset of overload, and not during normal operation. Since overload detection is non-trivial in itself, a detailed evaluation of detection schemes is outside this scope of this paper [22].

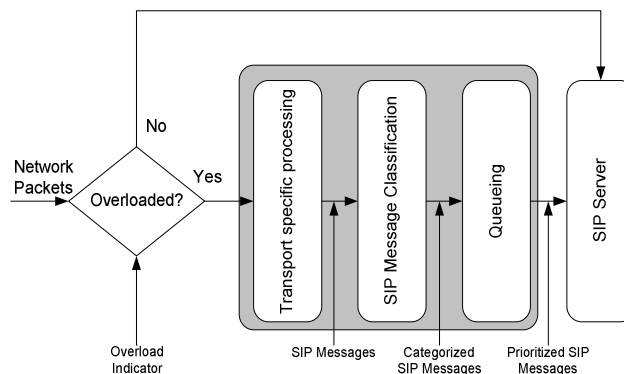


Figure 3: Overall context of classifier usage

A. Target Network Scenarios

We present three scenarios for our classification engine: (a) overload control for a SIP server servicing many clients, (b)

overload control for a SIP server interconnected with another SIP server, and (c) as a dispatcher for a SIP server farm. In (a), each client has a separate transport connection to a SIP server such as a proxy. The case for overload in this scenario arises when a proxy is supporting a large number of users, such as a VoIP service provider like Vonage, and clients come back up at roughly the same time after a regional loss of network connectivity (and thus requiring the clients to re-register). In (b), SIP servers in the core of a service provider network, for example, will receive requests from SIP servers in other administrative domains and/or the same domain, and this request volume is much higher than an access server which supports only user-agents. For a given server-to-server interconnection, all messages are sent over a common transport connection. In both of these cases, our classification engine would be collocated on each server, and programmed with a set of rules to provide overload control (maximizing revenue). Additionally, servers in (a) or (b) may be organized as a server farm front-ended by a SIP-aware dispatcher, whose job is to distributed requests amongst the servers, maintaining for example, session affinity. In this scenario (c), the dispatcher can also be realized through the same classification engine, but programmed with a different set of rules (to ensure for example, the INVITE, 200 OK and ACK messages constituting a session setup are sent to the same server). The key point to be made here that since the classification engine is programmable; it can be used in multiple ways with appropriate configurations. This paper will focus on use of the classification engine for overload control.

B. Impact of transport layer

As mentioned earlier, SIP can use multiple transport protocols and optionally security support. When SIP is used over UDP, each SIP message is completely contained within a UDP packet [6] and messages from multiple hosts arrive on the same socket, i.e. the classifier can pick up messages from this single socket. In case of TCP, data is delivered to the SIP proxy through a socket interface as a byte-stream. The overload protection mechanism thus needs to be interjected between the TCP implementation and in-kernel socket data structures, so that the byte-stream coming out of the TCP connection can be recognized as a series of messages which are then acted upon by the classifier (and re-ordered). Additionally, when TCP is used, each connection to the proxy results in a separate socket data structure. The classification of messages needs to be done *across* multiple connections. TCP connections could also be encrypted using SSL, especially for client-to-proxy connections. The classifier must inspect clear text (i.e., unencrypted) SIP headers to make its decisions, so SSL connections must be terminated by a component that is interposed between the user agent and the classifier. This component will terminate the SSL connections, and then forward unencrypted data to the classifier using a secure channel. Handling SSL connections is outside the scope of this paper: multiple options exist such as simple user-space SSL termination [11], in-kernel SSL termination, or a dedicated SSL termination server. For server-to-server interconnections, a transport connection (TCP or UDP) is likely to be secured by IPSec. Because IPSec is implemented in the kernel, the classifier can analyze the SIP messages in the kernel.

C. User-level vs. in-kernel design

The exact realization of this integration depends on whether it takes place as a user-space process, or in-kernel. However, executing the classifier in-kernel provides performance benefits because packets are not copied. In both cases, it should be noted that the classification *algorithm* remains exactly the same.

When applied to overload control, a kernel-level implementation is preferable, because overload control generally entails dropping messages in order to reduce load. Clearly, message dropping needs to happen as early in the processing path of a message to minimize the amount of processing (CPU, I/O, etc.) resources spent on a message that will ultimately be dropped. For this reason, an in-kernel realization of the classification algorithm can be expected to provide higher performance gains than say integrating an overload control module within a SIP server (in user-space). Additionally, a kernel-level implementation is more flexible and portable than a user-level implementation, because it can be applied to multiple unmodified SIP server implementations.

IV. OUR CLASSIFICATION ALGORITHM

The input to the classification algorithm is a set of rules, expressed as a conjunction of conditions. Our classification algorithm has a static and a runtime component. The static component consists of rule parsing and creating several tables and bitmaps that allow the runtime portion to operate efficiently. Our algorithm uses three tables: a *header* table, *header value* table and a *condition* table that store the required message headers, values of those headers and the conditions to evaluate, respectively. We use a list of rules that are expressed as bitmaps, where each bitmap represents the conditions that must be true for the rule to be matched. Each rule has an associated set of actions, one of which must be a priority for the packet.

The runtime component consists of extracting only those headers (and their values) from a SIP message that are present in the header table, evaluating the conditions in the condition tables, storing the results of the condition evaluations in a bitmap, and then comparing that bitmap with each rule. When a rule matches, a set of actions is applied. We make use of multi-pattern matching algorithms to extract only the necessary header values through a single scan of the message, i.e. we parse only the relevant headers and extract only the necessary sub-fields.

A key feature of our algorithm is that by supporting user-defined data types (such as associative arrays); we are able to maintain user-specified state that can be updated as part of rule actions.

In the remainder of this section, we describe our algorithm in detail. Section IV.A describes message headers and Section IV.B describes data types. Classification rules are described in Section IV.C, and their grammar in Section IV.D. The static and run-time phases of our algorithm are described in sections IV.F and IV.G, respectively. We illustrate the algorithm through an example in section IV.H.

A. Message Header Types and Specification

The header table consists of a list of SIP message headers, which we have classified according to three types:

- Simple headers, which exist as such in SIP message, e.g. From, Via, and Call-ID.
- Pseudo-Headers: These headers do not appear as such in a SIP message, but allow us to refer to certain strings or characteristics of the message. For example, Method is a pseudo-header that we define to represent the type of the SIP message such as an INVITE or a 200 OK.
- Derived Headers: These are constructed from one or more headers. The first kind of derived header is what we call as “sub-header” and comprise of pre-defined composite strings of the form X.Y where X represents a simple header, and Y represents either one of the specified parameters [6] than can appear in the value for header X (e.g., the “tag” parameter of the From header), or a list of tokens that we defined to represent specific values of interest in the header (e.g., “From.URI” represents the URI portion of the From header value). The second kind of derived headers consists of user-specified ordered list of Message headers (either simple headers or recursively, other derived headers). For example, a SIP session dialog comprises of the tag parameter values of the From and To headers, and the Call-ID header value:

```
Dialog= {From.tag,To.tag,Call-ID}
```

Here, From.tag and To.tag are sub-headers, and Call-ID is a simple header. Individual elements of a user-defined derived header are indicated by dotted notation, such as Dialog.From.tag or Dialog.Call-ID.

B. Data Types

In conjunction with user-defined derived headers, we also allow the user to specify complex data types such as structures and complex data variables such as associative arrays (hash tables), pointers and scalars. The basic data types are string and integer. Any time a user defines a derived header like Dialog, a type of the same name is also implicitly created. A structure is a data type consisting of a collection of data types.

```
For example: Struct Session = {Dialog Dialog1, String State}
```

The element “State” stores the state of a dialog which could be “*established*”, “*setup*”, or “*shutdown*”. To differentiate variables, we prefix each variable with “\$”, “*”, or “%” for scalars, pointers, and associative arrays, respectively. Associative arrays must be defined in terms of a structure. The first element of the structure is the key, and the remaining elements are the values. A list can be created by using an associative array with keys, but not values. Pointers can only reference elements within an associative array.

All variables are assumed to be global in scope unless explicitly specified to be of local scope. A variable with a global scope exists for the lifetime of the classification process, i.e. it retains its existence across messages and can be modified as a result of

classifying individual messages. Global variables are typically used to maintain state that is dynamically generated and modified by the classifier (e.g. an associative array of dialog-ids for ongoing SIP sessions). In contrast, a locally scoped variable does not share its value across messages, and in fact, retains its value only within the context of a specific rule execution. Moreover, multiple instances of a local variable can be concurrently instantiated (e.g., if the classifier is running on a multi-processor system it can process multiple packets at the same time). Local variables such as pointers are used to extract an element of a globally scoped list that matches with some set of header values in the message currently under classification.

For example, to maintain a count and an array of session information, three variables could be used.

```
Int: $ActiveSessionCount
Session: %ActiveSessions
Local Session: *CurrentSession
```

A scalar integer suffices for the count. An array (`ActiveSessions`) would use the `Session` structure previously defined, using a local reference (`CurrentSession`).

The classifier maintains a global variable table (GVT) and a per-message local variable table (LVT). Each entry in the GVT or LVT can store a basic type (string, integer), tuple, or reference to a list entry. The run-time classifier elements store only indices to the GVT or LVT (much like compiled code references variables by memory locations rather than names).

C. Rules

A rule consists of a conjunction (AND) of conditions resulting in an action, along with a priority for each rule. For example:

```
C1 AND C2 AND C3 → A1, Priority
```

Disjunctions (OR) do not need to be supported since disjunction of conditions can be expressed as separate rules, without loss of generality. Rules are applied sequentially, until a matching rule is found. If no match is found, the packet is given the lowest possible priority.

Conditions are of the form “Header op Literal.” The header may be a simple header, pseudo-header, or a derived header. The operator can be `==`, *subset*, *superset*, or *belongs-to*. The *belongs-to* operator also supports an optional assignment to a pointer. We also support negation for `==`, *subset*, *superset*, and *belongs-to* (without assignment). String equality can only be used for headers that have a single value. For multi-valued headers such as `Record-Route`, `Via`, or various authorization headers, the *subset* and *superset* operators are used. For example to match messages that traverse only “host1.watson.ibm.com,” “host2.watson.ibm.com,” or both the condition, “Via subset {host1.watson.ibm.com, host2.watson.ibm.com}” is used. To express set equality, the conjunction of *subset* and *superset* are used.

The belongs-to operator is used to find headers in the classifier’s state lists (i.e., a list of Dialog-IDs). For example, “Dialog-ID belongs-to %ActiveSessions” expresses the condition that Dialog-ID is a key in the associative array ActiveSessions. This operator returns a “true” value by returning a pointer to the element in the list that matches the Header; it returns a “false” value if no match is found. Thus, it serves a dual-use of evaluating the Boolean value of a condition and, in addition, returning a pointer value. Thus, we support a special assignment operator, =, that may prefix a belongs-to condition. For example, “*CurrentSession = (Dialog-ID belongs-to %ActiveSessions)” assigns the found item to the *CurrentSession pointer.

D. Actions

The actions that an implementation classifier provides necessarily includes the ability to color messages with a category (e.g., the priority for overload control), but our classification algorithm does not define what the set of actions is. We have chosen a general representation for our actions, an action table which is a series of actions represented as three-address code. Each action consists of a type, left-hand side, right-hand side, and a next pointer. The type of the action defines the instruction. For example, our in-kernel implementation defines priority assignment; variable assignment; arithmetic operations such as addition, multiplication, division, and modulus; tuple allocation, assignment, and extraction; variable assignment; array insertion and deletion, and more. The left-hand side of an action is a variable which may have side effects (e.g., for addition the left-hand-side serves as both an addend and where the result is stored). As not all instructions affect variables, the left-hand-side may not be specified (e.g., priority assignment does not alter any variable). The right hand-side is used as input and may be an immediate variable (e.g., the constant 1), a state variable (e.g., \$ActiveSessionCount), or a header, sub-header, or derived header (e.g., From.tag). For example, the action “ADD \$ActiveSessionCount 1” adds one to the value of \$ActiveSessionCount. Each rule has a pointer to an action that is executed when the rule is matched, and each action has a next pointer to another entry in the table. To prevent loops only backward pointers are allowed and the first entry (i.e. 0) terminates the action. For clarity, the rule compiler should allow actions to be specified using a richer and more complex syntax (e.g., “\$A = \$B + \$C”), but should convert the complex action into a series of simpler actions (e.g., “ASSIGN \$A \$B”, “ADD \$A \$C”).

E. Rules Set Specification Syntax

A complete set of rules begins with type definitions, a list of user-defined headers, variable declarations, and finally an ordered list of rules and actions. The BNF grammar of our rule language is shown in Figure 3. We use Italics for grammar symbols, bold characters for string literals, and roman type for alphanumeric strings (e.g., identifiers). The starting symbol is “RuleSet”, which is made up of type declarations, variable declarations, and one or more rules.

```

RuleSet = TypeDeclaration VarDeclaration* Rule*
TypeDeclaration = (UserHeader|Structure)*
UserHeader = TypeName = { (Header|TypeName) (,(Header|TypeName) ) * }
Structure = Struct TypeName = { TypeName FieldName (, TypeName FieldName) * }
VarDeclaration = (Local)? TypeName: Kind VarName (, Kind VarName)*
Kind = $|%/ *
Rule = Condition (AND Condition)* → Action (, Action)*
Condition = Header (==|!=) String | Header (subset|superset) {String (, String)*} |
            Header belongs-to %List | VarName = (Header belongs-to %List)
Action = (Assignment | Function | Priority)*
Assignment = (VarName = Value | VarName.FieldName = Value)

```

Figure 3: Rule Grammar

F. The Static Phase of the Algorithm

1. Extract a set C of unique conditions from the rule set specified.
2. From the set C, extract a set H of unique headers, which may be pseudo-headers, simple headers and derived headers. For each derived header in this set, recursively include the list of simple (or derived) headers that comprise that derived header. For example, the Dialog derived header defined in Section IV.B includes the derived headers From.tag and To.tag as well as the simple header Call-ID. The derived headers From.tag and To.tag would recursively lead to the inclusion of “From” and “To” in the set H.
3. Create a table, *Header Table*, whose each row consists of a header from the set H. The format of each row is $\langle \text{Header}, \text{Header-Type}, \text{List of indices}, \text{fn} \rangle$. Header is the actual string representation, such as “From”, or “From.tag”, “Dialog”, or “MSG_TYPE”. Header-Type refers to simple, derived or a pseudo-header. For derived headers, there is a corresponding ordered list of indices referring to the simple and pseudo-headers (or other derived headers) comprising a derived header, e.g. Dialog would refer to the indices for “From.tag”, “To.tag”, “Call-ID”. For pseudo-headers and derived headers, the element *fn* refers to a function that can extract the value of the header from the message (for pseudo-headers) or its component simple and derived headers (for derived headers). For example, the function pointer for Dialog encodes the necessary logic to create a dialog ID by walking through the associated list of indices. The entries in the header table are ordered such that the components of a derived header are always computed before the derived header.
4. Associated with the *Header Table* is another table, *Header Value Table*, which for every header (index) in the Header Table will eventually hold a value in the Header Value Table, e.g. the “From” header in the Header Table will contain “sip:xiping@us.ibm.com” in the associated entry for the *Header Value Table*. These values will be populated during runtime, i.e. when a message is being classified. For pseudo-headers and derived headers, the associated function *fn* when

executed will place the result (value) of the execution in the corresponding location/index in the *Header Value Table*. Each row, in addition to the value, also contains a type of value which could be a string, list of strings, tuple, integer, or NULL.

5. Create a table, Condition Table, whose each row represents a condition from the set C and consists of:
- <operator, header, literal, Assignment-variable>*

The aim is to efficiently represent conditions by storing pointers to header values of a SIP message under classification. At run time, the requisite header values can be referenced in constant time for efficient evaluation of these conditions. In general, *header* is an index to a header table, and the *literal* is a fixed operand that the *header* is compared to. The *operator* is one of the operators defined earlier:

- a) For string (in)equality operators, *literal* refers to the literal string that is being compared to a specific SIP message header value, which is specified by an index (*header*) in the Header-Table (HT). The fourth element of the row is unused. An example of this type of entry is *< ==, 0, "Charles", NULL >* representing the condition "*From == 'Charles'*". Where, HT[0] represents the header "From".
 - b) When the operation is *belongs-to*, *op1* refers to a header value and *op2* refers to a list. For example, the condition with assignment, "**S1=(Dialog belongs-to %L1)*", will be represented by a row in the condition table as *<belongs-to, 5, %L1, *S1>*. Where, HT[5] refers to the derived header Dialog (defined earlier).
 - c) For *subset* and *superset* operators, *header* is a message header that is list-valued such as "Via", while *literal* is a list of values. For example, a condition such as "*Via subset {proxy1, proxy2}*" will be represented as *<subset, 3, {proxy1, proxy2}, NULL>* where HT[3] is the entry in the header-table representing the "Via" header.
6. Bit vector representation: The set of conditions C is efficiently represented as a bit-vector (*Condition bit-vector*), where bit *i* refers to the *i*th condition in the Condition-Table, and will be set 1 iff that condition is true for a message being classified. Additionally, for each rule, create a *Rule bit-vector* where the *i*th bit is 1 iff the rule specification includes the *i*th condition.

G. Run-time: classification actions per-message

1. For each header in the header table, determine whether the header exists in the SIP message and if so, return a pointer to the header value in the message. This determination can be done by using an efficient multi-pattern string matching algorithm like SBOM [15] [16] or even a simple switch statement (as there are a limited number of pre-defined SIP headers). The advantage of doing this is (a) while the entire message needs to be *scanned* (with sliding window that can step over multiple characters at a time as dictated by the SBOM state machines), the entire message is not *parsed*. The end-result of this step is to populate each

header entry in the header value table with a pointer to the position in the SIP message corresponding to the value of the matching header. For pseudo-headers and derived headers, we perform limited parsing of the extracted simple headers and store the result in the corresponding index in the header table.

2. Walk through each row (condition) in the *Condition table* and set the corresponding bit in the *Condition bit-vector* to 1 if the condition is true or 0 if the condition is false. As was indicated earlier, our data structures were carefully designed to make evaluating the condition at run-time efficient: each entry in the *Condition Table* has pointers to the condition's operands (i.e. a literal and a header). For the *belongs-to* operator, the literal operand in the previous section is an associative array, i.e. a hash-table which enables us to evaluate this operation efficiently. We store references to assignment variables as integer indexes into the global or local variable table, thus retrieving and storing variables is efficient.

3. Next, we compare the *Condition bit-vector* (CBV) to each rule bit-vector. A rule matches the SIP message under classification, if and only if $(R \ \& \ CBV) == R$, where R is the rule bit vector for that rule and “&” is the bitwise *and* operation. In words, if the i^{th} bit is 1 in the bit-vector for rule R , then the same bit must be 1 in the condition-bit-vector. If R contains the i^{th} condition, then the value in the matching header of the SIP message must cause this condition to be true, in order for R to apply. Because the rule-bit-vectors are sorted according to priority, the matching process can be stopped after the first matching rule (because that is the highest-priority matching rule)

H. Data Structure Example

We now sketch an example using the following rules:

- `Method=="INVITE" AND To.tag==NULL` → High
- `Method=="INVITE" AND From.URI=="sip:carol"` → Medium
- `From.URI == "sip:alice@atlanta.com"` → Low

The set of headers H for this rule set is initially “Method”, “To.tag”, and “From.tag.” Because “To.tag” and “From.tag” are derived from “To” and “From”, respectively, the “To” and “From” headers are added to H forming the set “Method”, “To”, “From”, “To.tag,” and “From.tag.” The unshaded portion of the following table represents the static portion of the header table.

Index	Parent	Header	Value
0		Method	INVITE
1		To	Bob <sip:bob@biloxi.com>
2		From	Alice

			<sip:alice@atlanta.com>;tag=192
3	1	To.tag	NULL
4	2	From.URI	sip:alice@atlanta.com

The conditions are Method == "INVITE", To.tag == NULL, From.URI == "sip:carol", and From.URI == "sip:alice@atlanta.com". The unshaded portion of the table below represents the static portion of the condition table.

Index	Op	Header	Literal	Value
0	==	0 (Method)	INVITE	1
1	==	3 (To.tag)	NULL	1
2	==	4 (From.URI)	sip:carol	0
3	==	4 (From.URI)	sip:alice@atlanta.com	1

Note that headers and conditions that are duplicated within the rule set are only expressed in the header and condition tables a single time. Using the condition table, the rules can be expressed as a bitmap. For example Method == "INVITE" AND To.tag == NULL is expressed as 1100, with each bit corresponding to an index in the condition table. Similarly, Method == "INVITE" AND From.URI == "sip:carol" and From.URI == "sip:alice@atlanta.com" are expressed as 1010 and 1001, respectively.

The shaded columns in the tables represent the run time state for the message, using the sample message from Section II. The value column in the header table is the header value table, and the values from the condition table represent the condition bit vector, 1101. The condition vector is compared to each of the rules in turn. The first rule matches because $1101 \& 1100 == 1100$, so processing may stop.

This condition bit vector does not match the second rule, because bit 3 is not set ($1101 \& 1010 != 1010$). The third rule does match ($1101 \& 1001 == 1001$), but is not executed because the second rule already matched.

Next, we show a skeleton example of using state in the classifier. The following two rules describe how state consisting of the dialog-ids of ongoing sessions, is updated each time a new session starts or an existing session terminates. The first rule detects a new session and adds it to a dialog session list. Once the session is completed upon receipt of BYE message, the session is removed from the session list. (Clearly, additional rules will be needed for completeness.)

- Method=="INVITE" AND Dialog-ID !belongs-to %ActiveSessions) → addDialog(Dialog-ID, %ActiveSession)
- Method=="BYE" AND Dialog-ID belongs-to %ActiveSessions → delete(Dialog-ID, %ActiveSessions)

V. OVERLOAD CONTROL IMPLEMENTATION

A. Classifier Prototype

We have developed a prototype classifier for overload control that consists of three components: (1) the core of the classifier a Linux kernel module responsible for parsing and classifying messages, (2) a user-level rule parser, and (3) a kernel patch that provides an extensible priority queue for UDP sockets. The kernel module is 3,425 lines of code, provides support for parsing all defined SIP headers using a switch statement, associative arrays using linear hashing [21], all of the types, operators, and actions described in the previous section, and several additional actions. The user-level rule parser is 2,121 lines of C code that parses a set of rules and compiles them into a header table, condition table, rule list, and three-address-code action set. Our kernel priority queuing extension adds 626 lines of kernel code that adds a new socket option (SO_QDISC) which allows servers to specify which classification rule set should be used. When a packet is received over the network it is classified using our kernel module, and then inserted into one of n queues. When the server reads from the socket, higher priority messages are returned first. Additionally, if sufficient room is not available in the socket's buffer, lower priority messages are dropped in favor of higher-priority messages.

B. Testbed / Workload Used

We ran our classifier and SIP Express Router (SER) 0.9.6 on a dual 3.0 GHz Xeon with 4.5GB of RAM. We used two identical 1.7 GHz Pentium IVs with 512MB of RAM to send and receive the messages. All machines were connected via a 1Gbps Ethernet network.

The workload we used to evaluate our classifier is reflective of what would be typically used by a mobile service provider during call overload, which is to reduce the number of dropped calls due to handoffs in preference to new call setups. Call setup is modeled via SIP INVITE messages; while a call handoff is modeled via re-INVITE messages. The re-INVITE message is structurally same as an INVITE but with a non-empty value of the "tag" parameter in the To header. Prioritization of call-handoffs (re-INVITE) over call setup (INVITE) is accomplished via the actions associated with the rules that match one or the other message type.

We evaluated our classifier by sending a sequence of SIP messages to it that consists of INVITEs and re-INVITEs from one SIP client to another through SER with and without our classifier. We configured our classifier with the following set of rules that (a) differentiate re-INVITE vs. INVITE messages using the To.tag field and (b) prioritized re-INVITE messages over INVITE messages :

```
10: Method == "INVITE" && To.tag == NULL -> Color 1
20: Method == "INVITE" && NOT To.tag == NULL -> Color 0
```

```
30: NOT ReqResp == NULL -> Color 1
```

```
40: ReqResp == NULL -> Color 2
```

Rule 10 matches INVITE messages without a To.tag field (i.e., the INVITE messages) and assigns them to color 1 (the highest priority is zero). Rule 20 matches INVITE messages with a To.tag, and assigns them priority 0 (i.e. the highest priority). Rule 30 matches all other SIP messages and assigns them to the same priority as INVITE. Finally, rule 40 matches any non-SIP or malformed messages and assigns them the lowest priority.

When overload occurs, some messages are dropped. We record the number of messages of each type received. We compare the total number of messages received, the number of INVITES received, and the number of re-INVITES received. We measured our classifier using from 10,000-120,000 messages/sec (in increments of 10,000 messages/ second), which demonstrates its behavior both with and without overload.

C. Observed results

Fig. 4 shows the results our experiment with a mix of 75% INVITES and 25% re-INVITES (i.e. a handoff-ratio of 25%). We selected these ratios as in practice call handoffs happen less often than call setups. As can be seen in the figure, before overload is reached our classification engine has no impact on the number of messages processed by SER. However, after overload is reached, the number of messages/sec processed decreases by 8.8% for 40,000 messages/sec to between 18.7-23.3% for 60,000-120,000 messages/sec. The reason that throughput is decreased is that the classifier must process all incoming messages, and we observed that the classifier was able to handle 104,891 messages/sec at its peak, more than 2.6 times as many as SER could handle at its peak; at the same time SER was processing 31,616 messages/sec.

Because the classifier was able to process so many more messages than SER alone, it was able to select the high-value messages in this workload for processing, increasing the number of hand-offs processed by 50.9% for 60,000 messages/sec to 160.2% for 120,000 messages/sec. Of course this comes at a cost, a corresponding number of setup messages can no longer processed, and thus the call setup throughput decreased between 11.7% and 79.5%.

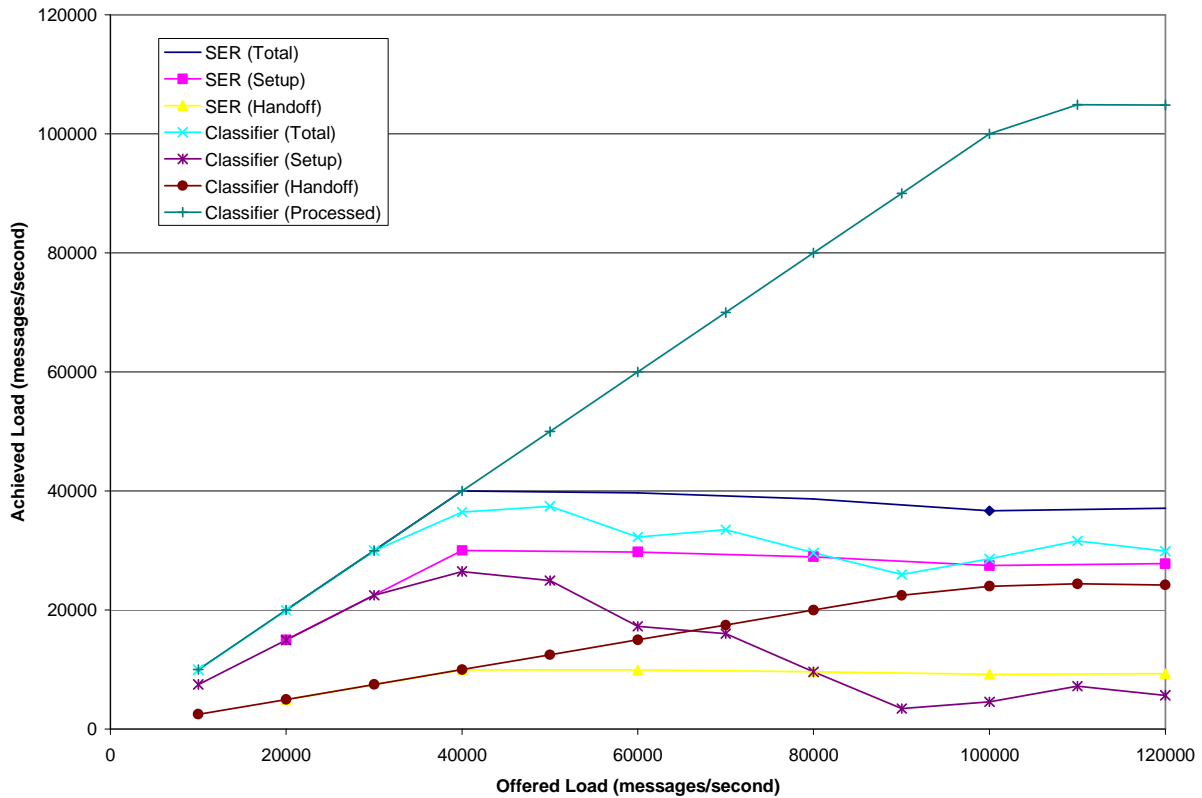


Figure 4: 75% INVITEs, 25% re-INVITEs

The performance of the classifier is dependent on the incoming stream of messages; therefore we also ran the same experiment using handoff ratios of 0%, 25%, 50% and 100% and recorded the maximum number of call handoffs that were processed. The results are shown in Figure 5, Maximum handoff operating region, as a function of the input stream. The region between the peak handoff capacity of SER and the peak handoff capacity of our classifier is the additional operating region provided by the classifier. As can clearly be seen, the peak number of handoffs scales linearly with the ratio for SER, because the messages SER processes are randomly selected from the input stream. The classifier is able to select the high-priority handoffs from the message stream and processes all of the available handoffs until the SIP server itself is saturated at 40,000 messages/sec, at which point no more capacity is available for additional hand-off messages to be processed.

It is important to note that the classifier can provide no benefit at ratios of either 0% or 100%, because there are no high-priority messages to select or no low priority messages to discard, respectively. As expected, the classifier is most effective, when the offered load is higher than the server's capacity, but number of offered high-priority messages is below the servers capacity (in this case a handoff ratio of 37.5%).

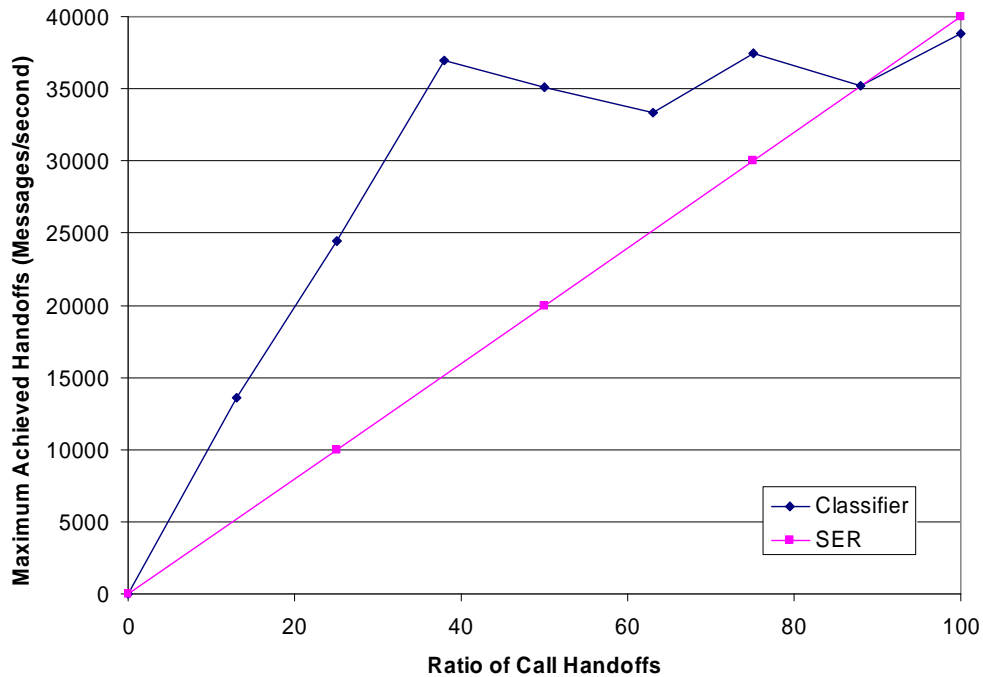


Figure 5: Maximum handoff operating region, as a function of the input stream.

VI. USE OF CLASSIFIER AS A SIP-AWARE DISPATCHER

As mentioned before, the programmability of the classifier through rule-sets makes it applicable to multiple scenarios. We have described one such scenario, namely overload control and demonstrated the performance gains achievable with the classifier. In this section, we briefly outline a second use-case for the classifier, and describe a partial rule-set to program the classifier.

In many configurations, it is necessary for functional correctness to dispatch all messages that belong to the same transaction to a common server¹ (e.g, messages that comprise the INVITE/OK/ACK three-way handshake for a call setup). Moreover, for correct accounting, it may be necessary to dispatch all messages related to the same call to the same server. This use-case consists of a set of SIP servers that is front-ended by a dispatcher, as shown in Fig. 5. A SIP-unaware dispatcher is not suitable, because for proxy-proxy interconnections many sessions will operate over a single connection. We construct a SIP-aware dispatcher using our classifier to determine which server a message should be sent to. After the message is tagged with a server by the classifier, the existing dispatcher framework can forward it as normal.

This is accomplished by programming the classifier that creates state in the classifier for the first message in a call or call setup transaction), based on the Call-ID message header. The call-ID header value is required to be unique for each call and is the same

¹ SIP message exchanges consist of transactions such as INVITE/OK/ACK, and such transactions can be either stateful or stateless. When stateful, all message for a transaction need to be handled by the specific server that maintains state for that transaction. A session consists of a sequence of transactions, such as INVITE/OK/ACK followed by BYE/ACK.

across all messages in a transaction [6] (and also across a session). This state is represented by the *Session* structure in the rules below, and includes the specific server id to which the INVITE will be forwarded. All succeeding messages for this transaction (e.g. RINGING, OK, ACK), matching this call-id, will be dispatched to the same server, thereby providing transaction affinity. The key aspect of the classifier that is being leveraged is the ability to dynamically maintain classification state through an in-kernel associative array [21].

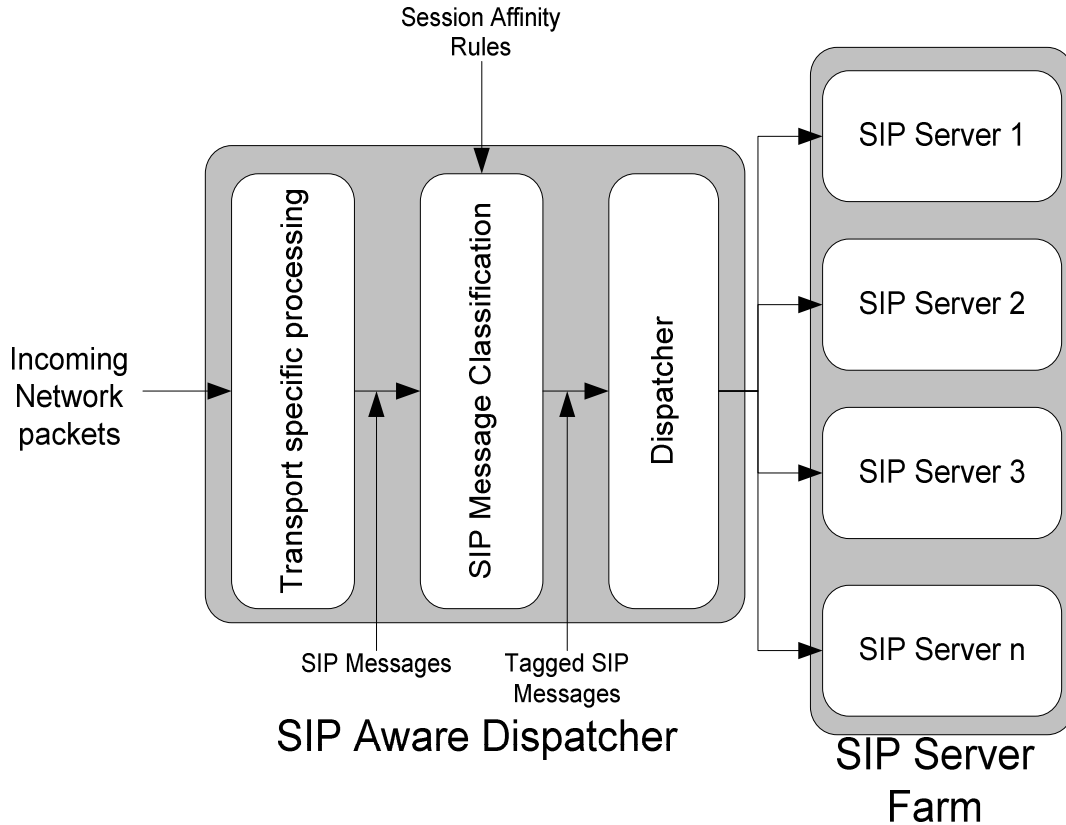


Figure 5 : SIP Dispatcher

The following annotated rule set demonstrates how state can be used. First, a session structure is defined and an associative array named `%ActiveSessions` is created with `ID` as the key. Two local variables are also required: (1) `$NewSession`, which is a temporary entry that can be inserted into `%ActiveSessions`, and (2) `*CurrentSession` which is used as a pointer into the `%ActiveSessions` Array.

```
Struct Session = { String ID, Int Server, Int Expire }
Global Session: %ActiveSessions
Local Session: $NewSession, *CurrentSession
```

We also define three integers: (1) `$MyServer` which is the server the current packet should be sent to, (2) `$CurrentServer` is a round-robin counter to evenly distribute the load across servers, and (3) `$nServers` is the number of servers in the server farm.

```
Local Int: $MyServer
```

```
Global Int: $CurrentServer, $nServers
```

The Init function is used to initialize variables (e.g., setting the number of servers to three), and create a kernel thread to remove array entries that are no longer needed.

```
Init -> $CurrentServer = 0, $nServers = 3, ExpiryThread(%ActiveSessions, Expire)
```

Rule 10 matches calls that have an entry in the %ActiveSessions array, and its action sets the Expire element of the entry to the current time plus 15 minutes (900 seconds). Finally, the message is colored with the Server element of the entry, so that the dispatcher knows to use the same server as the previous messages in this call.

```
10: *CurrentSession = Call-ID belongs-to %ActiveSessions
    -> *CurrentSession->Expire = Now() + 900
        Color *CurrentSession->Server
```

Rule 20 matches calls that do not have an entry in the %ActiveSessions array. First, a server is selected and stored in MyServer using a round-robin algorithm. Next, a new session structure is allocated and inserted into the %ActiveSessions array, using the Call-ID as the key. Finally, the \$MyServer variable is used to color the packet, thus indicating which server should be used by the dispatcher.

```
20: NOT Call-ID belongs-to %ActiveSessions
    -> $MyServer = $CurrentServer++ % $nServers
        $NewSession = (Call-ID, $MyServer, Now() + 900)
        Insert(%ActiveSessions, $NewSession)
        Color $MyServer
```

Rather than using expiration only, it is possible to modify the rules such that messages which indicate the end of a transaction trigger state deletion. For example, the following rule removes entries when a final response is received for non-invite transactions (INVITE transactions require an ACK as well):

```
5: Response >= 200 && CSeq.Method != "INVITE"
    && *CurrentSession = Call-ID belongs-to %ActiveSessions
    -> Remove(%ActiveSessions, *CurrentSession), Color *CurrentSession->Server
```

Of course, a timer to purge old state is still required as some transactions may never complete, and additional rules will be needed for corner cases (e.g., dropped messages, INVITE transactions, etc.).

VII. FUTURE PLANS

Programmability of the classifier lends itself for use in multiple scenarios. Our current work revolves around developing support for such scenarios, including design and evaluation of the corresponding rule sets, e.g. denial-of-service protection for VoIP. This

includes identifying scenarios where the classifier and SIP server work cooperatively. Additional systems-level work includes incorporating support for TCP connections as well as SSL[11]. We are also investigating better methods of detecting overload[22], and specializing the classifier for specific SIP servers such as a Presence server, which receive a narrower class of SIP messages, but with richer information in the payload (such as the Presence Information Document). We believe this work opens up a rich set of possibilities to enhance SIP server performance. DoS

VIII. RELATED WORK

Related work primarily includes IP packet classification, HTTP header inspection in web-proxies, and parsing of SIP messages within SIP proxies/servers such as SIP Express Router (SER) [12]. We are not aware of any earlier work on SIP message header classification per se. SIP proxies and libraries, use efficient parsing techniques such as lazy parsing which include parsing *up to* a required header and/or incremental parsing. However, in our case, the classification engine needs to extract a (small) subset of the header values and thus, multi-pattern matching algorithms are better suited [15] [16]. Additionally, extraction of information from the SIP message is only one aspect of our algorithm. We share a similar bit-vector representation for rules as in [7]; however, unlike [7], we operate on string-based header value pairs, with no predetermined ordering of headers, and our basis for creating rules by extracting a common set of conditions from the rules is conceptually different from creating numeric ranges of interest (e.g. port numbers, IP addresses). SIP messages are syntactically similar to HTTP headers; however the diversity, scope and semantics of SIP headers are much larger than HTTP. Web proxies typically use the content URL in a HTTP request for forwarding it to the right web server; since the number of headers in HTTP *requests* is fairly small, they usually simply inspect all headers.

A key feature of our classifier is its ability to create, manage and update state across multiple messages in a very general fashion. None of the related work cited support this feature since there was not any need for that. However since the notion of transactions and session being an integral part of SIP, this feature is an essential requirement for SIP classification.

There has been recent related work on overload control of SIP servers. The solution proposed in [17] uses queue-length thresholds within a SIP proxy to determine congestion, and during congestion, it separates INVITE messages from the rest and returns a 503 “service unavailable”. Current discussion in the IETF [18] [19] centers on creating an overload control framework and adopting appropriate new message headers to convey additional information beyond just sending a 503 response. The approach taken in [20] leverages the fact that throughput of a SIP proxy is higher when processing requests in a transaction-stateless manner, and thus their solution consists of handling a subset of requests statelessly during onset of congestion, thus trying to avoid overload. In contrast to the aforementioned work, our focus is on creating a mechanism in the form of a programmable engine that enables specific user-defined policies to be executed efficiently, without requiring any modification of

the SIP server per se.

IX. CONCLUSION

In this paper we have presented an algorithm for efficiently classifying SIP messages using a *programmable* set of rules, and applied it for overload control. Our algorithm consists of a static phase and a run-time phase. In the static phase, we define a header table that is a list of attributes to extract from the message and a condition table. These tables eliminate redundancy that is often found in classification rule sets. At run time, classification consists of directed parsing to extract only the relevant headers from the message, evaluating each unique condition, and efficient rule matching using bit vector representation for rules. We implemented an *in-kernel* Linux prototype of the algorithm and programmed the classifier prototype with rules to prioritize handoff messages over call setup messages. Our detailed performance evaluation shows that the in-kernel classification engine is able to process more than twice as many messages than the application-level SIP server, thus significantly extending the operating capacity of the server for high-value messages in a *transparent* manner.

ACKNOWLEDGMENTS

The authors wish to thank Erich Nahum, John M Tracey, Zhen Xiao and Vijay Balasubramanian for the many discussions on this paper and providing valuable input.

REFERENCES

- [1] B. Campbell et al. *The Message Session Relay Protocol*. SIMPLE Working Group Internet-Draft, Jan 2004.
- [1] B. Campbell, editor. *Session Initiation Protocol (SIP) Extension for Instant Messaging*. RFC 3428.IETF, Dec 2002
- [3] Handley, M. and V. Jacobson, *SDP: Session Description Protocol*, RFC 2327, IETF Apr 1998.
- [4] A. Niemi et al., *Session Initiation Protocol (SIP) Extension for Presence Publication*. SIMPLE Working Group Internet-Draft, June 2003.
- [5] H. Schulzrinne et al. *RTP: A Transport Protocol for Real-Time Applications*. RFC 1889.IETF, Jan 1996.
- [6] J. Rosenberg et al. *SIP: Session Initiation Protocol*. RFC 3261. IETF, June 2002.
- [7] T.V. Lakshman, and D. Stiliadis, *High-Speed Policy Based Packet Forwarding using Efficient Multi-Dimensional Range Matching*, Proceedings of ACM SIGCOMM'98
- [8] G. Camarillo and Miguel-Angel Garcia-Martin. *The 3G IP Multimedia Subsystem (IMS) : Merging the Internet and Cellular worlds*. John Wiley and Sons, 2004.
- [9] S. Donovan, The SIP INFO Method. RFC 2976. IETF, Oct 2000.
- [10] J. Rosenberg, The Session Initiation Protocol (SIP) UPDATE Method, RFC 3311. IETF, Sept 2002
- [11] <http://www.stunnel.org/>
- [12] <http://developer.berlios.de/projects/ser/>
- [13] Pankaj Gupta and Nick Mckeown, *Algorithms for packet classification*. IEEE Network, Mar/Apr 2001.

- [14] G. Varghese, *Network Algorithmics: An Interdisciplinary Approach to Designing Fast Networked Devices*. Morgan Kaufman, 2005.
- [15] Gonzalo Navarro and Mathieu Raffinot. *Flexible pattern matching in strings*. Cambridge University Press, 2002.
- [16] A-Ning Du et al. *Comparison of string matching algorithms: an aid to information content security*. Intl Conference on Machine Learning and Cybernetics, 2003.
- [17] M. Ohta, *Overload control in a SIP Signaling Network*. Transactions on Engineering, Computing and Technology V12, March 2006.
- [18] V. Hilt and I. Widjaja. *Hop-by-Hop Overload Control for the Session Initiation Protocol (SIP)*. IETF Internet-Draft, June 2006.
- [19] J. Rosenberg. *Requirements for Management of Overload in the Session Initiation Protocol*. IETF Internet-Draft, October 2006.
- [20] M. Cortes, J. O. Esteban and H. Jun. *Towards Stateless Core: Improving SIP Proxy Scalability*. IEEE Globecom 2006.
- [21] P. Larson. *Dynamic Hash Tables*. Communications of the ACM 31:4, April 1988. 446-457.
- [22] S. Kasera, J. Pinheiro, C. Loader, T. LaPorta, M. Karaul and A. Hari. *Robust Multiclass Signaling Overload Control*. In Proceedings of the 13th IEEE International Conference on Network Protocols (ICNP'05). November 2005. 246-258.