

Dynamic resource reservation for QoS support in IP networks

Autor(en): **Ben-Yacoub, Leila Lamti / Howlett, Claire**

Objektyp: **Article**

Zeitschrift: **Comtec : Informations- und Telekommunikationstechnologie = information and telecommunication technology**

Band (Jahr): **79 (2001)**

Heft 2

PDF erstellt am: **11.07.2024**

Persistenter Link: <https://doi.org/10.5169/seals-876517>

Nutzungsbedingungen

Die ETH-Bibliothek ist Anbieterin der digitalisierten Zeitschriften. Sie besitzt keine Urheberrechte an den Inhalten der Zeitschriften. Die Rechte liegen in der Regel bei den Herausgebern.

Die auf der Plattform e-periodica veröffentlichten Dokumente stehen für nicht-kommerzielle Zwecke in Lehre und Forschung sowie für die private Nutzung frei zur Verfügung. Einzelne Dateien oder Ausdrucke aus diesem Angebot können zusammen mit diesen Nutzungsbedingungen und den korrekten Herkunftsbezeichnungen weitergegeben werden.

Das Veröffentlichen von Bildern in Print- und Online-Publikationen ist nur mit vorheriger Genehmigung der Rechteinhaber erlaubt. Die systematische Speicherung von Teilen des elektronischen Angebots auf anderen Servern bedarf ebenfalls des schriftlichen Einverständnisses der Rechteinhaber.

Haftungsausschluss

Alle Angaben erfolgen ohne Gewähr für Vollständigkeit oder Richtigkeit. Es wird keine Haftung übernommen für Schäden durch die Verwendung von Informationen aus diesem Online-Angebot oder durch das Fehlen von Informationen. Dies gilt auch für Inhalte Dritter, die über dieses Angebot zugänglich sind.

Exploration Programmes:
Corporate Technology Explores Future Telecommunications

Dynamic Resource Reservations for QoS Support in IP Networks

Point-to-cloud SLA's for the provisioning of IP-VPN's with QoS guarantees are helpful for customers since they only have to specify the total traffic volume exchanged between the network and each of their sites with a set of end-to-end QoS guarantees. However, this type of SLA requires the network operator to assume the worst case traffic split between customer sites, where all agreed volume is sent to only one destination. Over-provisioning in this case could result in a huge waste of resources. A solution is proposed enabling the network operator to dynamically reserve capacity in an MPLS backbone by means of traffic prediction (based on measured traffic rates) and engineering techniques. We present results obtained with traffic prediction techniques run on real traffic traces collected from Swisscom voice and data networks.

The Exploration Programme "Customer Care and Service Management Platforms" deals with technologies and processes providing ways for Swisscom to differentiate both in Customer Relationship Management (CRM) and in service management. In customer care, the work is concentrating on multimedia web-based customer contact centre technology and on customer behaviour analysis using data mining techniques. Concerning service management, a platform prototype is built aiming at a radical simplification of 1) provisioning processes – particularly for QoS based SLA (Service Level Agreement) management – and 2) billing processes for IP services.

With its Exploration Programmes, Corporate Technology is exploring telecommunication technologies and new service possibilities with a long-term view of 2–5 years. Further, the expertise built up in the course of this activity enables active support of business innovation projects.

At present, Swisscom offers an IP-VPN service called "ng LAN-I" where only the Committed Access Rate IP (CARIP) is specified. This means that no end-to-end QoS is supported yet. However, the number of cus-

tomers requesting integrated services including QoS guarantees is currently exploding. A simple solution to the problem of offering and fulfilling SLA's for QoS-enabled IP-VPN services in a simple and efficient manner would be to establish IP virtual tunnels between each pair of customer sites. On each virtual tunnel, the needed amount of bandwidth as well as end-to-end QoS guarantees (delay, packet loss and jitter) would be defined. Unfortunately, this solution is not scalable since it requires the Network Operator (NO) to provide a fully meshed network of tunnels for each customer. We have explored an alternative solution called point-to-cloud SLA. In such an SLA, the total traffic volume entering the customer network, called Ingress Committed Rate (ICR), as well as the total traffic volume exiting the customer site, called Egress Committed Rate (ECR), are specified with some end-to-end QoS guarantees.

LEILA LAMTI BEN-YACOUB AND
CLAIRE HOWLETT

tomers requesting integrated services including QoS guarantees is currently exploding. A simple solution to the problem of offering and fulfilling SLA's for QoS-enabled IP-VPN services in a simple and efficient manner would be to estab-

A point-to-cloud SLA presents big advantages to IP-VPN customers: it enables them to send traffic with end-to-end QoS guarantees to all destination sites without the need to specify a complex traffic matrix. The support of this SLA by network operators requires efficient and intelligent resource management mechanisms so as to handle worst case traffic splits. To this end, we propose the use of simple traffic prediction methods coupled with MPLS traffic engineering techniques to support dynamic resource reservations and QoS guarantees.

To support point-to-cloud SLA's, we propose that the NO starts with an over-dimensioned backbone which takes into account the worst case traffic split. Then, future bandwidth requirements are predicted by measuring traffic rates in the network, and IP virtual tunnels are resized accordingly (resizing would be done every 15 min). This solution would lead to better balance load in the network.

Work done at Corporate Technology

Network equipment suppliers such as Cisco or Juniper are developing MPLS traffic engineering techniques, which allow to build IP virtual tunnels according to QoS criteria. We are studying the use of such techniques for the transport of traffic with end-to-end QoS guarantees. The approach is to collect traffic from the Swisscom network, determine the required capacity for each customer site and then assess the traffic split for each destination in the IP-VPN. Collected traffic volumes serve as input to prediction techniques which then allow to dynamically resize MPLS virtual tunnels.

Use of MPLS Traffic Engineering for IP-VPN's

MPLS traffic engineering techniques allow to optimise the routing of IP traffic given the constraints imposed by backbone capacity and topology. Packets are routed through the network according to resources still available, priority versus other flows, bandwidth, and QoS requirements [1]. Different constraint-based routing algorithms can be used. For a given traffic flow, those algorithms find the shortest path that meets the re-

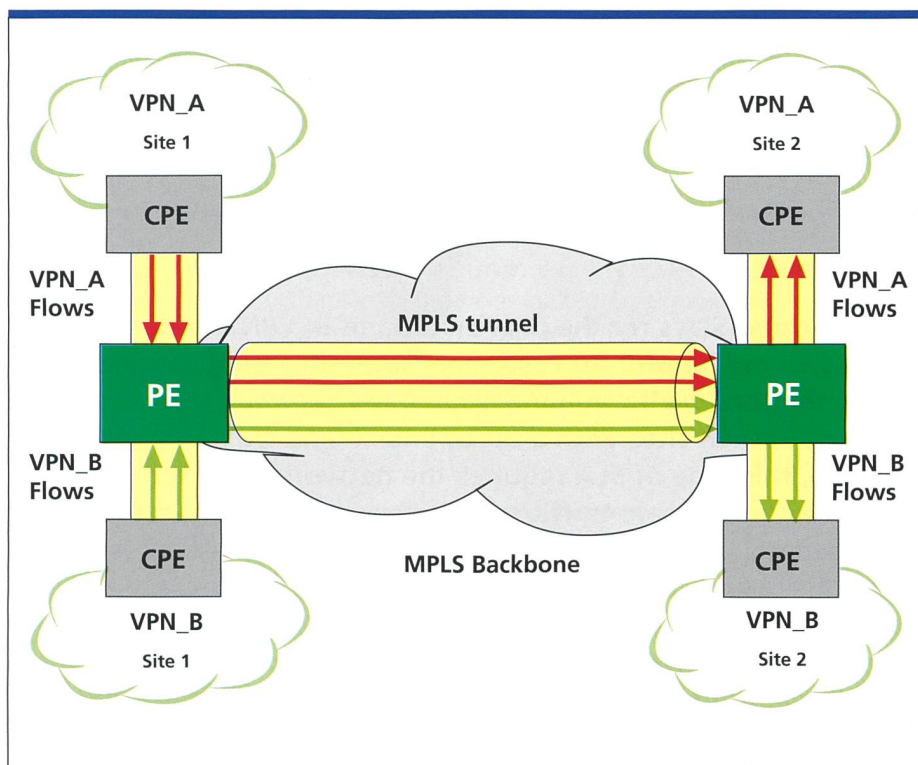


Fig. 1. IP-VPN implementation example in an MPLS backbone.

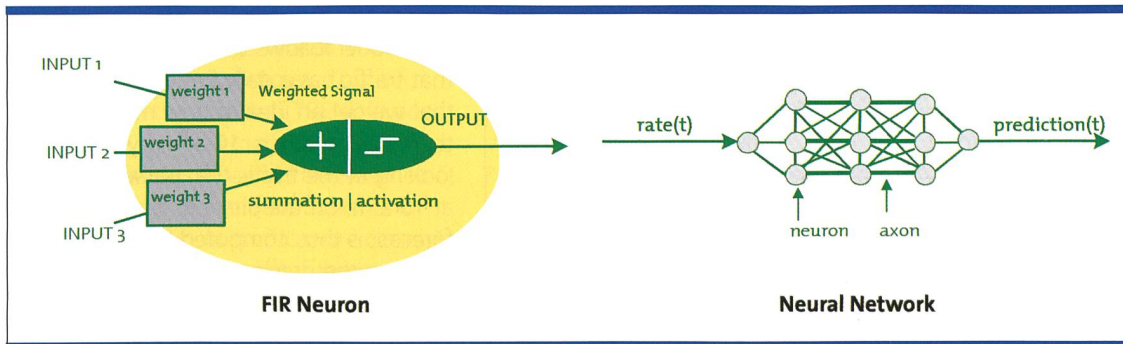


Fig. 2. A Neural Network Architecture.

quirements and then build an MPLS traffic engineered tunnel accordingly. The Differentiated Service framework (DiffServ) can be used to define flow priority in the network (more details on DiffServ can be found in [2]).

The use of MPLS tunnels to support a point-to-cloud SLA is as follows: each Customer Premise Equipment (CPE) is connected via one or more physical links to a Provider Edge (PE) router. A specific forwarding table is stored at each PE router to which members of the VPN are connected. Using these forwarding tables, packets are exchanged first between CPE and PE, then between PE's. In each region, the NO has a Point of Presence (PoP) physically implemented by the PE router. In order to allow interconnection between all the regions, a fully meshed network of MPLS tunnels is set between all PE's (fig. 1).

The advantage of this solution is that it results in a small number of MPLS tunnels to manage. Also, from the customer point of view, the burden of routing configuration, tunnel establishment and maintenance is outsourced to the NO. All flows entering the NO backbone are aggregated at the PE according to their needs in QoS and mapped to DiffServ Classes of Service (CoS). This procedure assumes that each packet is assigned to a CoS at the CPE by configuring the IP precedence field in its header. The PE router creates as many MPLS tunnels as existing CoS values to support service differentiation. A small number of CoS will be supported in the backbone. Therefore, flows sent by customers in the same region are multiplexed on the same MPLS tunnel if they belong to the same CoS and are destined to the same PE (fig.1).

The problem we address in this paper is the difficulty to estimate the needed amount of bandwidth for each MPLS tunnel. As mentioned above, a simple solution would be to provision MPLS tun-

nels with the highest capacity, assuming worst case traffic split situations, which leads to a lot of resource waste in the network. The solution we propose relies on traffic prediction to help dimension MPLS tunnels.

Traffic Prediction Techniques

Traffic prediction techniques belong to the mathematical domain of time series analysis. Their goal is to take as input a series of chronologically sorted values and give as output a forecast of one or several future values of the series. Many techniques, very different in complexity, can be used. Choosing a model in order to get the best results is a difficult task and requires a careful study of the series before trying to compute any forecast. In this article, we briefly describe three models that were tested using real traffic traces captured on Swisscom networks (IPSS with ng LAN-I traces and PSTN for voice traces).

Neural Networks

Neural networks have turned into a widely known term and one may wonder how such a model can be successful in so many different domains such as voice and writing recognition, optimisation problems, data mining, etc. The reason has to do with the fact that neural networks artificially recreate the structure of a human brain. The basic entity is the neuron which assumes two functions: The first, called summation function, weighs all inputs and sums them. The second, called activation function, computes a response according to this sum and sends it as output. A neural network is organised in layers of neurons, each one being linked to its neighbours by nerves (fig. 2).

Building such neural networks requires choosing the number of layers, the number of neurons per layer and the weights. The first two choices are human driven and are made once at the creation

of the network. Setting the weights is done during what is called "the learning phase". Samples of the time series we try to forecast are passed as inputs and weights are changed so that the response is the closest to the one we would like to get. We act like a teacher forcing his pupil to repeat a lesson until he knows it by heart. From this comparison, we get an idea of the issue inherent to neural networks: will it be clever enough to derive a general scheme from the examples it has learnt? The "cleverness" of the network highly depends on the samples that were used during the learning phase, as well as their number, the order and the number of times they were presented to the network. The importance of the learning phase renders neural networks uneasy to handle. A lot of parameters must be set and there is no algorithm leading to the best combination of these parameters: adjusting them altogether must be done via a heuristic approach.

Simple Estimators

At the lowest level on the complexity scale, we find a class of techniques which we call "Simple Estimators" because of their simplicity. Some research studies have tested them on real traffic traces, e.g. on the AT&T backbone for the data and voice traffic [3] and on LAN traces [4]. These references serve as a starting point for our work. Our purpose is to forecast an upper bound for future traffic volumes above the real value rather than predict less and cause under-allocation (i.e. call blocking in the context of voice and congestion in the context of data traffic). Instead of using a complex model to precisely forecast the traffic, we use an estimator such as the mean or the average rate computed over the last time period and then add a correction to avoid under-allocation. The only question about these estimators is how much correction we want to add. If

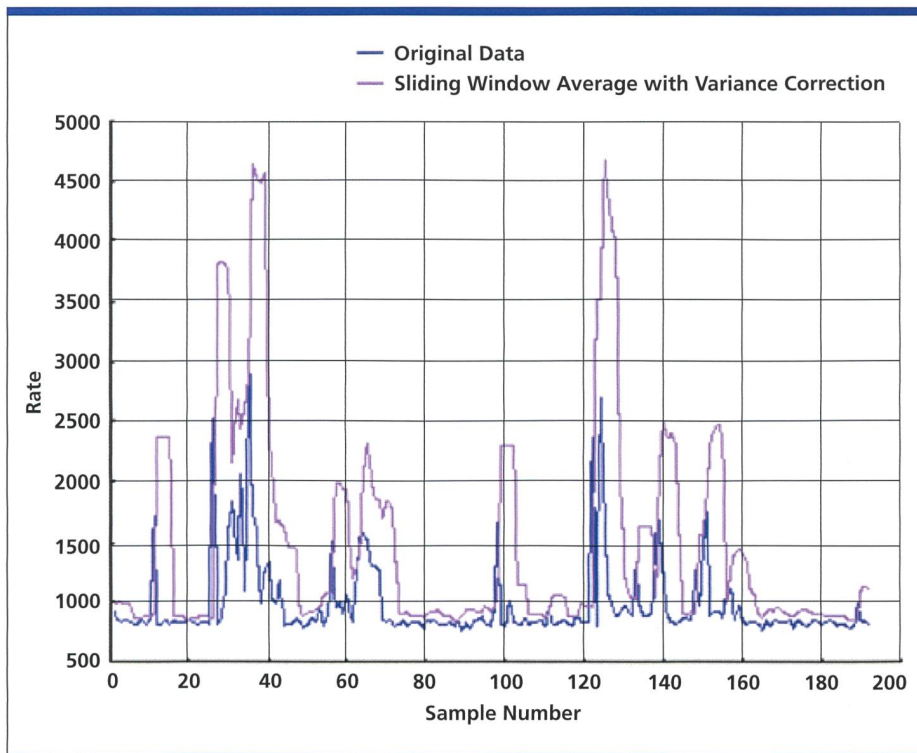


Fig. 3. Prediction Results with Data Traces.

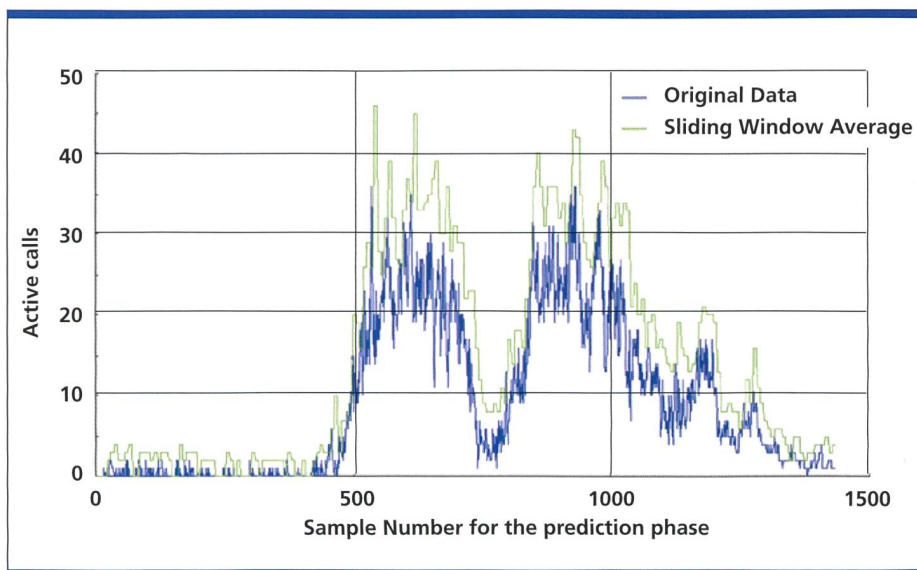


Fig. 4. Prediction Results with Voice Traces.

we add too much, then under-allocation is high and we do not gain much when dynamically resizing the reservations in the network, but predicting less than required is even worse because it results in congestion situations.

Among these estimators, the simplest are based on average and median with a correction based on the variance. Both take the last n traffic rates as input samples. The former computes the average over this set of data, while the latter computes the median. Correction is then

added as the variance over these samples multiplied by a coefficient α that we can modify so as to tune the influence of the correction. Such estimators can be enhanced by dynamically resizing the window over which the prediction is computed. Indeed, a highly variable traffic sees its bandwidth rate change so quickly that samples dating more than 10 minutes are not good for traffic prediction. On the contrary, a stable traffic will allow us to take more samples into account so as to improve accuracy.

Pattern Estimator

This model follows the common belief that traffic has a daily pattern, meaning that we get an idea of how much traffic there will be at e.g. 10 a.m. by simply looking at the traffic that passed through at 10 a.m. on the previous days. The forecast is thus computed by taking the value on the previous day or the average over the last days. Of course, as explained for the "Simple Estimators", such an estimation is not very precise and we have to correct it so as to reduce the risk of under-allocation. In our case, a double correction has been adopted: First, a value specified by the human user is automatically added and has the same value for each estimation. Second, a dynamic correction, which takes into account the behaviour of the real traffic, is added as follows. If under-allocation did occur during the last n minutes, the maximum traffic value is added to the prediction (n represents the period of validity of a forecast, i.e. if we resize the reservations every 15 minutes, n equals 15).

Data Collection and Prediction Results

The prediction tools have been tested over two types of traffic that Swisscom will have to handle over its IP backbone (namely the IPSS network). Data traces were collected from the ng LAN-I service whereas voice traces were computed from Call Data Records (CDR) captured on the PSTN network.

Data Traces

Each CPE using the ng LAN-I service is connected via one or more physical links to a PE router (fig. 1).

Catching traffic volumes was performed by TREND REPORTING®. This tool polls each PE router to retrieve performance statistics. For instance, it retrieves the count number of packets that passed through each interface (an interface corresponds to the link connecting the CPE to the PE), divides it by the number of seconds since the last poll and returns the traffic rate. At the moment, the number of customers for the ng LAN-I service is small which enables to map one VPN to an interface, so we were able to spread traffic traces according to each VPN.

Traffic Variability: Prior to any forecast, the traces were studied. Traffic turned out to be highly variable. During busi-

ness hours, traffic volumes sometimes steep from the average rate of 1000 bit/s to 2700 bit/s over 2 samples. This is partly due to the fact that TREND© does not provide samples more often than every 15 minutes (it would be better to have samples every minute so as to follow the evolution of traffic more closely). We also noticed that even if there was a busy period during office hours and a steady one the rest of the time, traces showed no pattern. For instance, traffic always decreases around mid-day because of lunch break but the phenomenon cannot be quantitatively described. For this reason, the pattern predictor technique was not adequate for data traces.

Results: Neural networks were trained with traces over two days and then tested over traces from the following day. As a first test, all traffic traces were included in the learning phase, and as a second we only considered traces from 8 a.m. to 5 p.m., i.e. the period which is the most difficult to forecast. Both tests resulted in very poor results, the neural network being obviously unable to predict the peaks.

Simple Estimators were tested with varying values of the window over which the prediction is computed and varying values of α . It turned out that the best results were achieved with a small number of samples. This is understandable with regard to the variability of the traces: the shape of the traffic changes so rapidly that old samples are quickly out of date. For the traces we present, the best performance was achieved by choosing a window of length 4 (corresponding to 1 hour). Setting α was more difficult. It is always a trade-off between over and under-allocation. If we cannot afford any under-allocation, then we have to forecast higher values, but in this case the gain, i.e. the amount of bandwidth that is spared when dynamically resizing the MPLS tunnel compared to a static reservation, is weaker. At present, we do not have enough information about the consequences of under-allocation to settle on a definitive value for α . For our tests, values of α were chosen so that under-allocations occurred in less than 10% of cases without causing too much over-allocation. For instance, the *average* estimator – run with variance correction, a window of 4 and $\alpha = 3$ – yields a most frequent error of 100 bits/s, and 13 un-

der-allocations out of 191 samples (7% of all cases). This is a good result compared to an average rate of 1000 bits/s and peaks up to 2700 bits/s. Fig. 3 shows the obtained results.

Performance of our Simple Estimators is quite satisfying. Due to its variability, forecasting data traffic is a difficult task and the behaviour of such simple models exceeded our expectations. All the more since the granularity of the traffic records was 15 min: if other monitoring tools are to be used in the IP backbone, we will be able to achieve even better performance.

Voice Traces

For the moment, Swisscom does not transport voice traffic over its IP network. However, related work is being done in the context of the SURPASS project with an aim to support voice over IPSS. Traffic traces were gathered by means of CDR's from the PSTN traffic which yielded the information <source number, destination number, date, start time, call duration>. The information relevant for our work, <date, hour, number of simultaneous calls> was derived from CDR's using PERL scripts. Unlike data traces, the time granularity of these records goes down to the second but we chose to process call counts every minute, which are the conditions under which we are likely to work later for data collection.

In addition to changing the format of the records from CDR's to traffic traces, they were spread according to the destination region. Note that our goal is to predict traffic on each MPLS tunnel, and therefore we have to consider traffic going on each pair (source region, destination region).

Traffic Variability: Traces captured over a week showed that voice traffic had a recognisable daily pattern. Traffic is very stable until 8 a.m., increases until 11 a.m. when it reaches its maximum, then decreases nearly to its lowest value. From 1 p.m. the volume increases again until it reaches a second maximum (always smaller than the first one) around 4 p.m., then decreases to reach its lowest level at 10 p.m. Such a profile is valid from Monday to Friday, whereas traffic has a different shape during weekends.

Results: For all graphs that we present here, it can be noticed that there is a maximum of 40 simultaneous calls per minute, which may not seem much. This

is due to the fact that CDR's from only one local switch have been gathered. As there are approximately 11 switches per region, we can suppose that there will be about 11 times as much traffic, which gives about 440 simultaneous calls for this example.

The pattern predictor could be tested thanks to the recurring profile of voice traffic. The best results were achieved with a window of 5, which means taking into account the traffic volumes over the last five minutes. Choosing the value of the standard correction raises the same dilemma as when choosing α for simple estimators. It is a trade-off between under-allocation and gain. For instance, considering 4 samples and adding 3 calls as standard correction resulted in 33 under-allocations out of 1440 (less than 3% of cases) and a most frequent error of 3 phone calls. For our tests, we only got traffic records over a week which meant computing a profile over 4 days and using the fifth day for testing the predictor (remember that weekends have a different profile and are therefore left aside). Simple estimators were also tested and the best results achieved with a window of width 7. When choosing $\alpha = 3$, the *average* predictor with variance correction achieved the result of 19 under-allocations out of 1435 (less than 2% of all cases) and a most frequent error of 0,5 phone calls. Fig. 4 presents the obtained results.

Results obtained with these Simple and Pattern Estimators are quite satisfying, showing that traffic prediction does not require complex models. This is an important aspect, especially when envisaging the implementation of a prediction tool: there are so many traffic engineering mechanisms to process in a network that such a tool must be easy to configure as well as low CPU consuming. With such models, the most cumbersome part of the work would be to compute the profile which could be done "off-line", for instance at night when traffic is stable and there is no need for traffic prediction.

Conclusions

The use of MPLS traffic engineered tunnels with end-to-end QoS guarantees is proposed for the provisioning of IP-VPN's with QoS guarantees based on point-to-cloud SLA's. The problem of dynamically estimating the needed amount of bandwidth for each MPLS tunnel can be

solved using traffic prediction techniques. The obtained results underline that a variety of tools exist for time series analysis and no a priori choice should be made. For instance, the simplest models, which we called "Simple Estimators", gave extremely satisfying results even though the traffic – especially for data traces – was highly variable.

Outlook

Since the number of IP-VPN customers is growing and many services are foreseen to be transported over a common IP infrastructure, efficient implementation of traffic prediction coupled with MPLS traffic engineering has to be studied. Towards that end, traffic traces at Provider Edge routers need to be collected and used as input to various traffic prediction techniques. The whole process of data collection is an important step that has to be analysed. Also, the most suitable measurement interval must be studied: short enough to track traffic variability, but long enough to avoid network instability.

These aspects will be investigated in a new Exploration Programme called "IP Business Support Issues". The objective is to propose a global implementation scenario for Swisscom IP network (namely IPSS) to enhance the ng LAN-I service with the support of point-to-cloud SLA's and QoS guarantees. 9.4, 7

Acknowledgements

The authors gratefully acknowledge the support of Franck Wyler and Dominique Moix (Corporate Technology), who provided them with the ng LAN-I traffic traces, using the TREND REPORTING© tool.

Leila Lamti Ben-Yacoub studied computer science at an engineering school in Tunisia from 1990 to 1995 and performed Ph.D. studies in ENST Bretagne France from 1995 to 1999 where she worked as a research assistant. Her Ph.D. work dealt with traffic management and QoS engineering in IP and ATM networks. In autumn 1999, she joined Swisscom Corporate Technology. She is working in the area of service provisioning and performance management for IP networks, with a specific focus on Voice over IP services and MPLS-based Virtual Private Networks.

Claire Howlett is in her last year at Télécom INT, a french grande école specialising in Telecommunications. Her diploma thesis currently in progress at Swisscom Corporate Technology deals with traffic prediction, end-to-end QoS and MPLS.

References

- [1] Cisco Systems, "MPLS Traffic Engineering – Cisco IOS Release 12.0(7)T", 1999.
- [2] V. Bucurescu, L. Lamti, J. Schneider, "Providing QoS-Enabled IP VPN Services", Comtec February 2000.
- [3] N. Duffield, P. Goyal, A. Greenberg, P. Mishra, K. Ramakrishnan, J. Van Der Merwe, "A Flexible Model for Resource Management in Virtual Private Networks", in ACM SIGCOMM'99, Computer Communication Review, Volume 29, No 4, October 1999, pp. 95–108.
- [4] R. Wolski, "Dynamically Forecasting Network Preference Using the Network Weather Service", 1998. SC97 Technical Paper – USCD Computer Science and Engineering Department, 1998.

Abbreviations

CoS	Class of Service
CARIP	Committed Access Rate IP
CDR	Call Data Record
CPE	Customer Premise Equipment
CPU	Central Processing Unit
DiffServ	Differentiated Services
ECR	Egress Committed Rate
ICR	Ingress Committed Rate
IPSS	IP Service Specifications
LAN	Local Area Network
MPLS	Multi-Protocol Label Switching
ng LAN-I	next generation LAN Interconnect
NO	Network Operator
PE	Provider Edge
PoP	Point of Presence
SLA	Service Level Agreement
QoS	Quality of Service
TE	Traffic Engineered
VPN	Virtual Private Network

Zusammenfassung

"Point-to-Cloud" SLA's erlauben es einem Kunden, IP-VPN Dienste mit garantierter Qualität (QoS) zu nutzen und dabei bloss das totale zwischen Netz und Kundenstandorten auszutauschende Verkehrsvolumen mit zugehöriger End-zu-End Qualität zu spezifizieren. Das zwingt den Netzwerkbetreiber zur Dimensionierung auf den schlimmst möglichen Fall, wo das gesamte Volumen an eine einzige Destination geschickt wird. Dies bedeutet aber eine bedeutende Ressourcenverschwendung.

Es wird eine Lösung vorgeschlagen, welche es dem Netzbetreiber erlaubt, die Kapazität in einem MPLS Backbone dynamisch zu reservieren. Die Lösung verwendet Verkehrsvorhersagen aufgrund von gemessenen Verkehrsraten. Die aus realem Verkehr gewonnenen Resultate zeigen, dass mehrere Methoden zum Ziel führen können.



Keep ahead in the Internet race.

Turn your switching investment into gold and gain a competitive advantage in the Internet race. Your customers need access to the staggering amounts of data on the Internet. Alcatel builds voice and data convergence solutions to provide the Internet flow your customers demand. Alcatel offers a unique migration path to the world of advanced multimedia services. Giving you the quality, security and speed to keep ahead in the race. **Alcatel, world leader in Switching and Routing Solutions.**



ARCHITECTS OF AN INTERNET WORLD