

Analysis of high performance software NAT design approaches

Author: Denis Plotnikov

Certified by: _____

Areg Melik-Adamyan, PhD, GNU Toolchain Team Manager, Intel

Vadim Sukhomlinov, Strategic Business Development Manager, Intel

Accepted by: _____

Mats Hanson, Dean of Education, Skoltech

Moscow

June, 2015

Copyright © 2015 Denis Plotnikov. All rights reserved.

Analysis of high performance software NAT design approaches

By

Denis Plotnikov

**Submitted to Skolkovo Institute of Science and Technology on April 24, 2015 in
Partial Fulfillment of the Requirements for the Degree of Master of Science**

ABSTRACT

These days each device tends to communicate with the rest of the world over the Internet. For message passing the stack of TCP/IP protocols is widely used. The protocol performing address resolving is called IP protocol. Today, the IPv4 is the internet communication standard.

There is a problem of address space exhaustion of IPv4. Although, IPv6 is designed to solve that problem and is available for using, transition to it is connected to a number of problems. The big Internet users, such as Telecom Carriers, Internet Service Providers and other enterprises with big data networks are not ready to attack these problems and prefer to solve these problems using Network Address Translators (NATs).

There are different NATs on the market that differs with functionality and performance. NAT, fulfilling the performance and functional requirements of the big Internet users is called CG-NAT.

The CG-NATs available on the market are specialized network devices requiring high level of investments but there are another ways to achieve high performance NAT functionality using a commodity computer.

This work is devoted to exploration of those ways aiming to fulfill the functionality and performance requirements for CG-NAT. The result of the work is experimentally tested approach description of high performance Network Address Translator developing which allows to reduce the CG-NAT price and to keep the performance at the level of specialized network devices.

Thesis advisor:

Areg Melik-Adamyan, PhD, GNU Toolchain Team Manager, Intel

Thesis co-advisor:

Vadim Sukhomlinov, Strategic Business Development Manager, Intel

TABLE OF CONTENTS

Part 1 NAT OVERVIEW	4
1.1 Introduction	4
1.2. Background and motivation.....	5
1.2.1 NAT Purpose and Motivation of Using	5
1.2.2 NAT operation	6
1.2.3 NAT behavioral requirements.....	8
1.2.4 Carrier grade NAT (CG-NAT)	13
1.2.5 NAT Performance metrics	15
1.3 NAT implementations comparison and analysis.....	17
1.3.1 Software NAT	17
1.3.3 Software and hardware NAT comparison.....	18
1.3.4 The comparison and analysis results.....	21
1.4 NAT improvements. The software CG-NAT	22
Part 2. ANALYSIS OF NAT DESIGN APPROACHES	24
2.1 Overview of high performance software design principles	24
2.2 Software NAT design exploration.....	27
2.2.1 NAT design overview	27
2.2.2 NAT bottleneck.....	27
2.2.3 The bottleneck overcoming. Exploration methodology	30
2.2.4 Data structures and algorithms exploration	30
2.3 Implementation	34
Part 3. RESULTS	38
3.1 Evaluation.....	38
3.1.1 Measurement setup.....	38
3.1.2 Experimental methodology	38
3.2 Results discussion.....	39
3.2.1 Results achieved	39
3.2.2 Final design solution	47
3.2.3 The software CG-NAT price estimation	48
3.3 Summary	49
REFERENCES	50
Appendix A.....	53

Part 1 NAT OVERVIEW

1.1 Introduction

These days each device tends to communicate with the rest of the world over the Internet. The internet is a gigantic data network used IP protocols for nodes identification and consists of big number of sub networks. These sub networks have a way of transparent communication between the nodes inside and outside the sub network. The transparent way is called NAT (*Network Address Translator*). The main its function is changing the source IP address and port number of the packet going from the inner network and changing the destination IP address and port number of coming to the inner network packets. This process is called *address translation*. There is a set of requirements for a NAT which stays the same but has some differences depending on the size of network and the use case.

The NAT solves an important problem of IPv4 address space exhaustion. Although, the next generation protocol, IPv6, is ready to use there are some reasons because of which the users of NAT are not ready to move to IPv6 from IPv4 right away [ref_ipchal, ref_depgog].

There are different types of NAT implementations: software and hardware. The software versions of NAT provide all the required functionality but have low performance level. The hardware versions provide full functionality and high performance but usually expensive. Nowadays, when the high performance needed the hardware NATs have no competitors.

This work is focused on finding the way of making a software version of NAT which would have the same performance as a hardware version. This is worth of making efforts because of the high prices of hardware NAT solutions.

The work organized as follows. Part 1 gives overview of NAT technology and the existing NAT solutions, formulates the requirements to the NAT in terms of functionality and the performance. Part 2 does the software design approach analysis and makes some assumptions about what means should be used in order to achieve desired results. Part 3 analyzes the results acquired from implementation of the assumptions described in part 2, makes a conclusion of what the best approach to use among described is and formulates the conclusion.

1.2. Background and motivation

1.2.1 NAT Purpose and Motivation of Using

NAT was invented as a way of using a single IP address for several network devices. Thus, the main reason of using it is to reduce the number of IP addresses used by a number of network devices. Besides of the main functionality the NAT gives some security benefits like internal network structure hiding while and ability to restrict the access to the outer network of an internal node. Some versions of NAT provide functionality of translating IPv4 address to IPv6 which is helpful when changing provider settings.

This time the number of network devices grows rapidly because personal devices with the Internet access become more popular and more people starts using it every day. Each of those devices needs a network address to communicate through the data network. Presently, the main protocol used in the Internet is IPv4. The problem with IPv4 is that at the time the number of address in the IPv4 has reached its limit and the organization affiliated to manage the addresses issues has stopped its free distribution at 2012 ^[ref_ripe_limit].

As the number of IPv4 addresses has reached its limit the addresses are turning into more and more valuable resource and the price of buying or renting it

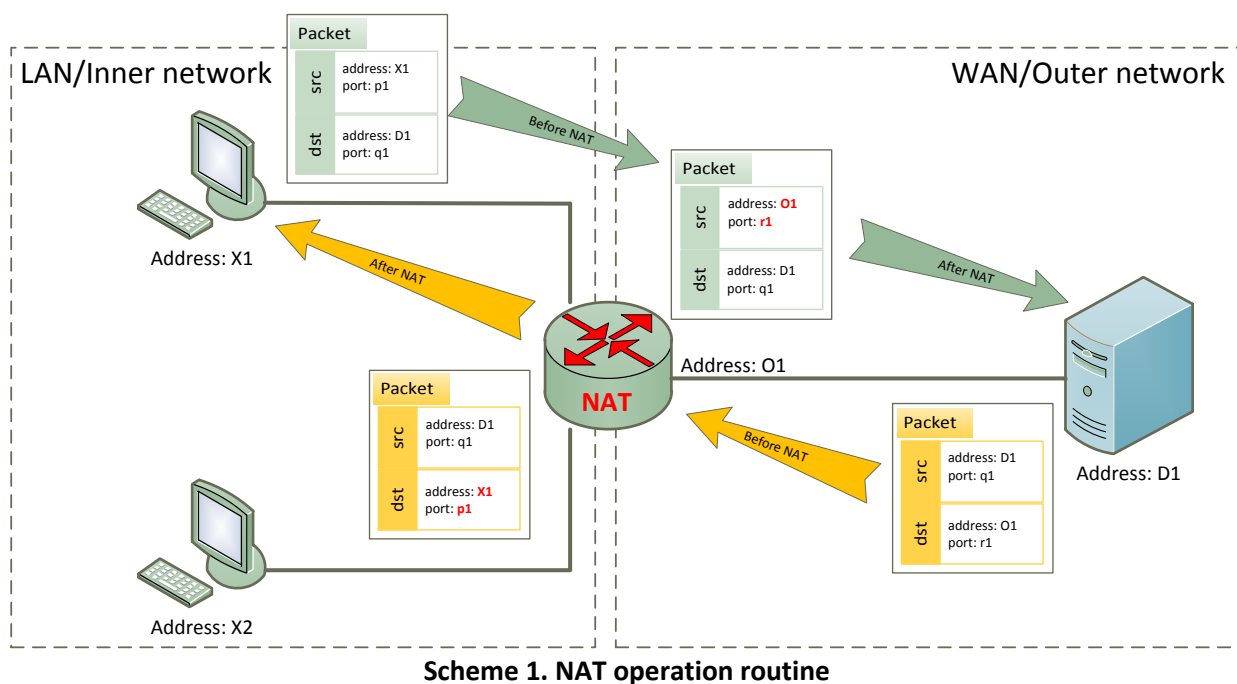
becomes higher and higher which means that more and more network devices will share the one IPv4 address for the Internet accessing ^[rfc_6888].

Although, IPv6 seems to be a solution of the lack of IPv4 addresses the process of switching to IPv6 is quite slow. The most likely reason for that is the whole network setting changing necessity as well as a need to remove all legacy network equipment and software which are not compatible with IPv6. This will take a lot of efforts from the Internet providers and they are not ready to put many of them right away because of the investments needed. Instead, they keep working on IPv4 using NATs to reduce the number of “white” IPv4 addresses used.

1.2.2 NAT operation

This document is focused in exploring of traditional NAT setup as most frequently used. Although, other NAT setups also have its own application case but they are quite rare in the real world and omitted from the consideration.

The traditional NAT setup ^[ref_rfc3022] looks as follows. There are two networks which considered by NAT as *inner* and *outer*. The inner network is the local area network served with NAT. The outer network is the wide area network which can communicate with inner network via the NAT only. There are two main NAT methods Basic NAT and NAPT (Network Address Port Translator). Basic NAT performs IP address translation only (i.e. changing a source IP address of a packet to an IP address allocated to NAT for translation) without changing of port number when NAPT does translation of a tuple {IP, port (TCP/UDP/ICMP)}. These days NAPT method is mostly used and furthers in this document it is implied when saying NAT.



Scheme 1. NAT operation routine

NAT operations as follows: a node from inner network sends a packet to a node in the outer network. The packet comes to the NAT. The NAT gets the packet, allocates an IP address and port number for translation. Then, the source IP address and port number stores in the NAT with respect to just allocated IP and port number to perform backwards translation. The packet source IP address and port number is replaced with the allocated tuple and checksums of IP and TCP/UDP/ICMP headers of the packet are recalculated. After that, the packet is sent to the destination node in the outer network. The outer network node receives the packet and sends a replying packet using the tuple of source IP address and port number from the just received packet as the destination IP address and port number in its packet. The replying packet comes to the NAT. The NAT using the destination address and port number from the replying packet looks for the corresponding tuple of IP address and port number saved previously. Having found the tuple, the NAT performs destination IP address and port number replacement in the replying packet as well as changing of checksums in IP and TCP/UDP/ICMP packets. Then the replying packet is sent to the node in the inner network which was an originator of connection.

Although, NAT supplies transparency to inner and outer nodes communication (i.e. the nodes know nothing about NAT existing) it has a serious disadvantage. There are a set of application protocols (FTP, DNS, PPTP, H.323, etc.) working onto TCP/UDP that store the connection data (i.e. IP address and port number) inside their packet on the level upper than L4 of ISO model where TCP/UDP works. This fact leads to inconsistency of IP addresses and port numbers in translated by NAT packets where source IP and port number in TCP/IP headers doesn't match to the source IP and port number in the upper level protocol headers. This problem is solved with ALGs (Application Layer Gateway) working with NAT. There are different protocol specific ALG. Each ALG protocol are able to distinguish and process the packets belongs to the protocol accordingly. This document is not focused on using ALGs in NAT and only core traditional NAT functionality is taken into consideration.

1.2.3 NAT behavioral requirements

NAT behavioral requirements for TCP, UDP and ICMP are clearly stated in [ref_rfc5382], [ref_rfc4787] and [ref_rfc5508] respectively. Here the generalized list of NAT behavioral requirements based on mentioned documents is shown.

1. A NAT must have an "Endpoint-Independent mapping" behavior. Endpoint-Independent mapping means that the NAT reuses port mapping for subsequent packets sent from the same internal IP address and port to any external IP address and port_[ref_rfc4787]
2. A NAT must support all valid sequences of TCP/UDP/ICMP packets for connections initiated both internally as well as externally when the connection is permitted by the NAT. In addition to handling the TCP 3-way

handshake mode of connection initiation, A NAT must handle the TCP simultaneous-open mode of connection initiation

3. If application transparency is most important, it is RECOMMENDED that a NAT have an "Endpoint-Independent Filtering" behavior for TCP. If a more stringent filtering behavior is most important, it is RECOMMENDED that a NAT have an "Address-Dependent Filtering" behavior.

Endpoint-Independent Filtering means that sending packets from the internal side of the NAT to any external IP address is sufficient to allow any packets back to the internal endpoint.

Address-Dependent Filtering means that for receiving packets from a specific external endpoint, it is necessary for the internal endpoint to send packets **first** to that specific external endpoint's IP address

4. A NAT must not respond to an unsolicited inbound TCP SYN packet for at least 6 seconds after the packet is received. If during this interval the NAT receives and translates an outbound TCP SYN for the connection the NAT must silently drop the original unsolicited inbound TCP SYN packet. Otherwise, the NAT should send an ICMP Port Unreachable error (Type 3, Code 3) for the original TCP SYN. The NAT must silently drop the original TCP SYN packet if sending a response violates the security policy of the NAT
5. If a NAT cannot determine whether the endpoints of a TCP connection are active, it MAY abandon the session if it has been idle for some time. In such cases, the value of the "established connection idle-timeout" must not be less than 2 hours 4 minutes. The value of the "transitory connection idle-timeout" must not be less than 4 minutes. The value of the NAT idle-timeouts may be configurable.

6. A NAT UDP mapping timer must not expire in less than two minutes, unless the port number is from the well-known port range of 0 -1023. In that case NAT may have shorter UDP mapping timers. The value of the NAT UDP mapping timer MAY be configurable. A default value of five minutes or more for the NAT UDP mapping timer is recommended.
7. A NAT must not have a "Port assignment" behavior of "Port overloading". Another words, the NAT must not always use port preservation even in the case of port collisions
8. It is recommended that a NAT have an "IP address pooling" behavior of "Paired". It means that the NAT use the same external IP address mapping for all sessions associated with the same internal IP address. This requirement is not applicable to NATs that do not support IP address pooling
9. It is recommended that a NAT have a "Port parity preservation" behavior which means that after the NAT processing an even UDP port will be mapped to an even UDP port, and an odd UDP port will be mapped to an odd UDP port.
- 10.If a NAT includes ALGs it is RECOMMENDED that all of those ALGs (except for FTP) be disabled by default.
11. A NAT must support "hairpinning" for TCP/UDP/ICMP. A NAT's hairpinning behavior must be of type "External source IP address and port". This means that two nodes are behind the same NAT both use for communication to each other different tuples of external IP address and port allocated by NAT
- 12.Unless explicitly overridden by local policy, a NAT device must permit ICMP Queries and their associated responses, when the Query is initiated from a private host to the external hosts

13. An ICMP Query session timer must not expire in less than 60 seconds. It is recommended that the ICMP Query session timer be made configurable
14. When an ICMP Error packet is received, if the ICMP checksum fails to validate, the NAT should silently drop the ICMP Error packet. If the ICMP checksum is valid, do the following: If the IP checksum of the embedded packet fails to validate, the NAT should silently drop the Error packet; If the embedded packet includes IP options, the NAT device must traverse past the IP options to locate the start of the transport header for the embedded packet; the NAT device should not validate the transport checksum of the embedded packet within an ICMP Error message, even when it is possible to do so; if the ICMP Error payload contains ICMP extensions, the NAT device must exclude the optional zero-padding and the ICMP extensions when evaluating transport checksum for the embedded packet.
15. If a NAT device receives an ICMP Error packet from an external realm, and the NAT device does not have an active mapping for the embedded payload, the NAT should silently drop the ICMP Error packet. If the NAT has active mapping for the embedded payload, then the NAT must do the following prior to forwarding the packet, unless explicitly overridden by local policy: revert the IP and transport headers of the embedded IP packet to their original form, using the matching mapping; leave the ICMP Error type and code unchanged; modify the destination IP address of the outer IP header to be same as the source IP address of the embedded packet after translation.
16. If a NAT device receives an ICMP Error packet from the private realm, and the NAT does not have an active mapping for the embedded payload, the NAT should silently drop the ICMP Error packet. If the NAT has active mapping for the embedded payload, then the NAT must do the following

- prior to forwarding the packet, unless explicitly overridden by local policy:
revert the IP and transport headers of the embedded IP packet to their original form, using the matching mapping;
17. leave the ICMP Error type and code unchanged; if the NAT enforces Basic NAT function, and the NAT has active mapping for the IP address that sent the ICMP Error, translate the source IP address of the ICMP Error packet with the public IP address in the mapping. In all other cases, translate the source IP address of the ICMP Error packet with its own public IP address
 18. While processing an ICMP Error packet pertaining to an ICMP Query or Query response message, a NAT device must not refresh or delete the NAT Session that pertains to the embedded payload within the ICMP Error packet.
 19. When a NAT device is unable to establish a NAT Session for a new transport-layer (TCP, UDP, ICMP, etc.) flow due to resource constraints or administrative restrictions, the NAT device should send an ICMP destination unreachable message, with a code of 13 (Communication administratively prohibited) to the sender, and drop the original packet.
 20. A NAT device MAY implement a policy control that prevents ICMP messages being generated toward certain interface(s).
 21. Receipt of any sort of ICMP message MUST NOT terminate the NAT mapping or TCP connection for which the ICMP was generated.

The requirements listed above do not guarantee the compliance of NAT with all application protocols but fulfilling them significantly improves the likelihood of successful processing of any kinds of packets.

1.2.4 Carrier grade NAT (CG-NAT)

The NAT supplies the ability to share a single external IP address among several nodes in the external network. Some of internet service providers have started offering this service long before IPv4 address space shortage problem has arisen showing that there is another driven force of using NAT. Each subscriber at the ISP's network assigned a private address and the NAT, situated at the customer edge, translates traffic between public and private addresses^[ref_frc6888].

Because of scales of ISP's networks NATs used there have some additional functional and determined performance requirements. The list of NAT functional requirements is added with following points (taken from ^[ref_frc6888]) extending the points shown previously in this chapter:

1. The CGN function should not have any limitations on the size or the contiguity of the external address pool. In particular, the CGN function must be configurable with contiguous or non-contiguous external IPv4 address ranges
2. A CGN MUST support limiting the number of external ports (or, equivalently, "identifiers" for ICMP) that are assigned per subscriber. Per-subscriber limits must be configurable by the CGN administrator. Per-subscriber limits may be configurable independently per transport protocol. Additionally, it is recommended that the CGN include administrator-adjustable thresholds to prevent a single subscriber from consuming excessive CPU resources from the CGN (e.g., rate-limit the subscriber's creation of new mappings).
3. A CGN should support limiting the amount of state memory allocated per mapping and per subscriber. This may include limiting the number of sessions, the number of filters, etc., depending on the NAT implementation. Limits should be configurable by the CGN administrator.

Additionally, it should be possible to limit the rate at which memory-consuming state elements are allocated.

4. It must be possible to administratively turn off translation for specific destination addresses and/or ports.
5. Once an external port is deallocate, it should not be reallocated to a new mapping until at least 120 seconds have passed, with the exceptions being If the CGN tracks TCP sessions TCP ports MAY be reused immediately. If external ports are statically assigned to internal addresses, the assignment remains constant across state loss, than ports may be reused immediately. If the allocated external ports used address-dependent or address-and-port-dependent filtering before state loss, they may be reused immediately. The length of time and the maximum number of ports in this state must be configurable by the CGN administrator.
6. A CGN must implement a protocol giving subscribers explicit control over NAT mappings. That protocol SHOULD be the Port Control Protocol [RFC6887]
7. CGN implementers should make their equipment manageable. Standards-based management using standards such as "Definitions of Managed Objects for NAT" [RFC4008] is recommended
8. When a CGN is unable to create a dynamic mapping due to resource constraints or administrative restrictions (i.e., quotas): it must drop the original packet; it should send an ICMP Destination Unreachable message with code 1 (Host Unreachable) to the sender; it should send a notification (e.g., SNMP trap) towards a management system (if configured to do so); it must not delete existing mappings in order to "make room" for the new one. (This only applies to normal CGN behavior, not to manual operator intervention.)

The requirements listed above do not guarantee the compliance of NAT with all application protocols and is based on the best practices. Fulfilling the requirements significantly improves the likelihood of successful processing of any kinds of packets while NAT working on the ISP edge.

Another important point is CGN performance. Although, there is nothing about it stated in the requirements, it worth a lot because of the nature of the network serving with CG-NAT. As CG-NAT is an ISP's appliance and ISP's network consist of tens of thousands of nodes the CG-NAT should provide sufficient performance level to satisfy the customer demands. In case of CG-NAT the performance seems to be a more critical characteristic than strict compliance with a standard. From the marketing point of view the performance could be more important than the fulfilling the standard functionality because no one would buy a system that fulfilled the functionality requirements and didn't perform them with proper pace and vice versa, if the performance are enough than one could possibly turn a blind eye to some lack of the features.

1.2.5 NAT Performance metrics

The goal of this work is to develop a working prototype of software defined carrier-grade network address translator (SD CG-NAT). To make sure that the SD CG-NAT is close to reality in terms of performance it is necessary to define the performance metrics and set their values. In order to get those metrics, a couple of sources are used. The first one is Rostelecom technical requirements for CG-NAT [ref_TT_ROS_TEL]. The second one is the performance specification claimed by one of the on-market available NAT device producers[ref_RDP.RU] which employ the same approach as this research does: **to use not task specific computer (a commodity server) to make a network specific solution using a mix of algorithmic and technological approaches. (our_approach)**

NAT performance metrics:

- **Packets processing rate** – (packets per second [PPS]) – the router's maximum rate of packet processing. This is the main metric describing the packet processing abilities of a NAT device.
- **Concurrent session support** – (number) – the maximum number of sessions produced by served network. It describes the maximum network size which can be served by the NAT device. As described later in this document than bigger the network than harder to maintain translations to its nodes.
- **Connections setups rate** – (connection setups per second [csps]) – the number of new NAT records to be created in a second. This metric shows the NAT ability to create new NAT records and could be a drawback of the NAT device in a certain modes of network work like when the networks nodes start creating of new connections actively, for example in the beginning working hours
- **Throughput** – (bit per second [bps]) – it isn't very clear metric of the NAT device because it is mostly defined with NIC (network interface card) performance used by NAT device. If the NAT device won't have enough of packet processing rate its throughput can't achieve the maximum throughput provided with NIC and vice versa. The main sense of having it in the metric list is to make sure that NAT device is able to transfer needed amount of information.
- **Latency** – seconds [sec] – time needed for one packet processing. This metric is important when evaluating the minimal time frame of one packet processing to know what part of runtime is needed for changing the packet data. This can be helpful when comparing performance growth and scalability.

This set of characteristics is usually used by equipment vendors while describing their competitive advantages. Thus, using it will allow one to be on the same page with all the professionals working in the field of computer network devices.

In this document for evaluation of the performance another characteristic is used:

Cycles per packet – [cps] – the amount of processors' cycles spent on processing of one packet. This characteristic seems to be more descriptive than others while describing the NAT performance because there are a lot different processors which differ to each other with CPU frequency and technologies used which makes it harder to compare the performance of the NAT on different processors using the set of metrics described earlier in this chapter. Cycles per packet characteristic gives clearer impression of the performance because at least it doesn't strongly depend on CPU frequency however there are other limiting factors influencing on the characteristic value such as system bus frequency and memory frequency. Another drawback is that this characteristic becomes quite confusing when trying to describe the performance on multiple cores. Thus, the main performance characteristic used in this work for assessing the performance is Cycles Per Packet and is used mainly for choosing the best working approach. The target metrics values are set in the following paragraph and are used as the requirements to the NAT settings and abilities.

1.3 NAT implementations comparison and analysis

There are two types of NAT: software and hardware.

1.3.1 Software NAT

Software NAT is a program that implements NAT functionality and works on the top of operating system. There are two kinds of software NAT implementation.

One is a user program working in the user mode totally like *NAT32 IP Router* [ref_nat32]

or *WinGate* for windows or *IPFilter* for Linux and using the OS resources to access the network. It can be installed or uninstalled by user's demands. There are commercial versions of this kind of NAT as well as free once.

Another kind of software NAT is a program working in the kernel space like Linux module or is a part of operation system like *ICS* in Windows or *iptables* in Linux. As this kind of NAT distributed along the operation systems it is usually free. Some of UNIX family software NATs have opened source code and are free for changing.

1.3.2 Hardware NAT

Hardware NAT is a specialized network device or a feature of specialized network device like firewall or router. These devices usually have their own specific operation system and interfaces for management. The hardware core of hardware NAT device consist of a specific processor designed for fast packet processing as well as associative memory and another chips that increase the performance in specific operations. Hardware NATs are produced by telecommunication producers like Cisco, Juniper and others. These devices have high performance and are expensive.

1.3.3 Software and hardware NAT comparison

Functionality: As hardware NAT is usually a part of industrial firewall or a router which is used at Internet Services Provider's or Data Center's facilities it is equipped with vast functionality where NAT is one of many. The functionality includes firewall, VPN support and crypto security features allowing a customer to have many of needed abilities in one box. Software NAT is often is an integral part of operation system and its functionality is not that advanced because it is usually

used for small offices or tiny private networks. Talking about NAT core functionality, both of them provide full support of NAT requirements.

Performance: Hardware NAT is specially designed for high performance used the cutting edge hardware for achieving it. In particular, for the performance improvement special memory units are used called CAM (Content-addressable memory a.k.a. associative memory). CAM(TCAM) is extremely fast in tasks of comparing input data against stored data and returning matched data as a result. Although CAM is fast it has low memory capacity and very high price. One unit with capacity of 80Kbytes costs around \$180 while hundreds of these units needed in a router working at an ISP's rack.

Software NAT is installed in commodity computers and has a serious performance limiting factor. As software NAT use OS system as the source of network resources it uses OS system calls to get them and, hence, is limited by the performance of that system calls which means that the NAT is not able to outperform the OS it uses. Thus, the main limiting factor is the network stack used by the operation system. The experiment revealed the packet performance rate for Linux 3.16 kernel around 260 Kpps.

Upgradability: The hardware NAT is a set of hardware mounted in some kind of chassis with a piece of specialized software pre-installed which is hardly be compatible with different set of hardware because of technical and vendor limitations. Another thing is that because of the software, controlling the hardware NAT is proprietary and upgrading might be an issue. This makes the only way of upgrading this equipment possible: buying a new software update, support plan, new hardware NAT device from the vendor. Unlike hardware NAT, software NAT can be installed at any system supporting the OS the software NAT specialized. Thus, increasing of performance is possible by updating the hardware where the software NAT installed. As software NAT is just a program than the upgrading/updating could be done by using the same approach as other pieces of

software used with paying no attention to the hardware used. This makes the software NAT more flexible in terms of modification and customization.

Price: The price of the system implementing NAT functions based on software NAT consists of several parts: price of the computer, price of the OS and the price of software NAT itself. The price of software NAT program was around \$2200 (Win gate_[ref_wingate]) for maximum functional version at the time of writing. There are free versions of software NAT working on free OS like iptables on Linux. In that case the price depends on the price of computer used only.

As for hardware NAT the prices are much higher. This is mainly due to pieces of hardware used in producing of such devices but the vendor interests also play a big role. The marked research has been done to evaluate the prices and the results are shown in Figure 1.

Performance/price

edge routers with NAT

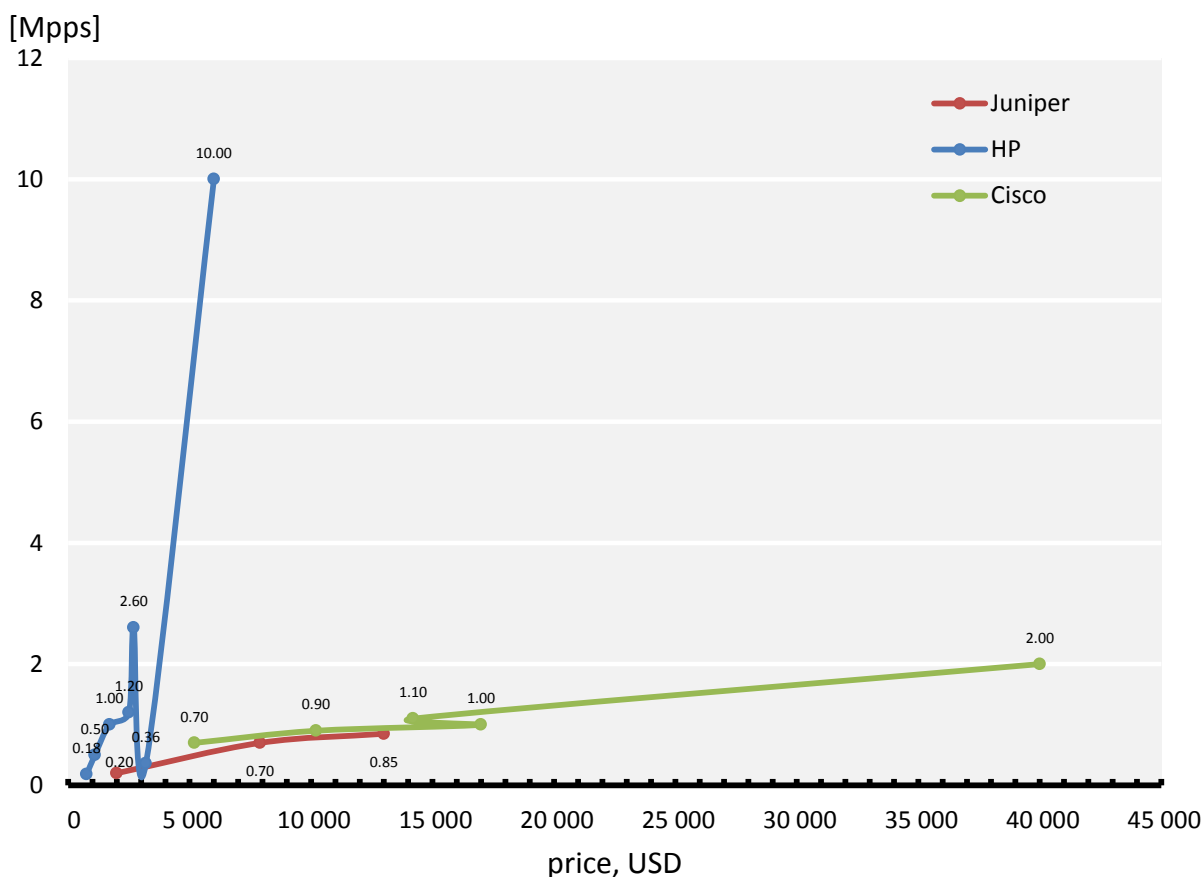


Figure 1. Performance/price relation plot of the modern edge routers supporting NAT

For evaluation a set of edge routers was taken produced by three well known telecommunication equipment vendors: HP, Cisco and Juniper. The performance of the NAT mode has not been found and firewall on mode was used as an approximation. The packet processing rate and related price for the set of HP devices seems to be not consistent with the same data for Cisco and Juniper devices. This is could be a consequence of difference in the packet processing rate definition and measuring methodology used by different vendors. When claiming performance, Juniper and Cisco use packet processing rate with firewall turned on. HP doesn't give any clarification about the mode they used for packet processing rate measuring. The data found for Juniper and Cisco routers shows pretty much the same trend and, thus, it can be considered as trustworthy and showing the current situation on the market.

It is seen from the chart that high performance costs a lot: roughly 15 000 USD for each 1 Mpps which makes buying of the router the matter of capital investments for ISPs.

The data for the chart was collected from the public resources on the Internet. The prices are relevant for Russian market, for other countries prices may vary. The list of models used is in Appendix A.

1.3.4 The comparison and analysis results

The results of comparison are shown in Table 1. The analysis revealed the situation as following. Hardware NAT provides vast functionality and high performance wherein the price is much higher than software NAT provides. The reason for the lower software NAT performance is limited packet processing ability of underlying OS. Also, Hardware NAT is hard to upgrade because of the high cost whereas Software NAT is easily updatable because it is just a program working on top of operation system.

	Hardware NAT	Software NAT
Additional functionality	High	Low
Performance	High	Low
Upgradeability	Low	High
Price	High	Low

Table 1. Advantages and disadvantages comparison of hardware and software NATs. The advantages marked with red color

1.4 NAT improvements. The software CG-NAT

How to improve NAT: Based on the analysis, the reasonable improvement would be increasing of performance of the Software NAT. In this case it could significantly increase the number of cases where Software NAT could be used instead of Hardware NAT. It is implied that software NAT would be cheaper than hardware NAT but this is a matter of exploration. If the software NAT was cheaper than the hardware one it would reduce the investments on building the network. Having these properties, the software NAT could be applied at ISP's facilities playing the role of a CG-NAT. In that case the software NAT should have a certain values of performance metrics.

This work is focused on exploring feasibility and reasonability of software NAT development that could be compared with hardware NAT in performance and could be used as a CG-NAT. Further in this work the name "the software CG-NAT" is used to denote the goal of exploration.

The performance metrics target values: As it was mentioned earlier the software CG-NAT should have the target metrics values. Based on the sources of information [ref_TT_ROS_TEL, ref_rdp.ru] it is reasonable to set the performance requirements of the software CG-NAT device to the following values:

- Packet processing rate: ≥ 5.5 Mpps
- Concurrent session support: ≥ 65.5 M (a B-class network with up to 1000 active ports for each node)
- Connection setups rate: ≥ 3 Mcsps
- Throughput: 10 Gbps
- Latency: ≤ 180 ns

These values are relevant to performance demands of the ISP's which are the users of CG-NAT.

Part 2. ANALYSIS OF NAT DESIGN APPROACHES

2.1 Overview of high performance software design principles

The essence of NAT is packet processing. To increase the ability of NAT to process the packets certain kinds of software design patterns should be employed. There are a number of works and books related to software design for fast packet processing] and underlying in chapter 5 of describing the technics capable to increase the packet processing performance [ref_ix, ref_uppc, ref_epssr, ref_click] and to boost overall system performance [chapter 5,6,7 of ref_cspp]. The list of these principles is following.

Zero-copy policy: Each data copying spends hundreds of processor cycles to perform the memory allocation for new copy of coping itself. To eliminate this unnecessary operation the zero-copy policy should be employed. Zero-coping means that while processing a packet the software doesn't make any copies of the data and for data manipulation the pointers are used. Using this technic allows making only one copy of data for the whole packet processing round when copying the data to transmission queue which significantly reduce the amount of possible memory allocations.

Cache optimized data access: Modern processors use multilevel cache memory based on fast but expensive SRAM. The cache represents a table consist of cache lines where the data stored. The cache line stores the chunk of continuous memory. To get the access to some memory address the processor, firstly, check if this address is currently in the cache and if not, than load this chunk of memory from the main memory. This situation is called a cache miss. If the processor finds the chunk of memory in its cache it is called a cache hit. A cache miss costs a processor around a hundred of cycles which increase overheads of runtime. To decrease the cache miss rate using principle of locality should be considered.

There are two types of localities. Due to the cache line nature of memory loading for reducing the reads from the main memory the data in memory have to be stored in continuous manner. This allows storing the data, to be used, in the cache in advance and then use it without additional readings. This is called spatial locality. Some data is needed to be used frequently during the runtime. This is called temporal locality. This kind of data has to stay in the cache but can be discarded from the cache because of a program's memory access pattern. Due to that while the program development this specialty should be taken into the account and the memory pattern access should be adopted accordingly.

Batching: It is a method that allows eliminating possible delays for reading packets from a network interface card as well as employs cache optimized data access by processing the packets in blocks. Processing the block of data gives increasing of data and instruction locality which reduces the overheads in comparison with sequential packet processing approach.

Parallel processing: Processing the data in parallel manner makes use of multiple cores available on modern processors. Theoretically, splitting the works into several threads gives maximum boost of performance directly proportional to the number of threads used. In practice, this boost is not achievable as stated in Amdahl's law [ref_amdahl]. Therefore, to benefit from the parallel technics the program should use designed so, that the sequential portion would be as small as possible. Another pitfall of parallel programming hides in the processor's hardware design. If some threads read from and write into the same shared data structure and the data structure is not adopted to multithreading use it is highly possible that the performance of the program will suffer from the consequences related to the cache coherence. The problem is that a thread can write to some value which is stored in some memory address which, in turn, is stored in cache of several other cores because of cache line properties. This write induces the other cores, having this address in their caches, to update the address's value. This operation takes

some time and is an unnecessary overhead when the value is not relevant to the algorithm essence. This property of the multi-core processors cache should be taken into consideration when using multi-thread programming technics.

Avoiding OS's networking facility: ... don't use Linux/Windows TCP/IP stack

Operation Systems implement a large number of network protocols providing vast functionality related to the networking. These network facilities are OS specific and optimized to be used by the operation system demands. The disadvantage is that this functionality is controlled by OS, accessed via system calls and is not designed and, thus, destined for using in high performance network applications. Instead of using the OS-supplied network functionality, alternative ways have to be considered like using a high performance network stack or a Data Plane framework.

Interruption avoiding: Both types of interruptions: software and hardware incur OS context switching and switches the processor from user mode to kernel mode where the interruption handler works which is fraught with doing additional work by processor increasing overheads and, therefore, program runtime. To avoid these overheads from the software interruptions point of view the system functions should be used with great carefulness or even be discarded if possible. From the hardware interruptions point of view the special network card drivers should be used which does not induce hardware interruptions.

2.2 Software NAT design exploration

2.2.1 NAT design overview

The most interesting part in the NAT system is the algorithm and data structure for storing the address translation information. There was a number of work related to this issue [ref_natimp, ref_ehms, ref_ecrc] Although, some technics described in these works can be employed but most of them are related to special network processors and equipment. In this work the implementation of NAT using a commodity computer is the matter of exploration. The commodity computer uses processors with the architecture different from the network ones. In the most cases it is x86 architecture which requires certain ways of high performance achieving.

The Network Address Translation process consists of several parts: packet receiving, packet translating (making a decision about packet changing, changing the packet headers) and sending the packet out.

Packet receiving and sending are related to the OS and hardware environment while the packet translating is related to internal program organization. To choose the design which would fulfill performance requirements it is necessary to determine the part of the system where performance reduction can occur.

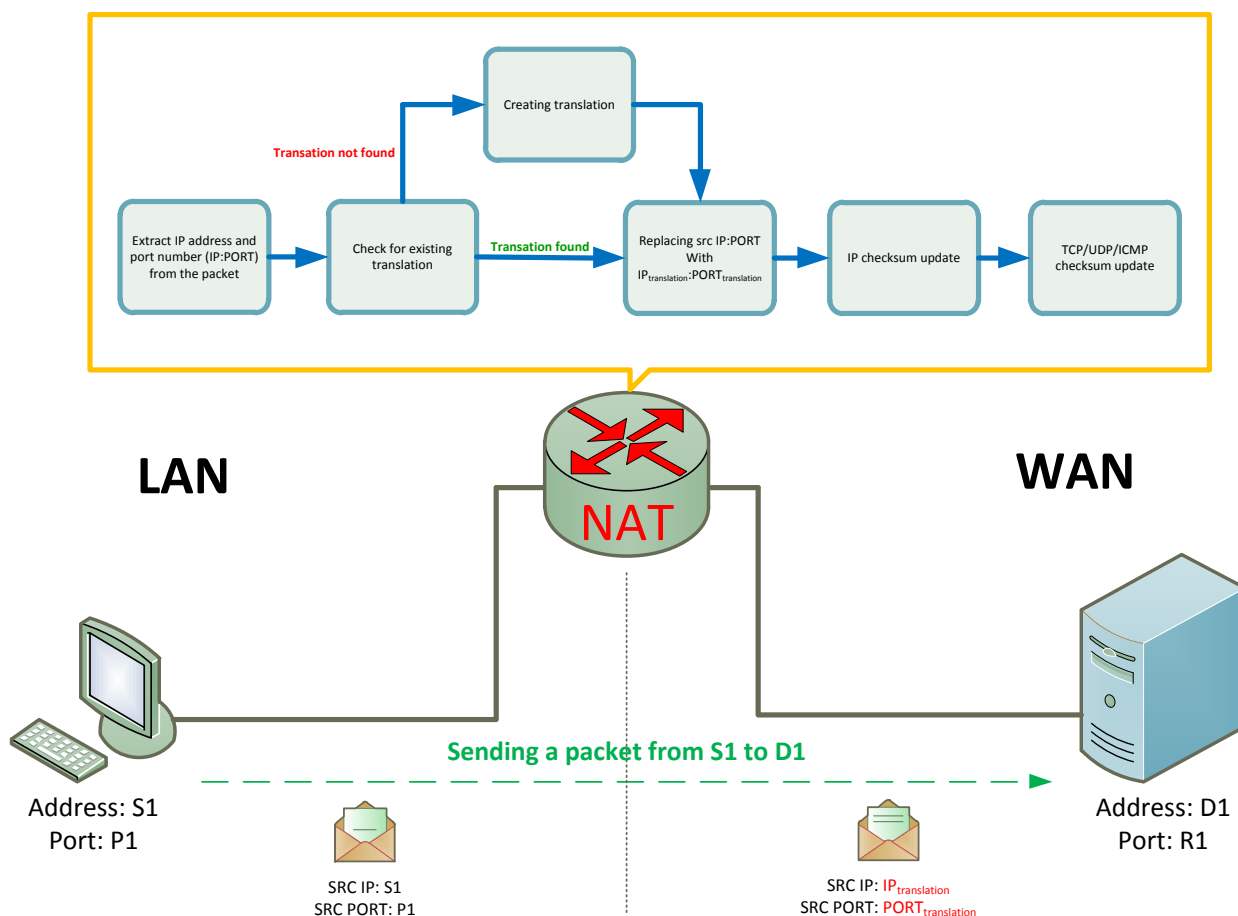
2.2.2 NAT bottleneck

For considering the bottlenecks the NAT system is conditionally split into 2 parts: transmitting part (packet receiving and sending) and processing part (translation). The transmitting part is closely related to hardware and OS. Modern hardware (CPUs and NICs) provide exceptional level of performance. The previous works [ref_coo] show that the maximum raw packet retransmitting rate is 23.06 Mpps

per 10 Gb port for 64-byte packet size. On bigger packet sizes the full throughput can be achieved. This limitation exists due to using PCIe Gen2 system bus. The processor is able to forward up to 92.22 Mpps. Having these numbers in mind, it is possible to claim that the hardware and OS (excluding its TCP/IP facilities) are able to exceed the required level of performance and, thus, are not the performance bottlenecks. This performance rate could be achieved by using Intel's DPDK framework [ref_dpdk].

The processing part consists of packet changing and translation managing parts. The packet changing part performs the packet headers modifications and has a constant runtime because it doesn't depend on input data. The tests, more closely described in part 3, shows that performing packet header changes only the packet processing rate can reach 20.1Mpps. This number also exceeds the required level of performance and is not a bottleneck.

The translation managing part is the most interesting one. It consists of several operations: storing the translation data, looking for the stored translation data, resources allocating and others. Storing the translation data and looking for the translation data are the core operations. They use a data structure and algorithm that perform storing and retrieving the data about translation to be done.



Scheme 2. Processing of packet going from LAN to WAN

The NAT translation is a critical part because of the required amount of data to be stored. The bigger the amount of data stored, the bigger the time of searching in this data chunk is. Thus, the NAT translation managing part is the bottleneck. To reduce the amount of runtime to be spent on packet translation, the translation storing data structure and searching algorithm should be found that provides higher or equal level of performance than required.

Having in mind that the NAT should be able to support 65.5M unique translations, it is easy to conclude that its lookup data structure has to be able to store 65.5M records and the search process will take a majority of packet processing time because of the size of this data structure. To achieve the target packet processing rate (5.5M pps) it is easy to calculate how much time is available for a packet processing. The time frame for a packet processing is around 180 ns. Considering

a processor working on frequency 2.4GHz we could spend no more than 436 cycles per packet. Hence, our target performance characteristic upper bound is 180 ns or 436 cycles per packet. This upper bound is used when evaluating the design approaches suggested further in this chapter in chapter 3.

2.2.3 The bottleneck overcoming. Exploration methodology

The bottleneck of the NAT is translation managing data structure and algorithm. It's not obvious from the first glance how to choose them properly. Furthermore, there are a lot of data structures and algorithms to test and trying each of them seems to be a too time consuming task. The scope of this work is to find a solution for the NAT that satisfies the certain requirements, so the goal is not to find the optimal solution and it is not required to cover all possible data structures and algorithms. Instead of this, the following methodology is used. The data structures and algorithms searching started from the simplest variants. If the solution doesn't satisfy the requirements, another, more complex, solution is suggested and tested. This process lasts while the suitable solution is found. The decision whether a solution is satisfying or not is made by testing it in the specially developed performance measurement program. The result, given by the tested algorithm, should be equal or lower the upper bound determined in the previous section.

2.2.4 Data structures and algorithms exploration

Linear search: So the first question to investigate is how fast the searching process is and does it really necessary to choose the algorithm and data structure. To answer that question the test has been performed which uses as a storing data

structure a simple linear array with linear search algorithm [ref_cormen]. This algorithm is known as having $O(n)$ search time and can be a good starting point of performance exploration. With high probability this approach is not suitable because of the linear nature of performance degradation. Having in mind the necessity to store the data for 65.5M instances the linear search is not able to provide with high performance. It can give some benefits when searching in small data sets consisting of tens of records.

Binary tree: To increase the performance another group of algorithms is taken into consideration. The group of algorithms based on tree-like data structures provide $O(\log N)$ theoretical search time. The simplest algorithm in that group is a simple binary search.

Its performance looks potentially promising but it consumes additional memory on tree node linking, in particular, each node uses 3 additional pointers to keep link with its parent node and 2 child nodes (left and right). Each of these links consumes at least 4 bytes of memory (12 in total) which leads to increasing of a NAT table entry size at least to 60%. So the overall memory overhead is more or equal than 60% depending on the CPU architecture and OS used.

The disadvantage of this algorithm is dependence from the input data storing sequence. In the worst case the binary tree can turn into a linked list. Searching in linked list has $O(n)$ searching time and has no benefits in comparison to the linear search. In fact, this situation is highly unlikely in the real world but has an implication of the binary tree being unbalanced. “Unbalanced” term means that the left and right branches of the binary tree have divers depth differing to each other more than 1. The consequence of this is that in the real world binary tree search cannot achieve $O(\log N)$ searching speed and the search time differs depending on the branch having the data to be searched.

RB-tree: There is an improved data structure and searching algorithm that fix the unbalancing property of the binary tree leaving the search time at the same level

of $O(\log N)$. The well-known representative of that group is red-black tree data structure. Searching the data saved in an rb-tree has the same workflow as searching in a binary tree. The distinction is in the storing (removing) logic. The main difference is that after each storing a piece of data in the rb-tree, the rb-tree is modified in a way to keep its structure balanced. This slightly increases leaving at the same order of $O(\log N)$ the time of storing the data but prevents the rb-tree from turning into a linked list.

Hash lists: Another group of the data structures and algorithms to be explored is one that uses a hash table based technics. The advantage of this group is a constant time ($O(1)$) of searching. For searching the data stored a key is used. The key is produced with a hash function from some data. The hash function has a deterministic value for a given input. It provides with good performance allowing to use simple (i.e. computationally cheap) hash functions.

The main issue while using a hash table is collision arising which occur when the hash function produce the same result for any two or more inputs. There are several schemes of this problem resolving. The most frequently used is separate chaining with linked lists. This means that each cell of the hash table is the first element of a linked list. If the hash function produces the same output for two or more different inputs all these values are stored in the same linked list. For searching the data scanning of the linked list should be done which might be a problem if the linked-list is long because of the $O(n)$ linked-list scanning time. Hence, the search time for the algorithm in case of using the modulo hash function is $O(n/\text{modulo_value})$. Other option of collisions resolving can be using a tree-like data structure instead of linked-list. Because of the $O(\log N)$ tree-based data structures searching time the search time for tree-like collision resolving in case of using the same modulo function as hash function is $O[\log(n/\text{modulo_value})]$. Theoretically the last collision resolving scheme gives better performance but in the real world the performance boost might be not so

good because of the CPU memory using schemes. To fix the memory issues it is possible to employ “cache-friendly” technics. It means that the memory for storing the collision resolving chains is allocated in chunks which are less or equal to the processors cache line size. This gives the elimination of cache misses improving the overall performance. The results of using this technic is discussed in Part 3.

Parallel processing: Almost all modern CPUs offer several cores on a single chip. Thus, using several cores for NAT routine seems to be a promising idea.

When employing multicores technics for software developing there are several issues to aware of which make significant influence on performance speed up.

The first issue is the cache coherence which arises when using shared data structures. Because of a copy of current processing data is stored in a core’s cache, changing the data by one core leads to updating the data in all cores currently use it. This means that the data have to overwritten in some common place of memory and then once again reloaded by other cores. This process is expensive and could cost hundreds of cycles which leads to considerable performance degradation.

The second issue is keeping the data in consistent state which is closely related to using special data structures known as locks. A lock also could be a problem because it makes the cores get access to the data in a sequence manner which can lead to core idling, decreasing the degree of parallelization. In the worst case it can lead to the result when the multicore code works with the same (or even less) performance it single core version.

Trying to avoid these issues the following approach was used. The biggest degree of parallelization can be achieved in case when a process running on core is fully independent from other processes (i.e. isn’t used the shared data). In case of our system this is possible because modern network interface cards provides multiple queues which can be used by different cores in associated manner when a given

queue is associated with one and only one core. In this case it can be seen as if a separate NAT process uses a core and a single network card and the amount of network cards installed in the system is more or equal to number of processors.

2.3 Implementation

For exploring the approach of building the NAT the testing application has been made. To simulate the NAT workflow several solutions have been implemented which use different data structures and software organization options described in paragraphs above. Conditionally the program is split into 3 parts: measuring part, generation part and simulating part.

The measuring part consists of the environment that performs testing routine and calculates the performance results. The metric produced by this part is cycles per packet. This metric is acquired by using **rdtsc** instruction which reads the internal processor tick counters. Then it recalculates the result from cycles to other metrics taking into account the CPU speed. The measuring part performs the number of tests, set by the user, and, produces the performance values.

The generation part generates a packet set to be processed by the simulation part. It imitates uniformly distributed network activity and stores generated packets in a one-dimensional array of structures which is the input to the simulation part. The packet set is fully stored in RAM of the computer. Time of packet set generation isn't taken into account when calculating the performance of the algorithm.

The simulation part is a core of the testing application and consists of NAT routine actions. There are several mandatory actions which must be performed by any NAT to perform proper packets translations. They are: calculation of the checksums, setting time stamp and saving/acquiring translation information in a NAT lookup data structure.

There are several necessary action to be performed by the NAT in order to perform address translation properly besides changing of packet's IP address and number of TCP/UDP port in the corresponding headers. They are: calculation of the checksum for IP and TCP/UDP/ICMP headers and storing the timestamp of the particular translation.

The checksum calculation is related to the packet processing. This action should be performed each time when the packet translation occurs and a packet IP and port number changed in order to be consistent with the requirements of the IP^[1.4 of ref_rfc791] and TCP/UDP/ICMP protocols ^[1.5 of ref_rfc793, rfc_768, rfc_5508].

The storing of the timestamp translation in the NAT translation data structure is necessary and cannot be eliminated because of the "Mapping Refresh" requirement for NATs ^[ref_rfc4787].

For the testing purposes in the NAT testing program the following function implementations are used. For checksum calculation *ip_fast_csum()* from Linux kernel is used. For getting/setting the timestamp *gettimeofday()* Linux system call is used. However, some different, faster, source of timestamp data can be used, for example CPU ticks counter which is faster but trickier when converting it to the physical time. These functions might be potential targets of performance optimization but are out of the scope of this document.

To store data about each address translation it is necessary to make the following record:

- source packet IP address [4 bytes] – IPv4
- source packet port number [2 bytes]
- IP address assigned by translation [4 bytes] – IPv4
- port number assigned by translation [2 bytes]
- translation timestamp – to calculate timeout [4 bytes]
- some additional service info (L4 protocol, flags) [4 bytes]

The total record size is 20 bytes. As it is needed to perform 2 translations for each connection it is necessary to save 2 records associated with that connection. So the total amount of data to save is 40 bytes per connection. This amount of data can be reduced. The memory reducing technic is described in the following chapters.

The target capacity of the NAT translation information data structure is 65.5M records which makes the NAT data structure space consuming. Having stored IP address and port number for each unique translation, minimum size of the data structure size is $65.5M * [4 \text{ (IP)} + 2 \text{ (port)} + 2 \text{ (timestamp)}] = 524 \text{ Mb}$. The NAT must be able to perform two translations for a single connection: from its inner network to its outer network and vice versa. So, it should have 2 similar data structures to store corresponding translation information. Thus, the amount of memory to be allocated is $2 * [\text{data structure size}]$ which is 1048 Mb in our case. This amount of memory is reduced by using a hash table with direct addressing for incoming packets. All translation specific data is stored in that hash table. The address for the hash table access is produced from incoming packet destination IP address and port using IP address as offset and port number as an index. The outgoing NAT table refers to that hash table and doesn't store any translation specific data.

There are implemented several version of the code: one for each approach suggested in 2.2.4 to get the data about the performance and make a decision about approach applicability.

The testing application was developed using C programming language with GCC<version> compiler(compiler flags were: -O3 -avx2) . The host operation system was native 64-bit Ubuntu 14.04(kernel version <version>).

For multi-threading **pthread** library was used. The approach of parallelizing the code is following. A process shares the same generated packet set but the set is split into parts and this parts associated to each thread so that no process packets

reads or writes interfere to other process packets reads or writes. Each thread is a separate network address translator with fully independent data structures (i.e NAT tables) which eliminates nearly all the pitfalls described 2.1 in the previous paragraph but leads to memory overheads.

Part 3. RESULTS

3.1 Evaluation

3.1.1 Measurement setup

All tests were done using a Lenovo laptop with the following characteristics:

- CPU: Intel i5-4210U (Haswell-ULT)
- RAM 8Gb Single-channel DDR3 @797MHz
- Motherboard: LENOVO Cherry 4A (U3E1)
- OS: Ubuntu 14.04

3.1.2 Experimental methodology

In getting metrics values the test setup plays a key role. The values of the metrics are highly dependent on the test methodology.

To get the values of interest the following setup was used. There were packets with mostly unique (more than 99%) tuples of IP and Port number in the packet set. The test routine gave this packet set as an input to the NAT system. The NAT system processed each packet and changed the values of IP, port number and checksums in the packet saving this data at the same packet set. Once all packets processed the routine performed backward translation simulating the response of the node from the outer network to the just translated and transmitted packet. After all packet processing has been done the check for translation correctness was performed.

This testing routine was used in order to simulate the most intensive regime of network working: the nodes of the network are constantly trying to communicate

with nodes in the outer network but the NAT device is offline, then the NAT device is switched on and right away starts serving the nodes connection, creating and performing new translations. This routine is more computationally intense than packet translation because the creating of a new connection costs more than a packet translating as it includes search for the connection translation data and if it was not found creating of a new translation entry is done.

This is the worst case scenario of network operating. The kind of testing used, allows getting the fair level of the NAT device performance.

3.2 Results discussion

The following paragraphs show sequential improvement attempts and results achieved.

From now and further in this document the NAT translation lookup data structure is called a NAT table and the translation record is called NAT table entry. The NAT table and the NAT table entry structures and sizes can differ in further described experiments depending on underlying data structure used in each particular experiment.

3.2.1 Results achieved

As it is seen from the Figure 1 the processing time of one packet excluding searching for translation data is a constant. This process is quite fast and can be compared with processors L3 cache miss penalty which is around 100 cycles.

So the first question to investigate is how fast the searching process is and does it really necessary to choose the algorithm and data structure. To answer that question the test has been performed which uses as a lookup data structure a simple linear array with linear search algorithm [ref_cormen]. This algorithm is known

as having $O(n)$ search time and can be a good starting point of performance exploration.

The results of testing the algorithm are shown in Figures 1 and 2. Linear search revealed the high linear performance degradation with increasing of the NAT records capacity: at size of 2000 entries the time of a packet processing is 3 times higher than at size of 500 entries and 3 times higher than the target performance. These results show that the translation data search is a real and quite serious bottleneck of the NAT performance and to solve this problem some effective algorithms and data structures are needed.

Nat Performance - Linear Search

CPU: i5-4210U, 1 core

[cycles/packet]

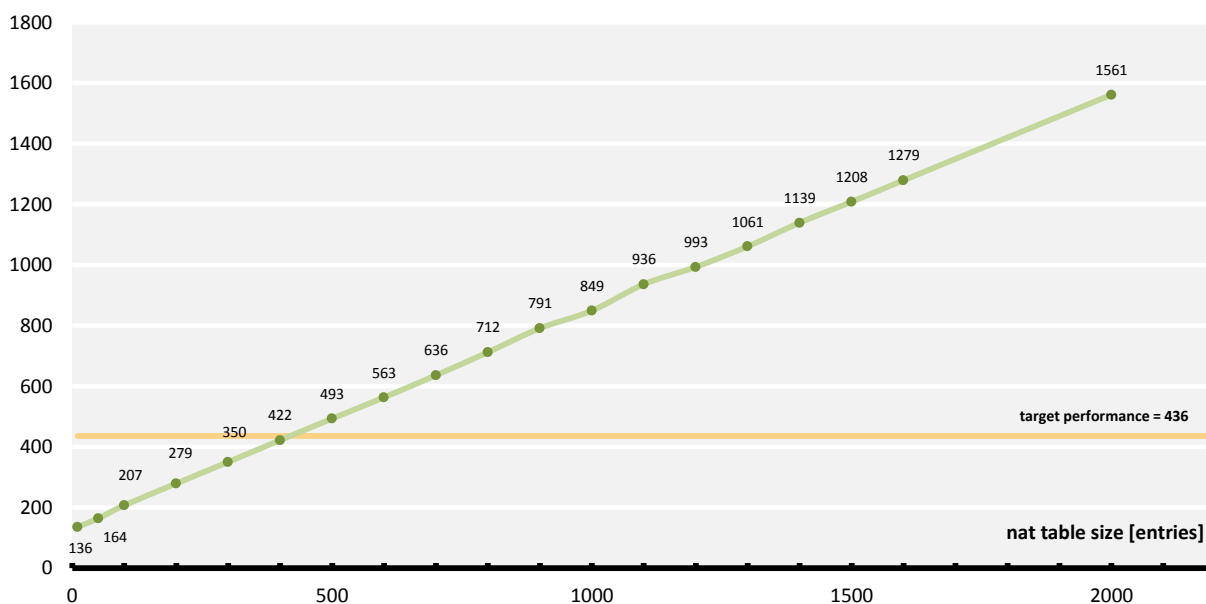
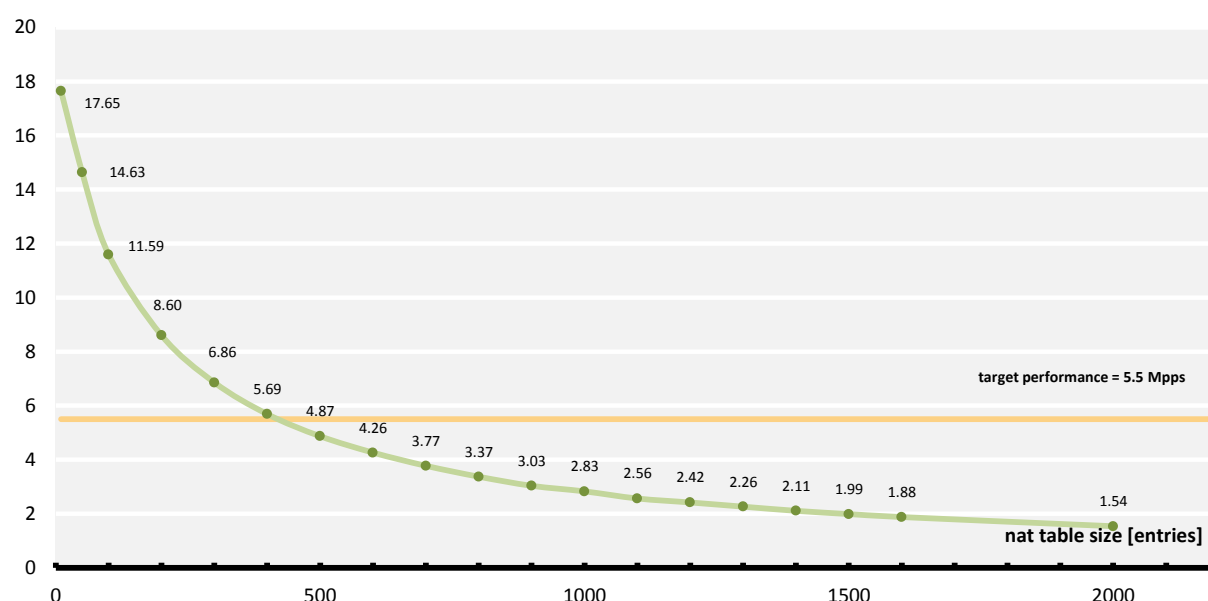


Figure 1. NAT Performance: Linear search in cycles per packet

Nat Performance - Linear Search

CPU: i5-4210U@2.4GHz, 1 core

[Mpps]

**Figure 2. NAT Performance: Linear search in packets per second**

The performance of tree-based data structures look potentially promising but it consumes additional memory on tree node linking, in particular, each node uses 3 additional pointers to keep the link with its parent node and 2 child nodes (left and right). Each of these links consumes at least 4 bytes of memory (12 in total) which leads to increasing of a NAT table entry size at least to 60%. So the overall memory overhead is more or equal than 60% depending on the CPU architecture and OS used.

The drawback of binary from the hardware point of view is low level of spatial locality^[ref_locality]. This happens because of the binary tree node's creation routine. The binary tree allows storing as many values as needed but, because of that, many memory allocations happen in different time frames are. Because these allocated chunks of memory are in the different parts of physical memory the CPU has to load each node in the memory separately instead of loading several of them at a time. As the NAT table size is much larger than a cache, even on a highly

advanced CPUs with a big cache size, only some of the most frequently used nodes stays in the cache. All other nodes are constantly removed and stored again which means that the CPU spends a majority of time on data transferring from the main memory to the cache and vice versa.

Nat Performance - Binary tree-based NAT table

CPU: i5-4210U, 1 core

[cycles/packet]

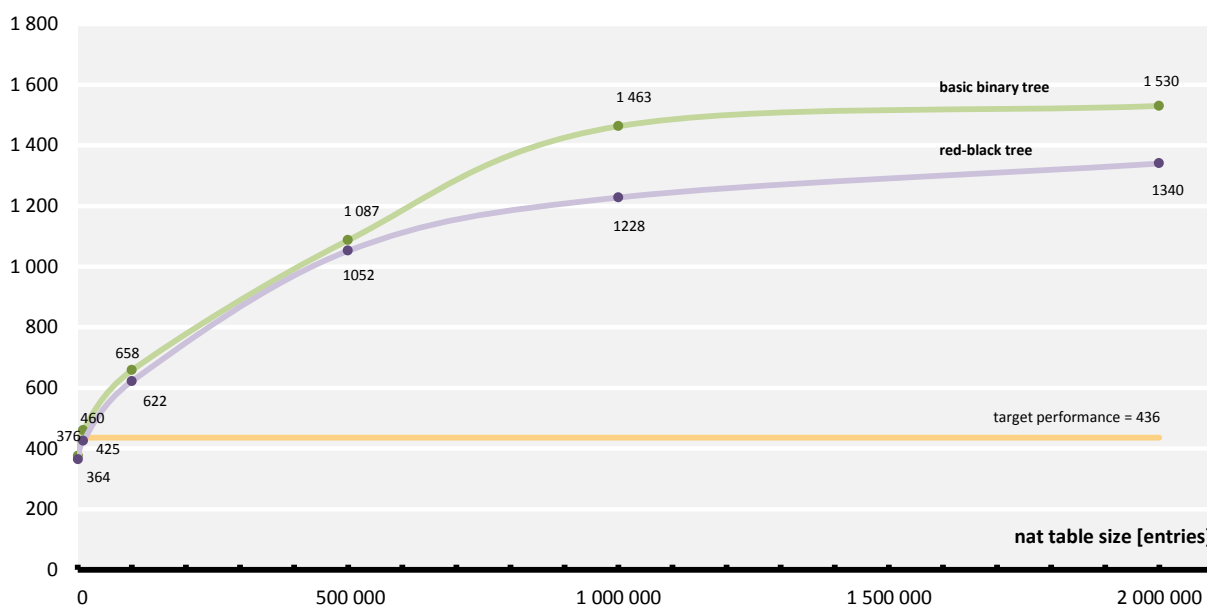


Figure 3. Nat Performance in cycles per packet. Binary tree-based NAT tables. Comparison between simple unbalanced and red-black binary trees.

Nat Performance - Binary tree-based NAT table + 1D array

CPU: i5-4210U@2.4GHz, 1 core

[Mpps]

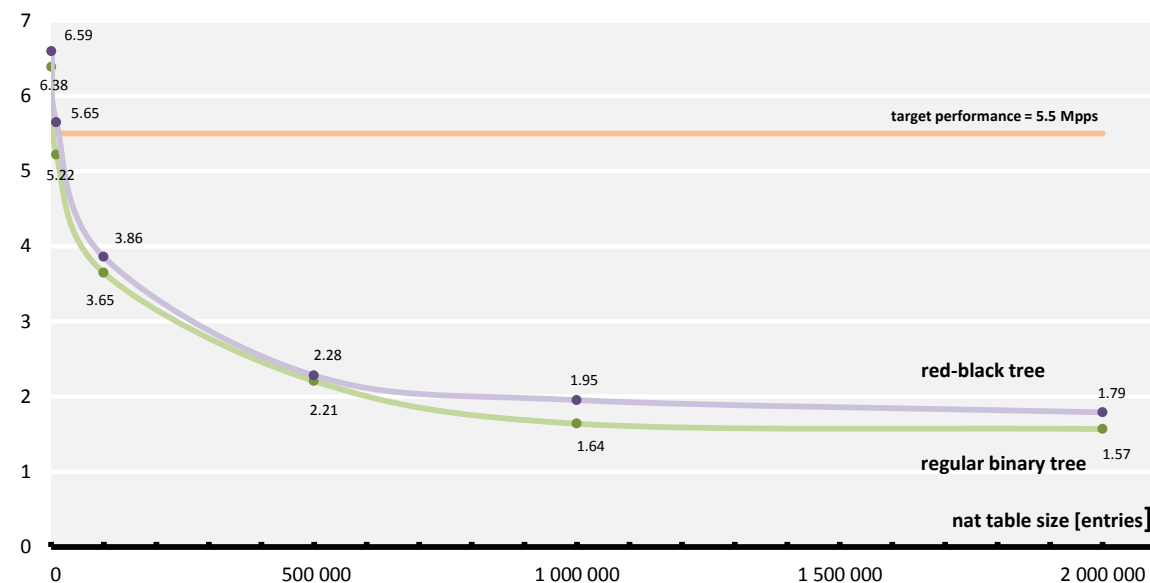


Figure 4. Nat Performance in packets per second. Binary tree-based NAT tables. Comparison between simple unbalanced and red-black binary trees.

The figure 3 and 4 show the results of using simple unbalanced and red-black binary trees. The performance improvement that gives red-black tree data structure is around 14% and still 3 times below the desired value.

The results of using tree-based data structures revealed the fact that the binary-trees group of data structures cannot be used for achieving of the target level of performance. The maximum result they could provide is 3 times slower than needed. For further investigation of the performance more fast data structures should be taken into consideration.

The next group of algorithms to be tested is hash-table based algorithms.

Nat Performance - Hash-based NAT table

CPU: i5-4210U@2.4 GHz, 1 core

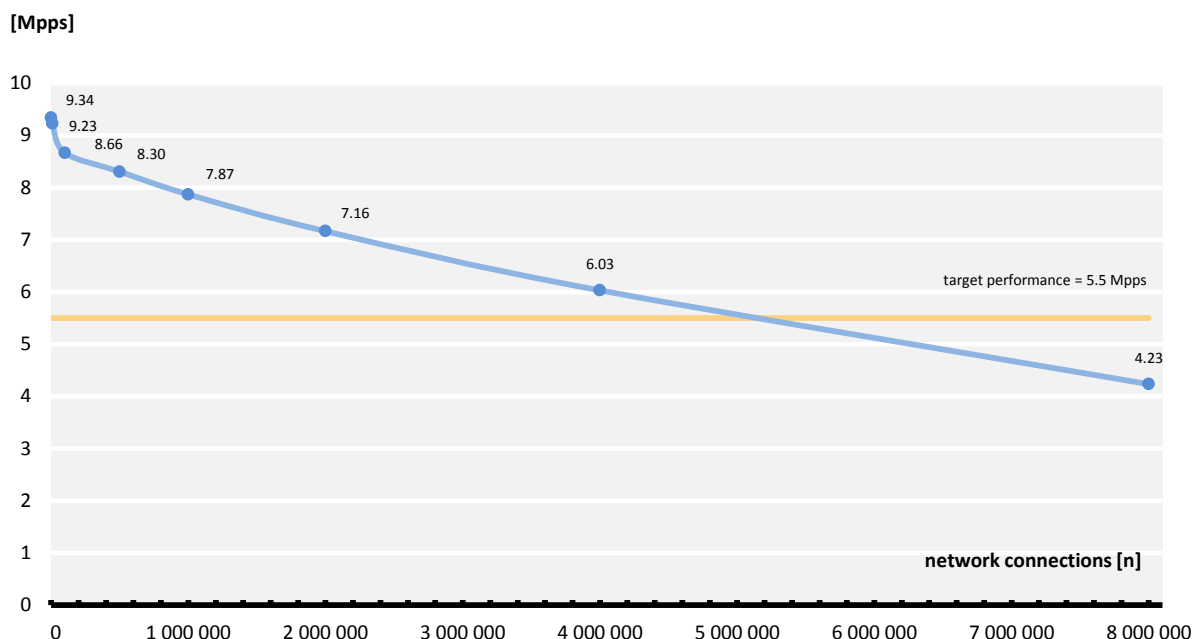
modulo value: 2^{23} 

Figure 5. Nat Performance in packets per second. Hash-based NAT table.

For implementing the NAT table modulo division operation was used as the hash function. The NAT table capacity is set to store 1000 translation for each of 65536 nodes of the supported network. The hash table cell includes supplementary data which includes a link to a corresponding cell in the incoming array for accelerating of new translation creating time. The results of using the hash table are shown in Figure 5. There is a significant performance increasing in comparison with the tree-based NAT tables: it is 4 times faster. Because of the test system limitations it is not seen from the chart what the NAT performance is at the target level of connections number. Although, it is not shown in Figure 5, it is possible to estimate the full load performance, using the essence of the separate chaining linked list hash table and the results of linear search. The worst case scenario search time for this data structure is determined by the maximum length of the linked list associated with the cell calculated by hash function. Thus, changing the divisor value we could adjust the search time. But with changing the initial value the size of the hash table varies: than bigger the divisor that bigger its initial size.

This leads to significant memory overheads. Using the big hash table significantly exceeding the number of possible connections significantly reduce the probability of collision (the same key appearing) which, in turn, reduces the length of the chain. This leads to big memory overheads. The overheads could be eliminated by using alternative schemes of hash table collision resolving. The 4 of them were tested. The best result were shown by hash table with cache optimized list chaining storing 3 list items in a processor's cache line. The increasing of list items amount could further speed up the performance but is limited with a cache line size of 64 bytes. The results of comparison are shown in Figure 6 and 7.

Taking into account the fact that the 3-item cache line optimized hash table supplies the target level of performance for number of connection 4 times bigger than its initial capacity, the estimated initial size (and modulo value) of the hash table capable to support the target value of connections (65.5M) is $2^{26}/4 = 2^{24}$ (around 300M). Relatively small memory consumption of the hash table becomes important when one trying to take advantage from parallelization.

Hash table performance comparison

CPU: i5-4210U, 1 core

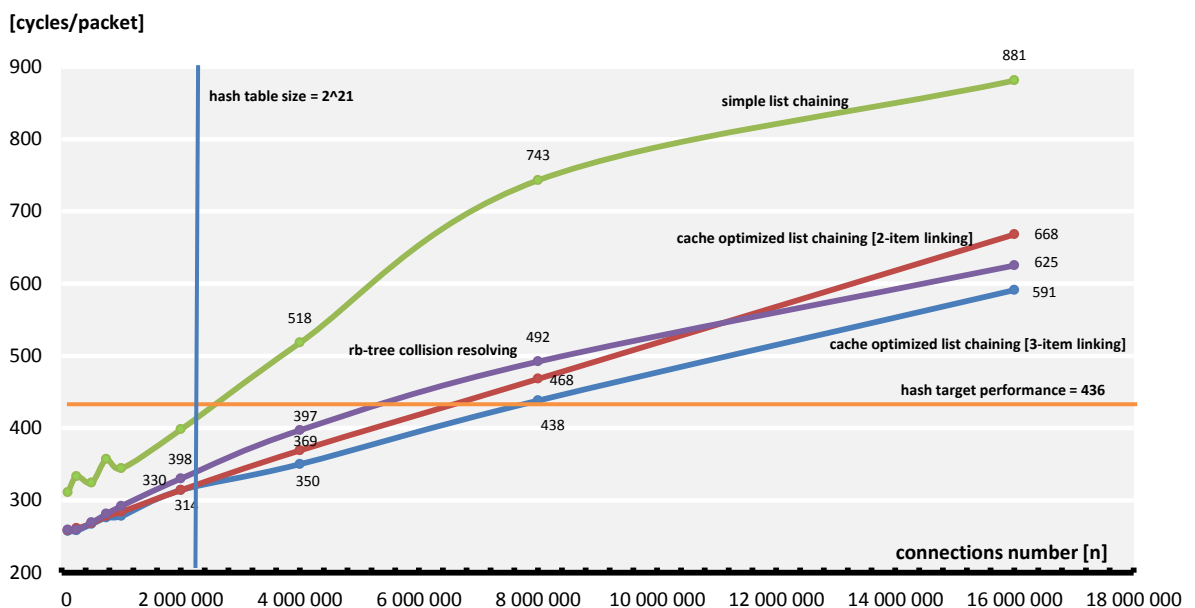


Figure 6. The comparison of hash tables with different types of collisions resolving in cycles per packet.

Hash table performance comparison

CPU: i5-4210U, 1 core

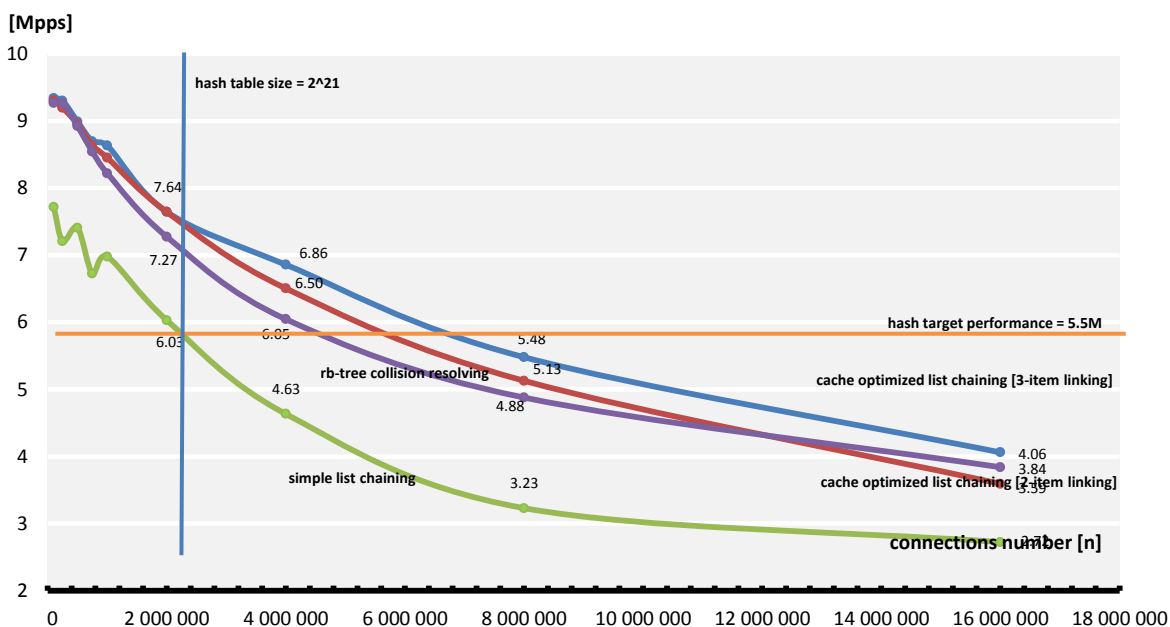


Figure 7. The comparison of hash tables with different types of collisions resolving in packets per second.

Further speed up can be gotten using parallelization of the translation process. All the results described previously in this document were gotten with 1 core at a multicore processor. Almost all modern CPUs offer several cores on a single chip. Thus, using several cores for NAT routine seems to be a promising idea.

In figure 8 the results are shown for different number of cores. The experiments were set for 4 threads maximum because the testing system has 4 hardware threads and setting experiment with more threads are meaningless because of the maximum number of simultaneously running cores is equal to number of cores available on the CPU chip and the OS scheduling issues, i.e. an operating system spends additional time on context switching and doesn't give any performance boost. The results shown reveal that the parallel approach gives the desired results: than more cores are available than grater performance improvement would be. Furthermore, the solution has a good scalability: adding an additional thread gives the same increasing of performance. This happens

because of the design used: each thread is an NAT with fully independent data structures which is used by a certain CPU core number only.

Nat Performance

CPU: i5-4210U, 4 cores

[Mpps]

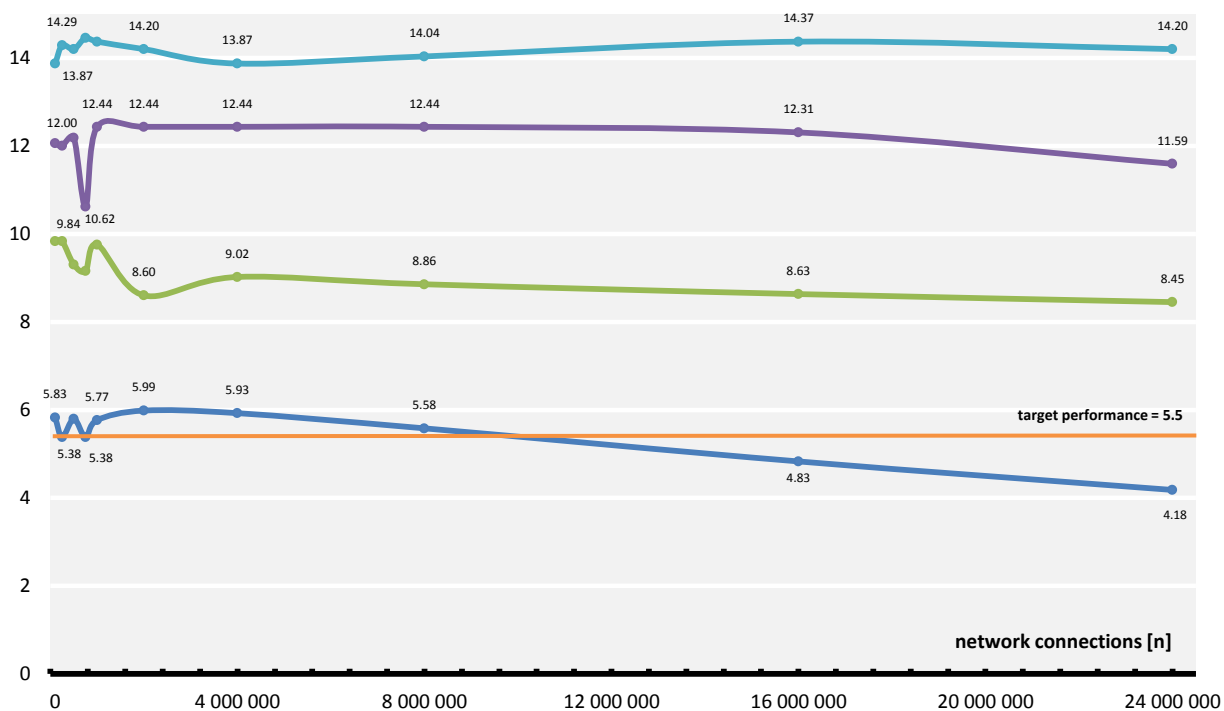


Figure 8. Nat Performance when using parallel 2-item optimized hash-based NAT table. Results for different number of simultaneously working cores.

This design is applicable in this case because the resources can be split evenly between the threads and modern Network Interface Cards provide the ability to split the incoming packets in several independent queues used by the threads independently.

3.2.2 Final design solution

Summing up the results, it is reasonable to conclude that using cache-friendly (3-item optimized) hash tables as a data structure for storing the translations

information which works in parallel manner is a reasonable choice for the high-performance NAT. It gives the performance rate 260% more than target level and has a significant reserve for further NAT algorithms complications. This data structure is used to store outgoing translations. For incoming translation a one-dimensional array is used.

The code for the application is written in C, except function for fast checksum calculation written using Assembler language and taken from Linux kernel. This approach promises to give more performance while using more advanced processors and compilers in the future.

Further modifications and optimizations of the NAT table structure based on hash-table could be done to improve the currently achieved performance values by using more sophisticated data structures, as well as other hardware-friendly optimization technics for example cuckoo hashing_[ref_coo], but this is the matter of further research.

3.2.3 The software CG-NAT price estimation

Having the achieved results and assuming that this level is feasible not only in the simulation but in the real world, it is reasonable to estimate the price of NAT system using the results of this research and compare it with competing systems.

As a base for the system the average prices on the Russian market is used. To estimate the price, the computer setup is used with similar characteristics that the computer has which was used for testing in this work. The setup is the following:

Component name	Model	Price, \$
Processor	Intel Core i5-4690 <small>[ref_pint]</small>	220
RAM Modules	SiliconPower SP016GXLYU16ANDA X2 <small>[ref_pmem]</small>	320
Motherboard	SuperMicro X10SLL-S (Intel C222) <small>[ref_pmb]</small>	200
Hard Disk	Intel DC S3610 <small>[ref_phard]</small>	275
Network Interface Card	Lenovo 10Gb X540-T2 <small>[ref_pnic]</small>	600
System Unit	SuperMicro CSE-732D2-500B <small>[ref_psysu]</small>	200

Total: \$1815

As it is seen from the estimation, one spending a little bit more than 1800 USD can build a NAT system which could work as a NAT having the same performance rate as specialized devices which cost at least 3 times higher than tested one (see Figure 1).

3.3 Summary

<to be done>

Only core functionality was considered

Fulfilling and exceeding all target characteristics

A laptop processor was used

REFERENCES

- [ref_TT_ROS_TEL] file:TT CGNAT 2014_26_06v1.doc
- [ref_rfc3022] Traditional IP address translator. <https://www.ietf.org/rfc/rfc3022.txt>
- [ref_RDP.RU] <http://rdp.ru/>
- [ref_rfc4787] Network Address Translation (NAT) Behavioral Requirements for Unicast UDP <https://tools.ietf.org/html/rfc4787#page-5>
- [ref_rfc5382] NAT Behavioral Requirements for TCP <https://tools.ietf.org/html/rfc5382>
- [ref_rfc5508] NAT Behavioral Requirements for ICMP <https://tools.ietf.org/html/rfc5508>
- [ref_rfc791] Internet protocol <https://www.ietf.org/rfc/rfc791.txt>
- [ref_rfc793] Transmission Control Protocol <https://www.ietf.org/rfc/rfc793.txt>
- [ref_rfc768] User Datagram Protocol <https://www.ietf.org/rfc/rfc768.txt>
- [ref_rfc792] Internet Control Message Protocol <https://tools.ietf.org/html/rfc792>
- [ref_cormen] Introduction to algorithms, By [Thomas H. Cormen](#), [Charles E. Leiserson](#), [Ronald L. Rivest](#) and [Clifford Stein](#) ISBN: 978026203384
- [ref_locality] http://en.wikipedia.org/wiki/Locality_of_reference
- [ref_ripe_limit] <https://www.ripe.net/publications/ipv6-info-centre/about-ipv6/ipv4-exhaustion>
- [ref_rfc6888] Common Requirements for Carrier-Grade NATs (CGNs) <http://www.rfc-base.org/txt/rfc-6888.txt>
- [rfc_3489] STUN - Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs) <https://tools.ietf.org/html/rfc3489>
- [ref_nat32] NAT 32 IP Router <http://v2.nat32.com/index.html>
- [ref_cam] <http://www.digikey.com/product-detail/en/EP20K160EBC356-2X/EP20K160EBC356-2X-ND/4160824>
- [ref_wingate] <http://www.wingate.com/purchase/wingate/purchase.php>
- [ref_IX] A Protected Dataplane Operating System for High Throughput and Low Latency. Adam Belay, Stanford University; George Prekas, École Polytechnique Fédérale de Lausanne (EPFL); Ana Klimovic, Samuel Grossman, and Christos Kozyrakis, Stanford University; Edouard Bugnion, École Polytechnique Fédérale de Lausanne (EPFL) <https://www.usenix.org/conference/osdi14/technical-sessions/presentation/belay>
- [ref_uppc] Understanding the Packet Processing Capability of Multi-Core Servers. Norbert Egi[‡], Mihai Dobrescu[†], Jianqing Du[†], Katerina Argyraki[†], Byung-Gon Chun[§], Kevin Fall[§], Gianluca Iannaccone[§], Allan Knies[§], Maziar Manesh[§], Laurent Mathy[‡], Sylvia Ratnasamy[§] § Intel Research, † EPFL, ‡ Lancaster University

- [ref_epssr] M. Dobrescu, N. Egi, K. J. Argyraki, B.-G. Chun, K. R. Fall, G. Iannaccone, A. Knies, M. Manesh, and S. Ratnasamy. RouteBricks: Exploiting Parallelism to Scale Software Routers. In Proceedings of the 22nd ACM Symposium on Operating Systems Principles (SOSP '09), pages 15–28, 2009
- [ref_click] E. Kohler, R. Morris, B. Chen, J. Jannotti, and M. F. Kaashoek. The Click Modular Router. ACM Trans. Comput. Syst., 18(3):263–297, 2000.
- [ref_cspp] Computer Systems: A Programmer's Perspective (2nd Edition) by Randal E. Bryant, David R. O'Hallaron, ISBN-13: 978-0136108047
- [ref_dirca] Direct Cache Access for High Bandwidth Network I/O. Ram Huggahalli, Ravi Iyer, Scott Tetrick, Intel Corporation Proceeding, [ISCA '05](#) Proceedings of the 32nd annual international symposium on Computer Architecture, Pages 50-59
- [ref_amdahl] Validity of the single processor approach to achieving large scale computing capabilities, Gene M Amdahl, IBM Sunnyvale California
- [ref_ipchal] IPv4 to IPv6: Challenges, solutions, and lessons, Stanford L. Levin a, Stephen Schmidt. Telecommunications Policy 38 (2014) 1059–1068
- [ref_depgog] Deploying IPv6 in the Google Enterprise Network. Lessons learned. Haythum Babiker, Irena Nikolova, Kiran Kumar Chittimaneni, Google, USENIX LISA, 2011
- [ref_ecrc] Exploiting a Computation Reuse Cache to Reduce Energy in Network Processors, Bengu Li, Ganesh Venkatesh, Brad Calder, and Rajiv Gupta, University of Arizona, University of California
- [ref_ehms] An Efficient Hardware-based Multi-hash Scheme for High Speed IP Lookup, Socrates Demetriades, Michel Hanna, Sangyeun Cho and Rami Melhem, University of Pittsburgh, 16th IEEE Symposium on High Performance Interconnects
- [ref_natimp] NAT implementation for the NetFPGA platform, Omar Choudary, David J. Miller, University of Cambridge, Cambridge, UK
- [ref_coo] Scalable, High Performance Ethernet Forwarding with CUCKOOSWITCH, Dong Zhou, Bin Fan, Hyeontaek Lim, David G. Andersen Carnegie Mellon University; Michael Kaminsky (Intel Labs), 10th USENIX Symposium on Networked Systems Design and Implementation, 2013
- [ref_dpdk] Intel Data Plane Development Kit (Intel DPDK) Overview.
<http://www.intel.com/content/dam/www/public/us/en/documents/presentation/dpdk-packet-processing-ia-overview-presentation.pdf>
- [ref_pint] http://ark.intel.com/ru/products/80810/Intel-Core-i5-4690-Processor-6M-Cache-up-to-3_90-GHz
- [ref_pmem] http://www.nix.ru/autocatalog/memory_modules_SiliconPower/Silicon_Power_SP016GXLY_U16ANDA_DDRIII_16Gb_8Gb_PC312800_CL9_204263.html?set_id=be30bbd1ed9211e4a20b002590c35102&vs_id=be1945beed9211e4a20b002590c35102&sort=0&from=sch&sch_id=34&sch_good=190779
- [ref_pmb] <http://www.srv-trade.ru/catalog/1824837288/CSE-732D2-500B/>

http://www.nix.ru/autocatalog/motherboards_supermicro/SuperMicro_X10SLLS_LGA1150_C222_PCIE_2GbLAN_SATARAID_microATX_2DDRIII_159355.html

[ref_pnic]

<http://shop.lenovo.com/us/en/itemdetails/0C19497/460/BF8C0B7DC9BB4C13B5EF4FE54E3ABB39>

[ref_psysu] <http://www.srv-trade.ru/catalog/1824837288/CSE-732D2-500B/>

[ref_phard][http://www.nix.ru/price/price_list.html?section=ssd_all&set_id=9e029764ed9711e4a20b002590c35102#cid=900&fn=900&set_id=a574a849ed9711e4a20b002590c35102&sort="+p1710&thumbnai](http://www.nix.ru/price/price_list.html?section=ssd_all&set_id=9e029764ed9711e4a20b002590c35102#cid=900&fn=900&set_id=a574a849ed9711e4a20b002590c35102&sort=)
l_view=1

Appendix A

List of router models used for market research

Vendor	Router Model	Packet processing rate, Mpps	Price, USD	URL
HP	MSR2021	0.18	760	http://www8.hp.com/ru/ru/products/networking-routers/product-detail.html?oid=5054094#!tab=specs
HP	MSR1002-4	0.50	1 100	http://www8.hp.com/us/en/products/networking-routers/product-detail.html?oid=6288749#!tab=specs
HP	MSR2003	1.00	1 700	http://www8.hp.com/ru/ru/products/networking-routers/product-detail.html?oid=5408900#!tab=specs
HP	MSR50-40	1.20	2 500	http://www8.hp.com/us/en/products/networking-routers/product-detail.html?oid=4199527#!tab=specs
HP	MSR3012	2.60	2 700	http://www8.hp.com/us/en/products/networking-routers/product-detail.html?oid=6288370&jumpid=reg_r1002_usen_c-001_title_r0002#!tab=specs
HP	MSR3040	0.36	3 200	http://www8.hp.com/ru/ru/products/networking-routers/product-detail.html?oid=4199541#!tab=specs
HP	MSR4060	10.00	6 000	http://www8.hp.com/us/en/products/networking-routers/product-detail.html?oid=5408896
CISCO	7301	1.00	10 000	http://www.cisco.com/web/RU/products/hw/routers/ps352/ps4972/index.html
CISCO	ASA5515-IPS-K9	0.50	4 800	http://ciscosales.ru/katalog_produkcii/cisco/mezhsetevye_ekrany_i_fil_try/cisco_asa_5500_series_accessories/asa5515-ips-k9/ http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-a
CISCO	ASA 5525-X	0.70	5 200	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-b
CISCO	ASA 5545-X	0.90	10 200	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-b
CISCO	ASA 5555-X	1.00	17 000	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-b

CISCO	5585-X SSP10	1.10	14 200	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-c http://ciscosales.ru/katalog_produkcii/cisco/mezhsetevye_ekrany_i_fil_try/cisco_asa_5500_series_firewall_edition_bundles/asa5585-s10-k9/
CISCO	5585-X SSP20	2.00	40 000	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-c http://ciscosales.ru/katalog_produkcii/cisco/mezhsetevye_ekrany_i_fil_try/cisco_asa_5500_series_firewall_edition_bundles/asa5585-s20-k8/http://ciscosales.ru/katalog_produkcii/cisco/mezhsetevye_ekrany_i_fil_try/cisco_asa_5500_series_firewall_edition_bundles/asa5585-s20-k8/
CISCO	5585-X SSP40	4.00	80 000	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-c http://ciscosales.ru/katalog_produkcii/cisco/mezhsetevye_ekrany_i_fil_try/cisco_asa_5500_series_firewall_edition_bundles/asa5585-s40-k8/
CISCO	5585-X SSP60	10.00	128 300	http://www.cisco.com/c/en/us/products/security/asa-5500-series-next-generation-firewalls/models-comparison.html#~tab-c http://ciscosales.ru/katalog_produkcii/cisco/mezhsetevye_ekrany_i_fil_try/cisco_asa_5500_series_firewall_edition_bundles/asa5585-s60-2a-k9/
JUNIPER	SRX240	0.20	2 000	http://www.juniper.net/us/en/local/pdf/datasheets/1000281-en.pdf http://www.srv-trade.ru/catalog/976735833/SRX240H/?gclid=CKXzk-_4ysMCFYPUcgodBmUALQ
JUNIPER	SRX550	0.70	7 900	http://www.juniper.net/us/en/local/pdf/datasheets/1000281-en.pdf
JUNIPER	SRX650	0.85	13 000	http://www.juniper.net/us/en/local/pdf/datasheets/1000281-en.pdf