



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



Nico Schottelius

High speed NAT64 with P4

Master Thesis MA-2019-19
February 2019 to August 2019

Tutors: Alexander Dietmüller, Edgar Costa Molero, Tobias Bühler
Supervisor: Prof. Dr. Laurent Vanbever

Abstract

Due to the lack of IPv4 addresses, IPv6 deployments have recently gained in importance in the Internet. Several transition mechanisms exist that include translating IPv6 packets into IPv4 packets, thus enabling the coexistence and interoperability of both protocols.

This thesis describes an implementation of the translation mechanism NAT64, implemented in P4. Using the P4 programming language a software emulated switch was created that translates IPv4 to IPv6 and vice versa. Due to the target independence of P4 the same code can be compiled for and deployed to the FPGA hardware platform "NetFPGA".

Within the NetFPGA the NAT64 implementation achieves a stable throughput of 9.28 Gigabit/s. Our solution allows in-network translations without a router or client configurations. Due to the nature of P4, the implementation runs at line speed and thus with different hardware the same code can run potentially at much higher speeds, for instance on 100 Gbit/s switches.

Contents

1	Introduction	9
1.1	IPv4 exhaustion and IPv6 adoption	9
1.2	Motivation	10
2	Background	13
2.1	P4	13
2.2	IPv6, IPv4 and Ethernet	13
2.3	ARP and NDP, ICMP and ICMP6	14
2.4	IPv6 Translation Mechanisms	15
2.4.1	Stateless NAT64	15
2.4.2	Stateful NAT64	15
2.4.3	Higher Layer Protocol Dependent Translation	16
2.4.4	Mapping IPv4 Addresses in IPv6	16
2.4.5	DNS64	17
2.5	Protocol Checksums	17
2.6	Network Designs	18
2.6.1	IPv4 only network limitations	19
2.6.2	Dualstack network maintenance	19
2.6.3	IPv6 only networks	20
2.7	NetFPGA	20
3	Design	23
3.1	P4/NAT64	23
3.2	P4/BMV2	24
3.3	P4/NetFPGA	24
3.4	Stateless NAT64	25
3.5	Stateful NAT64	26
3.6	NAT64 Verification	27
3.7	IPv6 and IPv4 configuration	27
4	Results	29
4.1	P4 based implementations	29
4.2	P4/BMV2	30
4.3	P4/NetFPGA	30
4.3.1	Features	30
4.3.2	Stability	31
4.3.3	Usability	32
4.4	Software based NAT64	34
4.5	NAT64 Benchmarks	34
4.5.1	Benchmark Design	34
4.5.2	Benchmark Summary	34
4.5.3	IPv6 to IPv4 TCP Benchmark Results	36
4.5.4	IPv4 to IPv6 TCP Benchmark Results	36
4.5.5	IPv6 to IPv4 UDP Benchmark Results	37
4.5.6	IPv4 to IPv6 UDP Benchmark Results	37
5	Conclusion and Outlook	39

A Resources and code repositories	41
A.1 Operating Systems	41
A.2 Master Thesis	41
A.3 Xilinx Toolchain	41
A.4 P4/NetFPGA support scripts	41
A.5 P4/NetFPGA compilation process	42
A.6 P4/NetFPGA Tests	42
A.6.1 Test 1: IPv4 Egress	42
A.6.2 Test 2: IPv6 egress	43
A.7 P4/BMV2 environment and tests	43
B NetFPGA Logs	45
B.1 NetFPGA Flash Errors	45
B.2 NetFPGA Flash Success	45
B.3 NetFPGA Kernel module	46
B.4 NetFPGA compile logs	47
C Benchmark Logs	55
C.1 Enabling hardware offloading	55
C.2 Tayga	56
C.3 Jool	58

List of Figures

1.1	RIR IPv4 rundown projection [26]	9
1.2	LACNIC Exhaustion projection [29]	10
1.3	Google IPv6 Statistics from [22]	10
1.4	Separated IPv6 and IPv4 network segments	11
2.1	P4 protocol independence [60]	13
2.2	ARP and NDP	14
2.3	ICMP6 option fields	15
2.4	IPv6 Header [18]	16
2.5	IPv4 Header [45]	17
2.6	Stateful NAT64	17
2.7	Representing an IPv4 address in an IPv6 prefix	18
2.8	IPv4 embedding depending on the prefix length	18
2.9	Illustration of DNS64	19
2.10	IPv6 Pseudo Header	19
2.11	IPv4 Pseudo Header	20
2.12	NetFPGA Board, [67]	20
3.1	P4 Switch Architecture	23
3.2	Standard NAT64 translation	24
3.3	In-network NAT64 translation	24
3.4	Our P4 Switch Architecture	25
3.5	P4/BMV2 checksumming	26
3.6	Calculating checksum based on header differences	27
4.1	Hardware Test NetPFGA card 1	32
4.2	Hardware Test NetPFGA card 2, [23]	32
4.3	Benchmark design for NAT64 in software implementations	34
4.4	NAT64 with NetFPGA benchmark	35

List of Tables

3.1	NAT64 match factors	25
3.2	NAT64 verification commands	28
3.3	IPv6 address and network overview	28
3.4	IPv4 address and network overview	28
4.1	P4/BMV2 feature list	30
4.2	P4/NetFPGA feature list	31
4.3	IPv6 to IPv4 TCP NAT64 Benchmark	36
4.4	IPv4 to IPv6 TCP NAT64 Benchmark	36
4.5	IPv6 to IPv4 UDP NAT64 Benchmark	37
4.6	IPv4 to IPv6 UDP NAT64 Benchmark	37



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

HIGH SPEED NAT64 WITH P4

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

SCHOTTELIUS

First name(s):

NICO

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

LINTHAL, 201908 19

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.

Chapter 1

Introduction

In this chapter we give an introduction about the topic of the master thesis, the motivation, and problems that we address. We explain the current state of IPv4 exhaustion and IPv6 adoption and describe how it motivates our work to support ease transition to IPv6 networks.

1.1 IPv4 exhaustion and IPv6 adoption

The Internet has almost completely run out of public IPv4 space. The 5 Regional Internet Registries (RIRs) report IPv4 exhaustion worldwide [50], [4], [29], [1], [5]. Figure 1.1 contains summarised data from all RIRs and projects complete IPv4 addresses depletion by 2021. The LACNIC project even predicts complete exhaustion for 2020 as shown in figure 1.2.

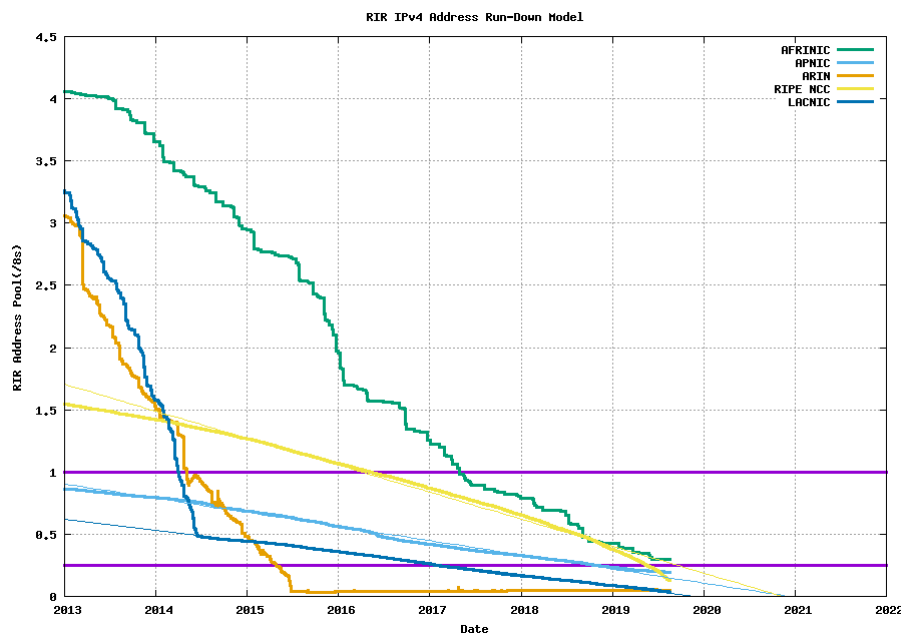


Figure 1.1: RIR IPv4 rundown projection [26]

On the other hand, IPv6 adoption grows significantly, with at least three countries (India, US, Belgium) surpassing 50% adoption [2], [62], [14]. Traffic from Google users reaches almost 30% as of 2019-08-08 [22], see figure 1.3.

We conclude that IPv6 is a technology strongly gaining importance. IPv4 depletion is estimated to be happening worldwide in the next years. Thus more devices will be using IPv6, while communication with legacy IPv4 devices still needs to be provided.

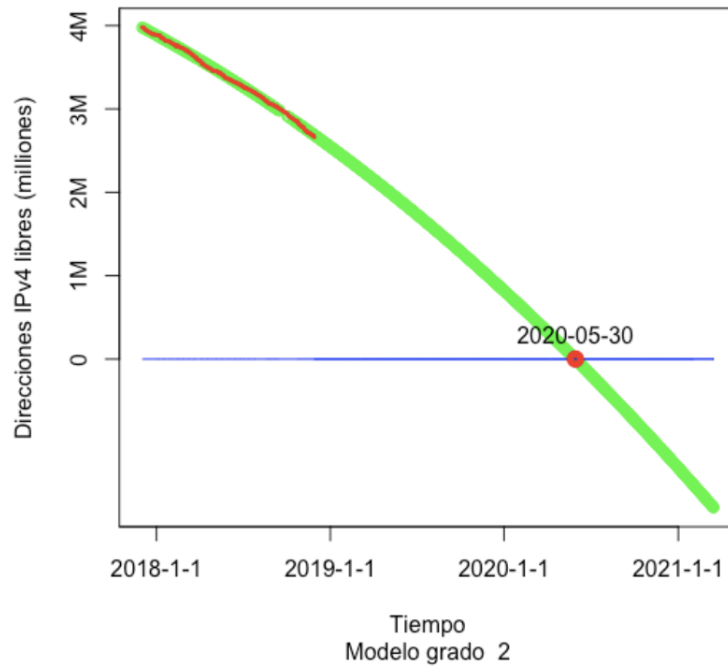


Figure 1.2: LACNIC Exhaustion projection [29]

1.2 Motivation

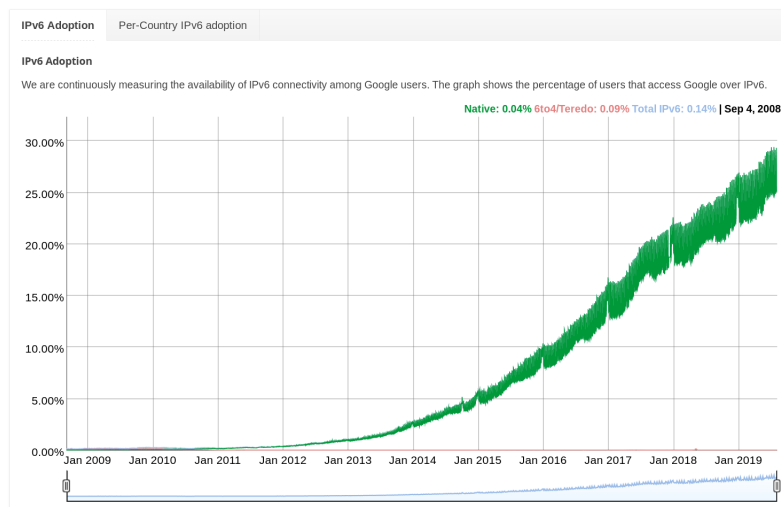


Figure 1.3: Google IPv6 Statistics from [22]

IPv6 hosts and IPv4 hosts cannot directly connect to each other, because the protocols are incompatible to each other. To allow communication between different protocol hosts, several transition mechanisms have been proposed [64], [41]. However installation and configuration of the transition mechanism usually require in-depth knowledge about both protocols and require additional hardware to be added in the network. In this thesis we show an in-network transition method based on NAT64 [6]. Compared to traditional NAT64 methods which require hosts to explicitly use an extra device in the network,¹ our proposed method is transparent to the hosts. This way the routing and network configuration does not need to be changed to support NAT64 within a network. Currently network operators have to focus on two network stacks when designing networks: IPv6 and IPv4. While in a small scale setup this might not introduce significant

¹Usually the default router will take this role.

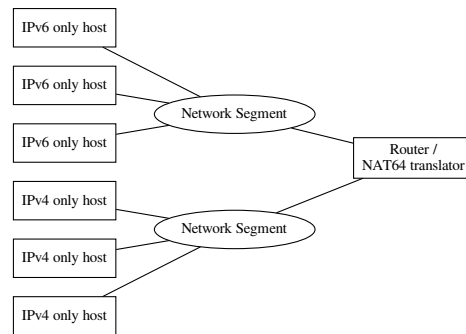


Figure 1.4: Separated IPv6 and IPv4 network segments

complexity, figure 1.4 shows how the complexity quickly grows even with a small number of hosts. The proposed in-network solution does not only ease the installation and deployment of IPv6, but it also allows line speed translation, because it is compiled into target dependent low level code that can run in ASICs [39], FPGAs [36] or even in software [10]. Figure 3.3 shows how the design differs for an in-network solution. Even on fast CPUs, software solutions like tayga [31] can be CPU bound (see section 4.4) and are incapable of translating protocols at line speed.

Chapter 2

Background

In this chapter we describe the key technologies involved and their relation to our work.

2.1 P4

P4 is a programming language designed to program inside network equipment. Its main features are protocol and target independence. The *protocol independence* refers to the separation of concerns in terms of language and protocols: P4, generally speaking, operates on bits that are parsed and then accessible in the self defined structures called headers. The general flow can be seen in figure 2.1: a parser parses the incoming packet and prepares it for processing in the switching logic. Afterwards the packets are output and deparsing of the parsed data might follow. In the context of NAT64 this is a very important feature: while the parser will read and parse in the ingress pipeline one protocol (f.i. IPv6), the deparser will output a different protocol (f.i. IPv4). The *target independence* is the second major feature of P4: it allows code

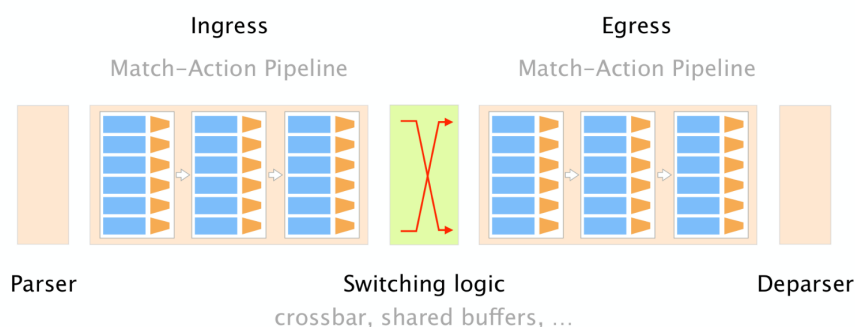


Figure 2.1: P4 protocol independence [60]

to be compiled to different targets. While in theory the P4 code should be completely target independent, in reality, there are some modifications needed on a per-target basis and each target faces different restrictions. The challenges arising from this are discussed in section 4.1. As opposed to general purpose programming languages, P4 lacks some features. Most notably loops, floating point operations and modulo operations. However within its constraints, P4 can guarantee operation at line speed, which general purpose programming languages cannot guarantee and also fail to achieve in reality (see section 4.4 for details).

2.2 IPv6, IPv4 and Ethernet

The first IPv6 RFC was published in 1998 [18]. Both IPv4 and IPv6 operate on layer 3 of the OSI model. In this thesis we only consider transmission via Ethernet, which operates at layer 2. Inside the Ethernet frame a field named “type” specifies the higher level protocol identifier.¹ This

¹0x0800 for IPv4 [25] and 0x86DD for IPv6 [16].

is important because Ethernet can only reference one protocol, which makes IPv4 and IPv6 mutually exclusive. In the figures 2.5 and 2.4 we show the packet headers of IPv4 and IPv6 for showing the in-protocol differences. The most notable differences between the two protocols for this thesis are:

- Different address lengths
 - IPv4: 32 bit
 - IPv6: 128 bit
- Lack of a checksum in IPv6
- Format of Pseudo headers (see section 2.5)

2.3 ARP and NDP, ICMP and ICMP6

While IPv6 and IPv4 are primarily used as a “shell” to support addressing for protocols that have no or limited addressing support (like TCP or UDP), protocols like ARP [42] and NDP [35] provide support for resolving IPv6 and IPv4 addresses to hardware (MAC) addresses. While both ARP and NDP are only used prior to establishing a connection and their results are cached, their availability is crucial for operating a switch, because without ARP or NDP no connection will every be established. Figure 2.2 illustrates a typical address resolution process. The major

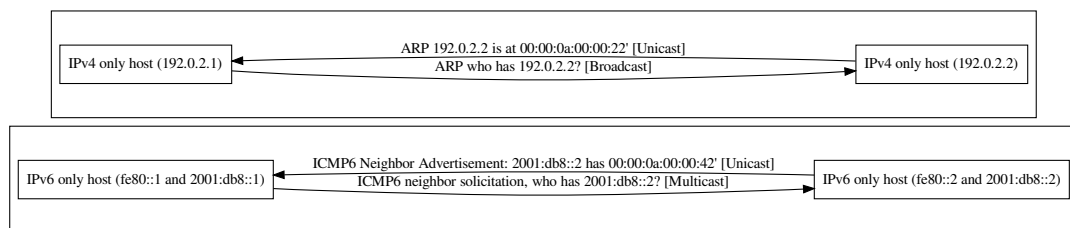


Figure 2.2: ARP and NDP

difference between ARP and NDP in relation to P4 are

- ARP is a separate protocol on the same layer as IPv6 and IPv4,
- NDP operates below ICMP6 which operates below IPv6,
- NDP contains checksums over payload,
- and NDP in ICMP6 contains optional, non-referenced option fields (specifically: ICMP6 link layer address option).

ARP is required to be a separate protocol, because IPv4 hosts don’t know how to communicate with each other yet, because they don’t have a way to communicate to the target IPv4 address (“The chicken and the egg problem”). NDP on the other hand already works within IPv6, as every IPv6 host is required to have a self-assigned link local IPv6 address from the IPv6 network `fe80::/10` (compare RFC4291 [24]). While ARP uses broadcasting for address resolution, NDP uses multicasting. IPv6 hosts automatically join multicast groups that embed parts of their IPv6 addresses [17], [66]. This way the collision domain is significantly reduced in IPv6, compared to IPv4.

As seen later in this document (compare section 4.3.1), the requirement to generate checksums over payload poses difficult problems for some hardware targets. Even more difficult is the use of options within ICMP6. The problem arises from the layout of the options, as seen in figure 2.3 and the following quote:

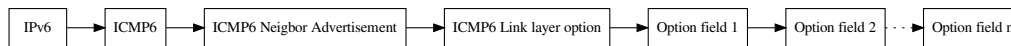


Figure 2.3: ICMP6 option fields

“Neighbor Discovery messages include zero or more options, some of which may appear multiple times in the same message. Options should be padded when necessary to ensure that they end on their natural 64-bit boundaries”.²

ICMP6 and ICMP are primarily used to signal errors in communication. Specifically, signalling that a packet is too big to pass a certain link and needs fragmentation is a common functionality of both protocols. For a host (or a switch) to be able to emit ICMP6 and ICMP messages, the host requires a valid IPv6 / IPv4 address. Without ICMP6 / ICMP support path MTU discovery [34], [32] does not work and the sender needs to determine different ways of finding out the maximum MTU on the path.

2.4 IPv6 Translation Mechanisms

While in this thesis we focus on NAT64 as a translation mechanism, there are a variety of different approaches, some of which we would like to portray here.

2.4.1 Stateless NAT64

Stateless NAT64 describes static mappings between IPv6 and IPv4 addresses. This can be based on longest prefix matching (LPM), ranges, bitmasks or individual entries.

NAT64 translations as described in this thesis modify multiple layers in the translation process:

- Ethernet (changing the type field)
- IPv4 / IPv6 (changing the protocol, changing the fields)
- TCP/UDP/ICMP/ICMP6 checksums

Figures 2.4 and 2.5 show the headers of IPv4 and IPv6. As can be seen in the diagrams not only are the addresses of different size, but fields have also been changed or removed when the version changed. Depending on the NAT64 translation direction, a translator will need to re-arrange fields to a different position, remove fields and add fields. This in turn causes the packet size for standard headers to differ by 160 Bit.³

2.4.2 Stateful NAT64

Stateful NAT64 as defined in RFC6146 [6] defines how to create 1:n mappings between IPv6 and IPv4 hosts. The motivation for stateful NAT64 is similar to stateful NAT44 [56]: while NAT44 allows translating many (private) IPv4 addresses to one (public) IPv4 address, NAT64 allows translating many IPv6 addresses to one IPv4 address. While the opposite stateful translation, mapping many IPv4 addresses to one IPv6 address, is also technically possible, the differences in address space size don't justify its use in general.

Stateful NAT64 in particular uses information in higher level protocols to multiplex connections: Given one IPv4 address and the TCP protocol, outgoing connections from IPv6 hosts can dynamically mapped to the range of possible TCP ports. After a session is closed, the port can be reused again. The selection of mapped ports is usually based on the availability on the IPv4 side and not related to the original port. To support stateful NAT64, the translator needs to store the mapping in a table and purge entries regularly.

Stateful NAT64 usually uses information found in protocols at layer 4 like TCP [46] or UDP [43]. However, it can also support ICMP [44] and ICMP6 [15].

²Quote from [35].

³IPv6: 320 Bit, IPv4 160 Bit

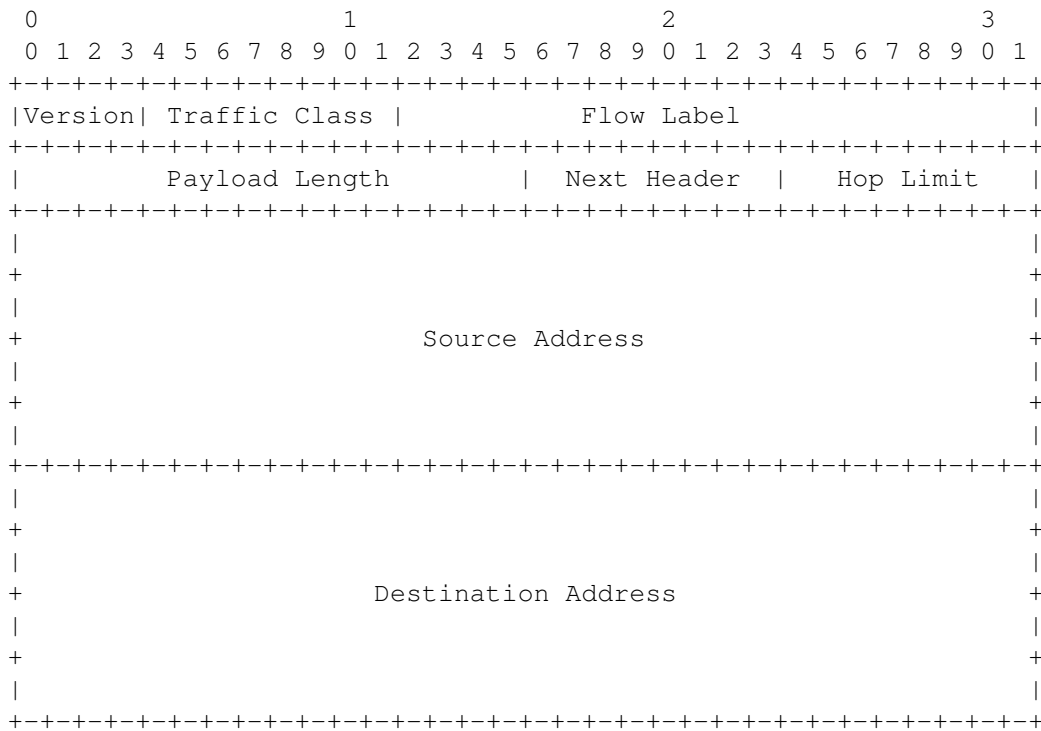


Figure 2.4: IPv6 Header [18]

2.4.3 Higher Layer Protocol Dependent Translation

Further translation can be achieved by using information in higher level protocols like HTTP [20] or TLS [9]. Application proxies like nginx [40] use layer 7 protocol information, like the requested hostname, to proxy towards backends.

Within this proxying method, the underlying IP protocol can be changed from IPv6 to IPv4 and vice versa. However, if using HTTPS with TLS 1.3 [48], the requested hostname that is usually used for selecting the backend can be encrypted, which poses a challenge for implementations. While protocol dependent translation has the highest amount of information to choose from for translation, complex parsers or even cryptographic methods are required for it. That reduces the opportunities of protocol dependent translations to run on devices with less sophisticated devices.

2.4.4 Mapping IPv4 Addresses in IPv6

As described in section 2.2, one of the major differences between IPv6 and IPv4 is the address length. As the whole IPv4 Internet can be represented in only 32 bits, it is a common practice to assign an IPv6 prefix for IPv6 hosts that represents a mapping to the whole IPv4 Internet. In RFC6052 [13] the well known prefix `64:ff9b::/96` is defined that can be used for this purpose. One possibility to map an IPv4 address into the prefix is by adding its integer value to the prefix, treating it as an offset. In figure 2.7 we show example python code of how this can be done. Network administrators can choose to use either the well known prefix or to use a network block of their own to map the Internet.⁴ While a /96 prefix seems a natural selection (it provides exactly 32 bit), other prefix lengths are defined in RFC6052 (see figure 2.8) that allow flexible embedding of the IPv4 address. RFC6146, which describes stateful NAT64, states that “IPv4 addresses of IPv4 hosts are algorithmically translated to and from IPv6 addresses by using the algorithm defined in [RFC6052]” [6] While this sentence does not use the typical RFC keywords like SHALL, REQUIRED, etc. [12], we interpret this sentence in the meaning of “a stateful NAT64 translator SHALL implement IPv4 address embedding as described in the algorithm of RFC6052”.

⁴For instance `2a0a:e5c0:0:1::/96` [59].

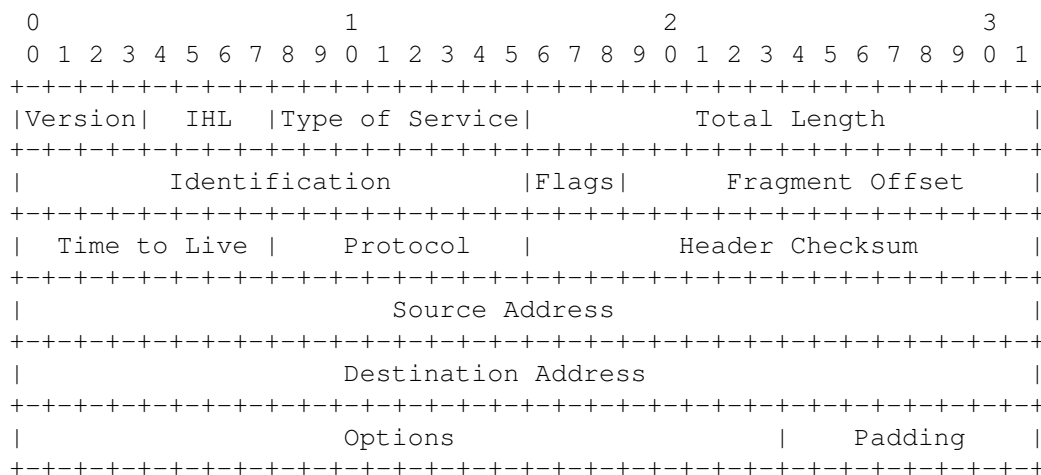


Figure 2.5: IPv4 Header [45]

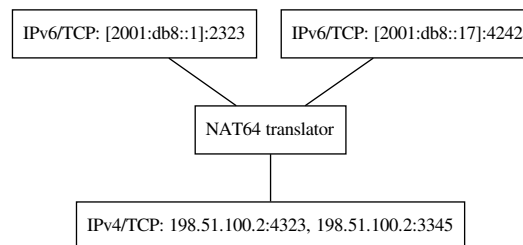


Figure 2.6: Stateful NAT64

2.4.5 DNS64

Tightly related to NAT64 is a technology known as DNS64 [28]. DNS64 tries to solve the problem of addressing IPv4 only hosts from IPv6 only hosts by adding a “fake” IPv6 (AAAA) DNS resource record, as shown in figure 2.9. The DNS64 DNS server will query the authoritative DNS server for an AAAA record. However as the host *ipv4onlyhost.example.com* is only reachable by IPv4, it also only has an A entry. After receiving the answer that there is no AAAA record, the DNS64 server will ask for an A record and gets an answer that the name *ipv4onlyhost.example.com* resolves to the IPv4 address *192.0.2.0*. The DNS64 server then embeds the IPv4 address in the configured IPv6 prefix (*64:ff9b::/96* in this case) and returns a fake AAAA record to the IPv6 only host (pointing to *64:ff9b::c000:200* in this case). The IPv6 only host then will use the address to connect to. The NAT64 translator recognises either that the address is part of a configured prefix or that it has a dedicated table entry for mapping this IPv6 address to an IPv4 address and translates it accordingly.

2.5 Protocol Checksums

One challenge for translating IPv6 to IPv4 are checksums of higher level protocols like TCP and UDP that incorporate information from the lower level protocols. The pseudo header for upper layer protocols for IPv6 is defined in RFC2460 [18] and shown in figure 2.10, the IPv4 pseudo header for TCP and UDP are defined in RFC768 and RFC793 and are shown in 2.11. When translating, the checksum fields in the higher protocols need to be adjusted. The checksums for TCP and UDP are calculated not only over the pseudo headers, but also contain the payload of the packet. This is important because some targets (like the NetFPGA) do not allow to access the payload (see section 3.3). The checksums for IPv4, TCP, UDP and ICMP6 are all based on

```
>>> import ipaddress
>>> prefix=ipaddress.IPv6Network("64:ff9b::/96")
>>> ipv4address=ipaddress.IPv4Address("192.0.2.0")
>>> int(ipv4address)
3221225984
>>> hex(3221225984)
'0xc0000200'
>>> prefix[int(ipv4address)]
IPv6Address('64:ff9b::c000:200')
```

Figure 2.7: Representing an IPv4 address in an IPv6 prefix

PL	0	32	40	48	56	64	72	80	88	96	104	
32	prefix	v4(32)				u	suffix					
40	prefix	v4(24)			u	(8)	suffix					
48	prefix	v4(16)		u	(16)	suffix						
56	prefix	(8)		u	v4(24)	suffix						
64	prefix	u		v4(32)	suffix							
96	prefix	v4(32)										

Figure 2.8: IPv4 embedding depending on the prefix length

the “Internet Checksum” [45], [11]. Its calculation can be summarised as follows:

The checksum field is the 16-bit one’s complement of the one’s complement sum of all 16-bit words in the header. For purposes of computing the checksum, the value of the checksum field is zero.⁵

2.6 Network Designs

In relation to IPv6 and IPv4, there are in general three different network designs possible: The oldest form are IPv4 only networks. These networks consist of hosts that are either not configured for IPv6 or are even technically incapable of enabling the IPv6 protocol. These nodes are connected to an IPv4 router that is connected to the Internet. That router might be capable of translating IPv4 to IPv6 and vice versa.

With the introduction of IPv6, hosts can have a separate IP stack active and in that configuration hosts are called “dualstack hosts”. Dualstack hosts are capable of reaching both IPv6 and IPv4 hosts directly without the need of any translation mechanism.

The last possible network design is based on IPv6 only hosts. While it is technically easy to disable IPv4, it seems that completely removing the IPv4 stack in current operating systems is not an easy task [58]. While the three network designs look similar, there are significant differences in operating them and limitations that are not easy to circumvent. In the following sections, we describe the limitations and reason how a translation mechanism like our NAT64 implementation should be deployed.

⁵Quote from Wikipedia [63].

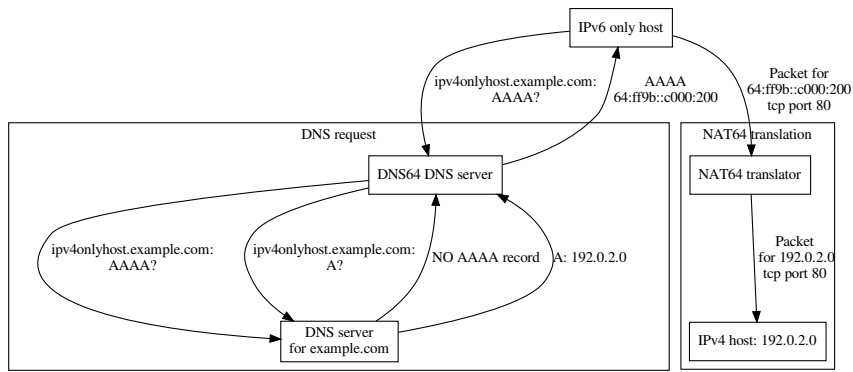


Figure 2.9: Illustration of DNS64

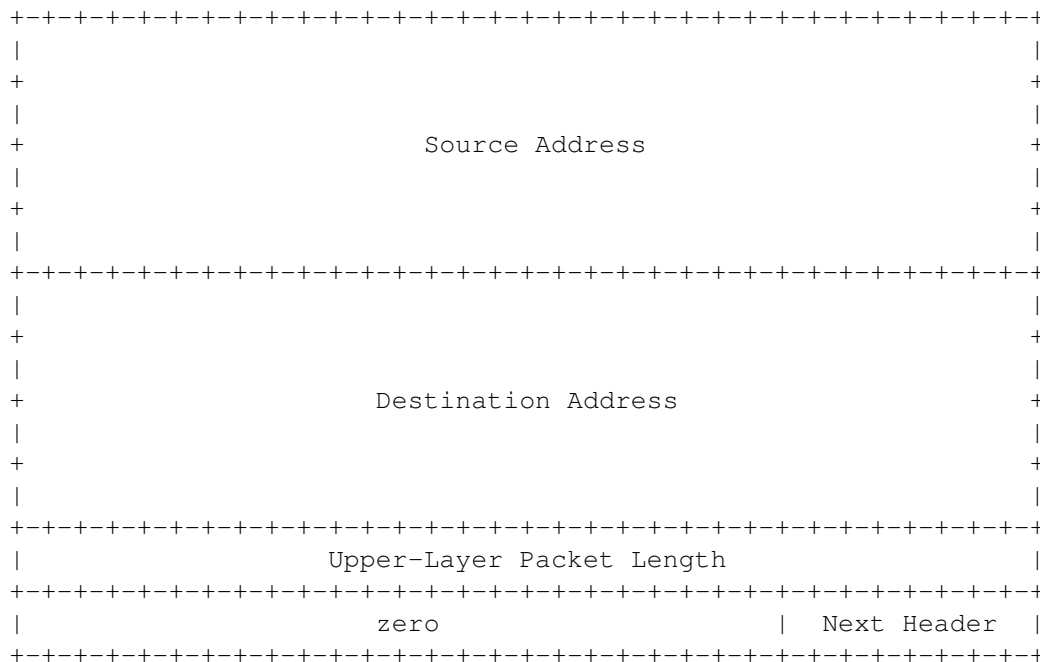


Figure 2.10: IPv6 Pseudo Header

2.6.1 IPv4 only network limitations

As shown in figures 2.5 and 2.4 the IPv4 address size is 32 bit, while the IPv6 address size is 128 bit. Without an extension to the address space, there is no protocol independent mapping of IPv4 address to IPv6⁶ that can cover the whole IPv6 address space. Thus IPv4 only hosts can never address every host in the IPv6 Internet. While protocol dependent translations can try to minimise the impact, accessing all IPv6 addresses independent of the protocol is not possible.

2.6.2 Dualstack network maintenance

While dualstack hosts can address any host in either IPv6 or IPv4 networks, the deployment of dualstack hosts comes with a major disadvantage: all network configurations double. The required routing tables double, the firewall rules roughly double⁷ and the number of network

⁶See section 2.2.

⁷The rule sets even for identical policies in IPv6 and IPv4 networks are not identical, but similar. For this reason we state that roughly double the amount of firewall rules are required for the same policy to be applied.

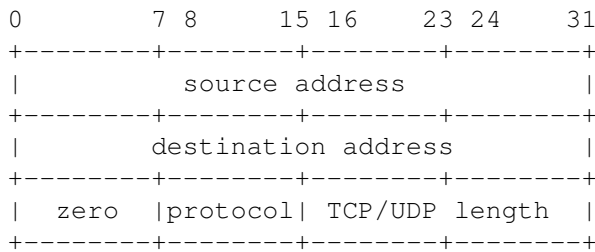


Figure 2.11: IPv4 Pseudo Header

supporting systems, (like DHCPv4, DHCPv6, router advertisement daemons, etc.) also roughly double. Additionally services that run on either IPv6 or IPv4 might need to be configured to run in dualstack mode as well and not every software might be capable of that. So while there is the instant benefit of not requiring any transition mechanism or translation method, we argue that the added complexity (and thus operational cost) of running dual stack networks can be significant.

2.6.3 IPv6 only networks

IPv6 only networks are in our opinion the best choice for long term deployments. The reasons for this are as follows: First of all hosts eventually will need to support IPv6 and secondly IPv6 hosts can address the whole 32 bit IPv4 Internet mapped in a single /96 IPv6 network. IPv6 only networks also allow the operators to focus on one IP stack.

2.7 NetFPGA

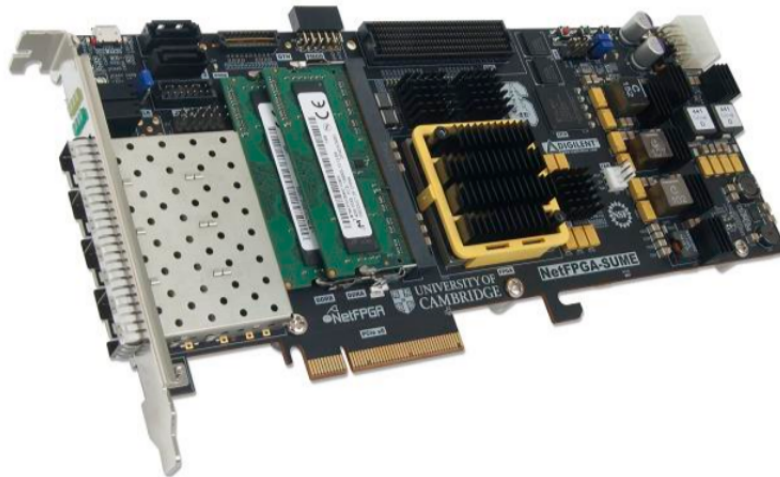


Figure 2.12: NetFPGA Board, [67]

The NetFPGA [67] is an FPGA card featuring four 10 Gbit/s SFP+ ports. It includes the Xilinx Virtex-7 690T FPGA on board, 27 MB of storage, allowing to save table data, and 8 GB of DDR3 RAM. The NetFPGA can be run inside a host (connected by PCI-E, gen 3) or as a standalone card.

It can be used as a “traditional” FPGA, with the focus on designing the logic. However, the NetFPGA also supports the P4 programming language [36] and thus abstracts away the low level logic by providing a higher level interface. For the purpose of this thesis we treat the NetFPGA as a standard P4 target, similar to other available P4 targets [39], [38], [37]. In particular, we

treat the NetFPGA as a P4 capable, four port 10 Gbit/s network switch that allows us to process packets at line speed.

Chapter 3

Design

In this chapter we describe the architecture of our solution and our design choices. We first introduce the general design of NAT64 in the P4 architecture. Afterwards we describe the design differences for the BMV2 and NetFPGA P4 architectures. Afterwards we discuss the design of stateless and stateful NAT64 in relation to P4 as well as two existing software NAT64 solutions. Lastly we discuss how we verify NAT64 functionality present and the network configurations that we use.

3.1 P4/NAT64

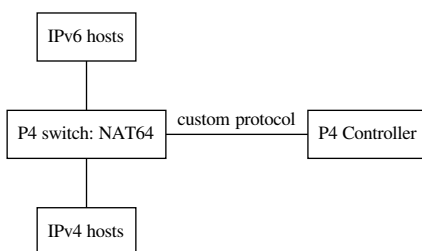


Figure 3.1: P4 Switch Architecture

In section 2.4 we discussed different translation mechanisms for IPv6 and IPv4. In this thesis we focus on the translation mechanisms stateless and stateful NAT64. While higher layer protocol dependent translations are more flexible, this topic has already been addressed in [55] and the focus in this thesis is on the practicability of high speed NAT64 with P4. The high level design can be seen in figure 3.1: a P4 capable switch is running our code to provide NAT64 functionality. A P4 switch cannot manage its tables on its own and needs support for this from a controller. The controller also has the role to handle unknown packets and can modify the runtime configuration of the switch. This is especially useful in the case of stateful NAT64. If only static table entries are required, they can usually be added at the start of a P4 switch and the controller can also be omitted. However stateful NAT64 requires the use of a controller to create session entries in the switch tables. The P4 switch can use any protocol to communicate with the controller, as the connection to the controller is implemented as a separate Ethernet port. Software NAT64 solutions typically require routing to be applied to transport the packet to the NAT64 translator as shown in figure 3.2.

Our design differs here: while routing could be used like described above, NAT64 with P4 does not require any routing to be setup. Figure 3.3 shows the network design that we realise using P4. This design has multiple advantages: first it reduces the number of devices to pass and thus directly reduces the RTT, secondly it allows translation of IP addresses within the same logic network segment. P4 switches in general look very similar to regular switches, however

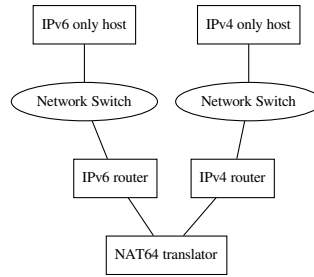


Figure 3.2: Standard NAT64 translation

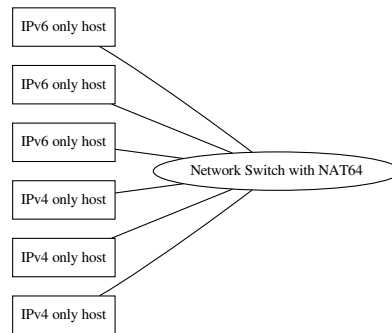


Figure 3.3: In-network NAT64 translation

support executing logic while the packet passes through the switch. Figure 3.4 illustrates how our solution is implemented and translates packets.

3.2 P4/BMV2

The software emulated switch that is implemented using Open vSwitch [21] and the behavioral model [10] offers the fastest and easiest way of P4 development. All NAT64 features are tested first on P4/BMV2 and in a second step ported to P4/NetFPGA and modified, where necessary. The development follows closely the general design shown in section 3.1. As outlined in section 2.5, checksums inside higher level protocols need to be adjusted after translation. Within the software emulation checksums can be computed with two different methods:

- Recalculating the checksum by inspecting headers and payload
- Calculating the difference between the translated headers

The BMV2 model is sophisticated and provides direct support for calculating the checksum over the payload. This allows the BMV2 model to operate as a full featured host, including advanced features like responding to ICMP6 Neighbor discovery requests [35] that include payload checksums. Sample code that calculates the required checksum for answering NDP queries is shown in figure 3.5. The code shows how the field `hdr.icmp6.checksum` is updated with the `csum16` method depending on the IPv6 and ICMP6 headers as well as the payload. The second option of using the differences is described in section 3.3.

3.3 P4/NetFPGA

While the P4-NetFPGA project [36] allows compiling P4 to the NetFPGA, the design slightly varies due to limitations in the available toolchain. In particular, the NetFPGA P4 compiler does

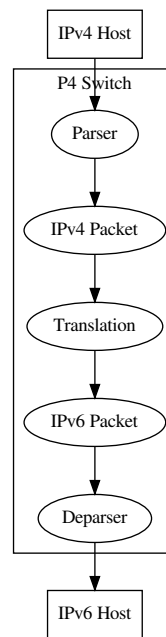


Figure 3.4: Our P4 Switch Architecture

not support reading the payload.¹ For this reason it also does not support creating the checksum based on the payload. To support checksum modifications in NAT64 on the NetFPGA, the checksum is calculated using differences between the IPv6 and IPv4 headers.

As the checksum calculation only depends on the 1-complement sums of headers and the payload (compare section 2.5) and only headers are modified during NAT64 translations, the higher level protocol checksums can be corrected based on the sum of differences of both headers. Thus our P4/NetFPGA implementation first calculates the sum of the relevant IPv4 headers (`v4sum()`), the sum of the relevant IPv6 headers (`v6sum()`) and then calculates the difference including a possible carry bit and adjusts the higher level protocol by this difference (`delta_tcp_from_v6_to_v4()`). Figure 3.6 shows an excerpt of the code used for adjust the checksum when translating TCP from IPv6 to IPv4. It is notable that not the full headers are used, but only a “pseudo header” is (compare figures 2.10 and 2.11).

3.4 Stateless NAT64

As seen in section 2.4.1, stateless NAT64 can be implemented using various factors. Our design for the stateless depends on the capabilities of the environment and is summarised in table 3.1. When using LPM for translating from IPv6 to IPv4, a /96 IPv6 network is configured for

Implementation	NAT64 match
P4/BMV2	LPM (both directions) and individual entries (both directions)
P4/NetFPGA	Individual entries
Tayga	LPM (IPv6 to IPv4) and individual entries (IPv4 to IPv6)
Jool	LPM (both directions)

Table 3.1: NAT64 match factors

covering the whole IPv4 Internet and the individual IPv4 address is appended to the prefix

¹This feature could be implemented in theory, but isn't available at the moment, see [53].

```

/* checksumming for icmp6_na_ns_option */
update_checksum_with_payload(meta.chk_icmp6_na_ns == 1,
{
    hdr.ipv6.src_addr,          /* 128 */
    hdr.ipv6.dst_addr,        /* 128 */
    meta.cast_length,         /* 32 */
    24w0,                      /* 24 0's */
    PROTO_ICMP6,              /* 8 */
    hdr.icmp6.type,           /* 8 */
    hdr.icmp6.code,           /* 8 */

    hdr.icmp6_na_ns.router,
    hdr.icmp6_na_ns.solicited,
    hdr.icmp6_na_ns.override,
    hdr.icmp6_na_ns.reserved,
    hdr.icmp6_na_ns.target_addr,

    hdr.icmp6_option_link_layer_addr.type,
    hdr.icmp6_option_link_layer_addr.ll_length,
    hdr.icmp6_option_link_layer_addr.mac_addr
},
hdr.icmp6.checksum,
HashAlgorithm.csum16
);

```

Figure 3.5: P4/BMV2 checksumming

(compare section 3.7). We also use LPM to match on an IPv4 sub network that translates to an IPv6 sub network. Individual entries are configured differently depending on the implementation: Limitations in the P4/NetFPGA environment require to use table entries. Jool supports individual entries as a special case of LPM, with a network mask matching only one IP address. Tayga supports LPM to translate from IPv6 to IPv4, but requires individual entries for translating from IPv4 to IPv6. Our P4/BMV2 offers the highest degree of flexibility, as it provides support for individual entries based on table entries and LPM table entries.

3.5 Stateful NAT64

Similar to stateless NAT64, the design of stateful NAT64 depends on the features of the individual implementation. As pointed out in section 2.4.2, stateful NAT64 is very similar to stateless NAT64, with the main difference being an additional stateful table that helps to create 1:n mappings. We use different approaches within the implementations to solve this problem:

- For P4/BMV2 and P4/NetFPGA a python controller handles packets that don't have a table entry, sets the table entry in the P4 switch and inserts the original packet afterwards back into the switch.
- With tayga we rely on the Linux kernel NAT44 capabilities
- Jool implements its own stateful mechanism based on port ranges

All methods though operate in a very similar fashion: A "controller" inspects the IPv6 packet and depending on the source address, destination address, protocol (TCP, UDP, ICMP, ICMP6, etc.) and the protocol ID (source / destination TCP/UDP port, ICMP identifier) it selects an outgoing IPv4 address, and source port or ICMP identifier. In case of Jool and Tayga this decision is based on a session table inside the Linux kernel, in case of P4 this decision is based on a session table inside the python controller. While the Jool and Tayga both support cleaning up old session entries, our P4 based solution does not support this feature at the moment.

```

action v4sum() {
    bit<16> tmp = 0;

    tmp = tmp + (bit<16>) hdr.ipv4.src_addr[15:0];           // 16 bit
    tmp = tmp + (bit<16>) hdr.ipv4.src_addr[31:16];         // 16 bit
    tmp = tmp + (bit<16>) hdr.ipv4.dst_addr[15:0];           // 16 bit
    tmp = tmp + (bit<16>) hdr.ipv4.dst_addr[31:16];         // 16 bit

    tmp = tmp + (bit<16>) hdr.ipv4.totalLen -20;             // 16 bit
    tmp = tmp + (bit<16>) hdr.ipv4.protocol;                 // 8 bit

    meta.v4sum = ~tmp;
}

/* analogue code for v6sum skipped */

action delta_tcp_from_v6_to_v4()
{
    v6sum();
    v4sum();

    bit<17> tmp = (bit<17>) hdr.tcp.checksum + (bit<17>) meta.v4sum;
    if (tmp[16:16] == 1) {
        tmp = tmp + 1;
        tmp[16:16] = 0;
    }
    tmp = tmp + (bit<17>) (0xffff - meta.v6sum);
    if (tmp[16:16] == 1) {
        tmp = tmp + 1;
        tmp[16:16] = 0;
    }

    hdr.tcp.checksum = (bit<16>) tmp;
}

```

Figure 3.6: Calculating checksum based on header differences

3.6 NAT64 Verification

We use socat [49] to verify basic operation of the NAT64 gateway and iperf [19] to test stability of the implementation and measure bandwidth. In particular we use the commands listed in table 3.2. The socat commands allow interactive testing on TCP and UDP connections, while the iperf commands fully utilise the available bandwidth with test data. The socat and iperf commands are used to verify all three NAT64 implementations (p4, tayga, jool).

3.7 IPv6 and IPv4 configuration

The following sections refer to host and network configurations. In this section we describe the IPv6 and IPv4 configurations as a basis for the discussion.

All IPv6 addresses are from the documentation block *2001:DB8::/32* [27]. In particular the following sub networks and IPv6 addresses are used:

We use private IPv4 addresses as specified by RFC1918 [47] from the 10.0.0.0/8 range as follows:

Command	Example	Description
socat - TCP6:HOST:PORT	socat - TCP6:[2001:db8:42::a00:2a]:2345	Connect via IPv6/TCP to IPv4 host
socat - UDP6:HOST:PORT	socat - UDP6:[2001:db8:42::a00:2a]:2345	Connect via IPv6/UDP to IPv4 host
socat - TCP:HOST:PORT	socat - TCP:10.0.1.42:2345	Connect via IPv4/TCP to IPv6 host
socat - UDP:HOST:PORT	socat - UDP:10.0.1.42:2345	Connect via IPv4/UDP to IPv6 host
socat - UDP6-LISTEN:PORT	socat - UDP6-LISTEN:2345	Listen on IPv6/UDP
socat - TCP6-LISTEN:PORT	socat - TCP6-LISTEN:2345	Listen on IPv6/TCP
socat - UDP-LISTEN:PORT	socat - UDP-LISTEN:2345	Listen on IPv4/UDP
socat - TCP-LISTEN:PORT	socat - TCP-LISTEN:2345	Listen on IPv4/TCP
iperf3 -PROTO -p PORT -B IP -s	iperf3 -4 -p 2345 -B 10.0.0.42 -s iperf3 -6 -p 2345 -B 2001:db8:42::42 -s	IPv4 iperf server IPv6 iperf server
iperf3 -PROTO -p PORT -O IGNORETIME -t RUNTIME -P PARALLEL -c IP iperf3 -PROTO -p PORT -O IGNORETIME -t RUNTIME -P PARALLEL -c IP -u -b0	iperf3 -6 -p 2345 -O 10 -t 190 -P20 -c 2001:db8:23::2a	Connect to iperf server Run for 190 seconds, skip first 10 seconds with 20 sessions connecting to 2001:db8:23::2a Same as above, but connect via UDP

Table 3.2: NAT64 verification commands

Address	Description
2001:db8:42::/64	IPv6 host network
2001:db8:23::/96	IPv6 mapping to the IPv4 Internet
2001:db8:42::42	IPv6 host address
2001:db8:42::77	IPv6 router address
2001:db8:42::a00:2a	In-network IPv6 address mapped to 10.0.0.42 (p4)
2001:db8:23::a00:2a	IPv6 address mapped to 10.0.0.42 (tayga)
2001:db8:23::2a	IPv6 address mapped to 10.0.0.42 (jool)

Table 3.3: IPv6 address and network overview

Address	Description
10.0.0.0/24	IPv4 host network
10.0.1.0/24	IPv4 network mapping to IPv6
10.0.0.77	IPv4 router address
10.0.0.66	In-network IPv4 address mapped to 2001:db8:42::42 (p4)
10.0.1.42	IPv4 address mapped to 2001:db8:42::42 (tayga)
10.0.1.66	IPv4 address mapped to 2001:db8:42::42 (jool)

Table 3.4: IPv4 address and network overview

Chapter 4

Results

This section describes the achieved results and compares the P4 based implementation with real world software solutions.

We distinguish the software implementation of P4 (BMV2) and the hardware implementation (NetFPGA) due to significant differences in deployment and development. We present benchmarks for the existing software solutions as well as for our hardware implementation. As the objective of this thesis is to demonstrate the high speed capabilities of NAT64 in hardware, no benchmarks were performed on the P4 software implementation.

4.1 P4 based implementations

We successfully implemented P4 code to realise NAT64 [54]. It contains parsers for all related protocols (IPv6, IPv4, UDP, TCP, ICMP, ICMP6, NDP, ARP), supports EAMT as defined by RFC7757 [3], and is feature equivalent to the two compared software solutions *tayga* [31] and *jool* [33]. Due to limitations in the P4 environment of the NetFPGA environment, the BMV2 implementation is more feature rich.

For this thesis the parsing capabilities of P4 were adequate. However P4, at the time of writing, cannot parse ICMP6 options in general, as the upper level protocol does not specify the number of options that follow and parsing of an undefined number of 64 bit blocks is required, which P4 does not support.

The language has some limitations on the placement of conditional statements (*if/switch*).¹ Furthermore P4/BMV2 does not support for multiple LPM keys in a table, however it supports multiple keys with ternary matching, which is a superset of LPM matching.

When developing P4 programs, the reason for incorrect behaviour we have seen were checksum problems. This is in retrospective expected, as the main task of our implementation is modifying headers on which the checksums depend. In all cases we have seen Ethernet frame checksum errors, the effective length of the packet was incorrect.

The tooling around P4 is somewhat fragile. We encountered small language bugs during the development [52], (compare section B.4) or found missing features [51], [57]: it is at the moment impossible to retrieve the matching key from table or the name of the action called. Thus if different table entries call the same action, it is impossible within the action, or if forwarded to the controller, within the controller to distinguish on which match the action was triggered. This problem is very consistent within P4, not even the matching table name can be retrieved. While these information can be added manually as additional fields in the table entries, we would expect a language to support reading and forwarding this kind of meta information.

While in P4 the P4 code and the related controller are tightly coupled, their data definitions are not. Thus the packet format definition that is used between the P4 switch and the controller has to be duplicated. Our experiences in software development indicate that this duplication is a likely source of errors in bigger software projects.

The supporting scripts in the P4 toolchain are usually written in python2. However, python2 “is legacy” [61]. During development errors with unicode string handling in python2 caused changes to IPv6 addresses.

¹In general, if and switch statements in actions lead to errors, but not all constellations are forbidden.

4.2 P4/BMV2

The software implementation of P4 has most features, which is mostly due to the capability of creating checksums over the payload. It enables the switch to act as a “proper” participant in NDP, as this requires the host to calculate checksums over the payload. Table 4.1 references all implemented features. The switch responds to ICMP echo requests, ICMP6 echo requests,

Feature	Description	Status
Switch to controller	Switch forwards unhandled packets to controller	fully implemented ^a
Controller to Switch	Controller can setup table entries	fully implemented ^b
NDP	Switch responds to ICMP6 neighbor solicitation request (without controller)	fully implemented ^c
ARP	Switch can answer ARP request (without controller)	fully implemented ^d
ICMP6	Switch responds to ICMP6 echo request (without controller)	fully implemented ^e
ICMP	Switch responds to ICMP echo request (without controller)	fully implemented ^f
NAT64: TCP	Switch translates TCP with checksumming from/to IPv6 to/from IPv4	fully implemented ^g
NAT64: UDP	Switch translates UDP with checksumming from/to IPv6 to/from IPv4	fully implemented ^h
NAT64: ICMP/ICMP6	Switch translates echo request/reply from/to ICMP6 to/from ICMP with checksumming	fully implemented ⁱ
NAT64: Sessions	Switch and controller create 1:n sessions/mappings	fully implemented ^j
Delta Checksum	Switch can calculate checksum without payload inspection	fully implemented ^k
Payload Checksum	Switch can calculate checksum with payload inspection	fully implemented ^l

^aSource code: `actions_egress.p4`

^bSource code: `controller.py`

^cSource code: `actions_icmp6_ndp_icmp.p4`

^dSource code: `actions_arp.p4`

^eSource code: `actions_icmp6_ndp_icmp.p4`

^fSource code: `actions_icmp6_ndp_icmp.p4`

^gSource code: `actions_nat64_generic_icmp.p4`

^hSource code: `actions_nat64_generic_icmp.p4`

ⁱSource code: `actions_nat64_generic_icmp.p4`

^jSource code: `actions_nat64_session.p4, controller.py`

^kSource code: `actions_delta_checksum.p4`

^lSource code: `checksum_bmv2.p4`

Table 4.1: P4/BMV2 feature list

answers NDP and ARP requests. Overall P4/BMV is very easy to use, even without a controller a fully functional network host can be implemented.

This P4/BMV implementation supports translating ICMP/ICMP6 echo request and echo reply messages, but does not support all ICMP/ICMP6 translations that are defined in RFC6145 [30].

4.3 P4/NetFPGA

In the following section we describe the achieved feature set of P4/NetFPGA in detail and analyse differences to the BMV2 based implementation.

4.3.1 Features

While the NetFPGA target supports P4, compared to P4/BMV2 we only implemented a reduced features set on P4/NetFPGA. The first reason for this is missing support of the NetFPGA P4 compiler to inspect payload and to compute checksums over payload. While this can (partially) be compensated using delta checksums, the compile time of 2 to 6 hours contributed to a significant slower development cycle compared to BMV2. Lastly, the focus of this thesis is to implement high speed NAT64 on P4, which only requires a subset of the features that we

realised on BMV2. Table 4.2 summarises the implemented features and reasons about their implementation status.

Feature	Description	Status
Switch to controller	Switch forwards unhandled packets to controller	portable ^a
Controller to Switch	Controller can setup table entries	portable ^b
NDP	Switch responds to ICMP6 neighbor solicitation request (without controller)	portable ^c
ARP	Switch can answer ARP request (without controller)	portable ^d
ICMP6	Switch responds to ICMP6 echo request (without controller)	portable ^e
ICMP	Switch responds to ICMP echo request (without controller)	portable ^f
NAT64: TCP	Switch translates TCP with checksumming from/to IPv6 to/from IPv4	fully implemented ^g
NAT64: UDP	Switch translates UDP with checksumming from/to IPv6 to/from IPv4	fully implemented ^h
NAT64: ICMP/ICMP6	Switch translates echo request/reply from/to ICMP6 to/from ICMP with checksumming	portable ⁱ
NAT64: Sessions	Switch and controller create 1:n sessions/mappings	portable ^j
Delta Checksum	Switch can calculate checksum without payload inspection	fully implemented ^k
Payload Checksum	Switch can calculate checksum with payload inspection	unsupported ^l

^aWhile the NetFPGA P4 implementation does not have the clone3() extern that the BMV2 implementation offers, communication to the controller can easily be realised by using one of the additional ports of the NetFPGA and connect a physical network card to it.

^bThe p4utils suite offers an easy access to the switch tables. While the P4-NetFPGA support repository also offers python scripts to modify the switch tables, the code is less sophisticated and more fragile.

^cNetFPGA/P4 does not offer calculating the checksum over the payload. However delta checksumming can be used to create the required checksum for replying.

^dAs ARP does not use checksums, integrating the source code `actions_arp.p4` into the netpfga code base is enough to enable ARP support in the NetPFGA.

^eSame reasoning as NDP.

^fSame reasoning as NDP.

^gSource code: `actions_nat64_generic_icmp.p4`

^hSource code: `actions_nat64_generic_icmp.p4`

ⁱICMP/ICMP6 translations only require enabling the `icmp/icmp6` code in the netpfga code base.

^jSame reasoning as "Controller to switch".

^kSource code: `actions_delta_checksum.p4`

^lTo support creating payload checksums, either an HDL module needs to be created or to modify the generated the PX program. [53]

Table 4.2: P4/NetFPGA feature list

4.3.2 Stability

Two different NetPFPGA cards were used during the development of this thesis. The first card had consistent ioctl errors (compare section B.4) when writing table entries. The available hardware tests (compare figures 4.1 and 4.2) showed failures in both cards, however the first card reported an additional "10G_Loopback" failure. Due to the inability of setting table entries, no benchmarking was performed on the first NetFPGA card. During the development and benchmarking, the second NetFPGA card stopped to function properly multiple times. In these cases the card would not forward packets anymore. Multiple reboots (up to 3) and multiple times re-flashing the bitstream to the NetFPGA usually restored the intended behaviour. However due to this "crashes", it was impossible for us to run a benchmark for more than one hour. Similarly, sometimes flashing the bitstream to the NetFPGA would fail. It was required to reboot the host containing the NetFPGA card up to 3 times to enable successful flashing.²

²Typical output of the flashing process would be: "fpga configuration failed. DONE PIN is not HIGH"

Run Auto Test	TestID	Result	Description
DDR3B_RW	DDR3B_RW	Passed	Read/Write on DDR3B SODIMM
DDR3B_IIC	DDR3B_IIC	Passed	IIC R/W on DDR3B SODIMM
Test DDR3B	DDR3A_RW	Failed	Read/Write on DDR3A SODIMM
Test DDR3A	DDR3A_IIC	Passed	IIC R/W on DDR3A SODIMM
Test CPLD and FLash	CPLD	Passed	CPLD, Flash and Configuration
Test GPIO	FMC_Clocks	Failed	Clock Signals on FMC Connector
Test FMC	SD_Card	Failed	SD Card (4-bit SDIO)
Test QDRII+ A	GPIO_Test	Failed	GPIO Walking 1/0 on FMC and Pmod
Test QDRII+ C	FMC	Failed	FMC Connector GTH Transceiver (12.5Gbps) Lo...
Test QDRII+ B	QDRA_RW	Passed	QDR II+ A Read/Write
Test PCI-E Gen3 x8	QDRC_RW	Passed	QDR II+ C Read/Write
Test 10G Loopback	QDRB_RW	Passed	QDR II+ B Read/Write
Test SATA III	PCIE	Failed	PCI-Express Gen3 (8Gbps) Loopback
Test QTH	10G_Loopback	Failed	10G Ethernet Loopback
	SATA	Failed	SATA III (6Gbps) Loopback
	QTH	Failed	QTH Connector GTH Transceiver (12.5Gbps) Lo...

Figure 4.1: Hardware Test NetFPGA card 1

Run Auto Test	TestID	Result	Description
Show Test Summary	DDR3B_...	Passed	Read/Write on DDR3B SODIMM
Test DDR3B	DDR3B_I...	Passed	IIC R/W on DDR3B SODIMM
Test DDR3A	DDR3A_...	Failed	Read/Write on DDR3A SODIMM
Test CPLD and FLash	DDR3A_I...	Passed	IIC R/W on DDR3A SODIMM
Test GPIO	CPLD	Passed	CPLD, Flash and Configuration
Test FMC	FMC_Clo...	not tested	Clock Signals on FMC Connector
Test QDRII+ A	SD_Card	not tested	SD Card (4-bit SDIO)
Test QDRII+ C	GPIO_Test	not tested	GPIO Walking 1/0 on FMC and Pmod
Test QDRII+ B	FMC	not tested	FMC Connector GTH Transceiver (12.5Gbps) Loopback
Test PCI-E Gen3 x8	QDRA_RW	Passed	QDR II+ A Read/Write
Test 10G Loopback	QDRC_RW	Passed	QDR II+ C Read/Write
Test SATA III	QDRB_RW	Passed	QDR II+ B Read/Write
Test QTH	PCIE	not tested	PCI-Express Gen3 (8Gbps) Loopback
	10G_Loo...	Passed	10G Ethernet Loopback
	SATA	not tested	SATA III (6Gbps) Loopback
	QTH	not tested	QTH Connector GTH Transceiver (12.5Gbps) Loopback

Figure 4.2: Hardware Test NetFPGA card 2, [23]

Performance

The NetFPGA card performed at near line speed and offers NAT64 translations at 9.28 Gbit/s (see section 4.5 for details). Single and multiple streams performed almost exactly identical and have been consistent through multiple iterations of the benchmarks.

4.3.3 Usability

The handling and usability of the NetFPGA card is rather difficult. In this section we describe our findings and experiences with the card and its toolchain.

To use the NetFPGA, the tools Vivado and SDNET provided by Xilinx need to be installed. However a bug in the installer triggers an infinite loop, if a certain shared library³ is missing on the target operating system. The installation program seems to be still progressing, however never finishes.

While the NetFPGA card supports P4, the toolchains and supporting scripts are in an immature

³The required shared library is libncurses5.

state. The compilation process consists of at least 9 different steps, which are interdependent.⁴ Some of the steps generate shell scripts and python scripts that in turn generate JSON data.⁵ However incorrect parsing generates syntactically incorrect scripts or scripts that generate incorrect output. The toolchain provided by the NetFPGA-P4 repository contains more than 80000 lines of code. The supporting scripts for setting table entries require setting the parameters for all possible actions, not only for the selected action. Supplying only the required parameters results in a crash of the supporting script.

The documentation for using the NetFPGA-P4 repository is very distributed and does not contain a reference on how to use the tools. Mapping of egress ports and their metadata field are found in a python script that is used for generating test data.

The compile process can take up to 6 hours and because the different steps are interdependent, errors in a previous stage were in our experiences detected hours after they happened. The resulting log files of the compilation process can be up to 5 MB in size. Within this log file various commands output references to other logfiles, however the referenced logfiles do not exist before or after the compile process.

During the compile process various informational, warning and error messages are printed. However some informational messages constitute critical errors, while on the other hand critical errors and syntax errors often do not constitute a critical error.⁶ Also contradicting output is generated.⁷

Programs or scripts that are called during the compile process do not necessarily exit non zero if they encountered a critical error. Thus finding the source of an error can be difficult due to the compile process continuing after critical errors occurred. Not only programs that have critical errors exit “successfully”, but also python scripts that encounter critical paths don’t abort with `raise()`, but print an error message to `stdout` and don’t abort with an error.

The most often encountered critical compile error is “Run ‘impl_1’ has not been launched. Unable to open”. This error indicates that something in the previous compile steps failed and can refer to incorrectly generated testdata to unsupported LPM tables.

The NetFPGA kernel module provides access to virtual Linux devices (`nf0...nf3`). However `tcpdump` does not see any packets that are emitted from the switch. The only possibility to capture packets that are emitted from the switch is by connecting a physical cable to the port and capturing on the other side.

Jumbo frames⁸ are commonly used in 10 Gbit/s networks. According to [65], even many gigabit network interface card support jumbo frames. However according to emails on the private NetFPGA mailing list, the NetFPGA only supports 1500 byte frames at the moment and additional work is required to implement support for bigger frames.

Our P4 source code required contains Xilinx annotations⁹ that define the maximum packet size in bits. We observed two different errors on the output packet, if the incoming packets exceed the specified size:

- The output packet is longer than the original packet.
- The output packet is corrupted.

While most of the P4 language is supported on the `netpfga`, some key techniques are missing or not supported.

- Analysing / accessing payload is not supported
- Checksum computation over payload is not supported

⁴See source code `bin/do-all-steps.sh`.

⁵One compilation step calls the script “`config_writes.py`”. This script failed with a syntax error, as it contained incomplete python code. The scripts `config_writes.py` and `config_writes.sh` are generated by `gen_config_writes.py`. The output of the script `gen_config_writes.py` depends on the content of `config_writes.txt`. That file is generated by the simulation “`xsim`”. The file “`SimpleSumeSwitch_tb.sv`” contains code that is responsible for writing `config_writes.txt` and uses a function named `axi4_lite_master_write_request_control` for generating the output. This in turn is dependent on the output of a script named `gen_testdata.py`.

⁶F.i. “CRITICAL WARNING: [BD 41-737] Cannot set the parameter `TRANSLATION_MODE` on `/axi_interconnect_0`. It is read-only.” is a non critical warning.

⁷While using version 2018.2, the following message was printed: “WARNING: command ‘`get_user_parameter`’ will be removed in the 2015.3 release, use ‘`get_user_parameters`’ instead”.

⁸Frames with an MTU greater than 1500 bytes.

⁹F.i. “`@Xilinx_MaxPacketRegion(1024)`”

- Using LPM tables can lead to compilation errors
- Depending on the match type, only certain table sizes are allowed

Renaming variables in the declaration of the parser or deparser lead to compilation errors. Function syntax is not supported. For this reason our implementation uses `#define` statements instead of functions.

4.4 Software based NAT64

Both solutions Tayga and Jool worked flawlessly. However as expected, both solutions are CPU bound. Under high load scenarios both solutions utilise one core fully. Neither Tayga as a user space program nor Jool as a kernel module implement multi threading.

4.5 NAT64 Benchmarks

In this section we give an overview of the benchmark design and summarise the benchmarking results.

4.5.1 Benchmark Design

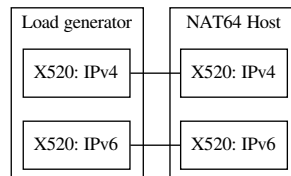


Figure 4.3: Benchmark design for NAT64 in software implementations

We use two hosts for performing benchmarks: a load generator and a NAT64 translator. Both hosts are equipped with a dual port Intel X520 10 Gbit/s network card. Both hosts are connected using DAC without any equipment in between. TCP offloading is enabled in the X520 cards. Figure 4.3 shows the network setup. When testing the NetPFGA/P4 performance, the X520 cards in the NAT64 translator are disconnected and instead the NetPFGA ports are connected, as shown in figure 4.4. The load generator is equipped with a quad core CPU (Intel(R) Core(TM) i7-6700 CPU @ 3.40GHz), enabled with hyperthreading and 16 GB RAM. The NAT64 translator is also equipped with a quad core CPU (Intel(R) Core(TM) i7-4770 CPU @ 3.40GHz) and 16 GB RAM. The first 10 seconds of the benchmark are excluded to avoid the TCP warm up phase.¹⁰

4.5.2 Benchmark Summary

Overall **tayga** has shown to be the slowest translator with an achieved bandwidth of **about 3 Gbit/s**, followed by **Jool** that translates at about **8 Gbit/s**. **Our solution** is the fastest with an almost line rate translation speed of about **9 Gbit/s**.

The TCP based benchmarks show realistic numbers, while iperf reports above line rate speeds (up to 22 gbit/s on a 10gbit/s link) for UDP based benchmarks. For this reason we have summarised the UDP based benchmarks with their average loss instead of listing the bandwidth details. The “adjusted bandwidth” in the UDP benchmarks incorporates the packets loss (compare tables 4.6 and 4.6).

¹⁰iperf -O 10 parameter, see section 3.6.

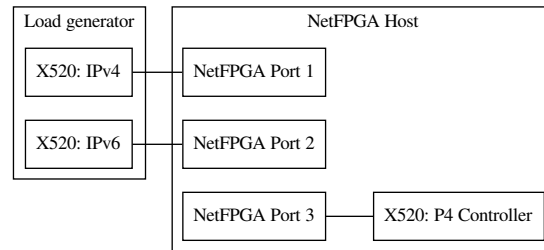


Figure 4.4: NAT64 with NetFPGA benchmark

Both software solutions showed significant loss of packets in the UDP based benchmarks (tayga: up to 91%, jool up to 71%), while the P4/NetFPGA showed a maximum of 0.01% packet loss. Packet loss is only recorded by iperf for UDP based benchmarks, as TCP packets are confirmed and resent if necessary.

Tayga has the highest variation of results, which might be due to being fully CPU bound even in the simplest benchmark. Jool has less variation and in general the P4/NetFPGA solution behaves almost identical in different benchmark runs.

The CPU load for TCP based benchmarks with Jool was almost negligible, however for UDP based benchmarks one core was almost 100% utilised. In all benchmarks with tayga, one CPU was fully utilised. And as the translation for P4/NetFPGA happens within the NetFPGA card, there was no CPU utilisation visible on the NAT64 host.

We see lower bandwidth for translating IPv4 to IPv6 in all solutions. We suspect that this might be due to slightly increasing packet sizes that occur during this direction of translation. Not only does this vary the IPv4 versus IPv6 bandwidth, but it might also cause fragmentation to slowdown the process.

During the benchmark with 1 and 10 parallel connections, no significant CPU load was registered on the load generator. However with 20 parallel connections, each of the two iperf processes¹¹ partially spiked to 100% cpu usage, which with 50 parallel connections the cpu load of each process hit 100% often.

While tayga's performance seems to reduce with the growing number of parallel connections, both Jool and our P4/NetFPGA implementations vary only slightly.

Overall the performance of tayga, a Linux user space program, is as expected. We were surprised about the good performance of Jool, which, while slower than the P4/NetFPGA solution, is almost on par with our solution.

¹¹One for sending, one for receiving.

4.5.3 IPv6 to IPv4 TCP Benchmark Results

Implementation	min/avg/max in Gbit/s			
Tayga	2.79 / 3.20 / 3.43	3.34 / 3.36 / 3.38	2.57 / 3.02 / 3.27	2.35 / 2.91 / 3.20
Jool	8.22 / 8.22 / 8.22	8.21 / 8.21 / 8.22	8.21 / 8.23 / 8.25	8.21 / 8.23 / 8.25
P4 / NetPFGA	9.28 / 9.28 / 9.29	9.28 / 9.28 / 9.29	9.28 / 9.28 / 9.29	9.28 / 9.28 / 9.29
Parallel connections	1	10	20	50

Table 4.3: IPv6 to IPv4 TCP NAT64 Benchmark

4.5.4 IPv4 to IPv6 TCP Benchmark Results

Implementation	min/avg/max in Gbit/s			
Tayga	2.90 / 3.15 / 3.34	2.87 / 3.01 / 3.22	2.68 / 2.85 / 3.09	2.60 / 2.78 / 2.88
Jool	7.18 / 7.56 / 8.24	7.97 / 8.05 / 8.09	8.05 / 8.08 / 8.10	8.10 / 8.12 / 8.13
P4 / NetPFGA	8.51 / 8.53 / 8.55	9.28 / 9.28 / 9.29	9.29 / 9.29 / 9.29	9.28 / 9.28 / 9.29
Parallel connections	1	10	20	50

Table 4.4: IPv4 to IPv6 TCP NAT64 Benchmark

4.5.5 IPv6 to IPv4 UDP Benchmark Results

Implementation	avg bandwidth in gbit/s / avg loss / adjusted bandwidth			
Tayga	8.02 / 70% / 2.43	9.39 / 79% / 1.97	15.43 / 86% / 2.11	19.27 / 91% / 1.73
Jool	6.44 / 0% / 6.41	6.37 / 2% / 6.25	16.13 / 64% / 5.75	20.83 / 71% / 6.04
P4 / NetPFGA	8.28 / 0% / 8.28	9.26 / 0% / 9.26	16.15 / 0% / 16.15	15.8 / 0% / 15.8
Parallel connections	1	10	20	50

Table 4.5: IPv6 to IPv4 UDP NAT64 Benchmark

4.5.6 IPv4 to IPv6 UDP Benchmark Results

Implementation	avg bandwidth in gbit/s / avg loss / adjusted bandwidth			
Tayga	6.78 / 84% / 1.06	9.58 / 90% / 0.96	15.67 / 91% / 1.41	20.77 / 95% / 1.04
Jool	4.53 / 0% / 4.53	4.49 / 0% / 4.49	13.26 / 0% / 13.26	22.57 / 0% / 22.57
P4 / NetPFGA	7.04 / 0% / 7.04	9.58 / 0% / 9.58	9.78 / 0% / 9.78	14.37 / 0% / 14.37
Parallel connections	1	10	20	50

Table 4.6: IPv4 to IPv6 UDP NAT64 Benchmark

Chapter 5

Conclusion and Outlook

The objective of implementing high speed NAT64 in P4 has been achieved.

Our implementation has been shown to be portable between 2 different P4 targets and we expect it to be portable to other P4 targets, potentially at much higher speeds. Our in-network solution allows novel translations without involving external routers, without involving external routers.¹ We expect this to support migration to IPv6 only networks.

P4 has been proven to us as a suitable programming language for network equipment with great potential. However in the current state the tooling and frameworks are still immature and need significant work to become usable for solving day-to-day challenges or supporting large scale projects. Even with the current state drawbacks, P4 is a very convincing language that has wide range of applications due to its protocol independence and easy to understand architecture. The availability of protocol independent programmable network equipment opens up many possibilities in network programming. While this thesis focused on NAT64, the accompanying technology DNS64 [7] could also be implemented in P4, thus completing the translation mechanism.

In our opinion, the P4/NetFPGA platform is a good showcase for the capabilities of P4, demonstrating near line speed P4 programs with good capabilities of demonstrating scientific research. However, the supporting code toolchain shows strong weaknesses that render productive deployments difficult.

While the project concluded successfully, there are a variety of possible improvements to our work as well to the used toolchains.

The implementation of our algorithm uses the IPv4-Compatible IPv6 Address [24] to embed IPv4 addresses. However RFC6052 [8] defines different embeddings depending on the prefix size. A future version should not only support more flexible embeddings, but also consider more in depth translation like ICMP/ICMP6 specifics.

The P4 language has shown maturity, but the usability and ease of use of the provided toolchains can be significantly improved. Additionally, we envision a stronger tie between the different tools in the P4 environment, like a collection of libraries and modules that could form something on the line of a “P4OS”. This operating system could spawn over network switches with P4, provide a coherent library and define data definitions that can be used in various programming languages bindings.

The NetFPGA, from the hardware point of view, is a very interesting hardware platform. Reducing the difficulties we experienced with the surrounding toolchain and making development experience more consistent has the potential to not only make NetFPGA, but also the whole set of P4 hardware more interesting for developers.

¹ Compare figures 3.2 and 3.3.

Appendix A

Resources and code repositories

The following sections describe how to acquire the resources to reproduce the test results.

A.1 Operating Systems

All P4 compilations were made on Ubuntu 16.04 with kernels

- 4.15.0-54-generic (Supporting Desktop)
- 4.4.0-143-generic (BMV2 test VM)
- 4.15.0-55-generic (Desktop with NetFPGA card)

A.2 Master Thesis

The master thesis including all self developed source code is available by git via

- `git clone git@gitlab.ethz.ch:nicosc/master-thesis.git`
- `git clone git@gitlab.ethz.ch:nsg/student-projects/ma-2019-19_high_speed_nat64_with_p4`

It can be browsed online on <https://gitlab.ethz.ch/nicosc/master-thesis> and on https://gitlab.ethz.ch/nsg/student-projects/ma-2019-19_high_speed_nat64_with_p4.

A.3 Xilinx Toolchain

A prerequisite for building the NetFPGA source code is the installation of

- Xilinx_SDNNet_2018.2_1005_9
- Xilinx_Vivado_SDK_20182_0614_1954

Both tools need to be installed to `/opt/Xilinx/`, as paths are hardcoded in various places.

A.4 P4/NetFPGA support scripts

To be able to compile P4 source code to the NetFPGA the collection of scripts, Makefiles and sample code of P4-NetFPGA is required. The repository `git@github.com:NetFPGA/P4-NetFPGA-live.git` needs to be cloned to “projects” subdirectory as “P4-NetPFGA” of the user that wants to compile the source code. Access to the repository is granted after applying for access as described on <https://github.com/NetFPGA/P4-NetFPGA-public/wiki>. After that the variable `P4_PROJECT_NAME` in `/projects/P4-NetFPGA/tools/settings.sh` needs to be modified to read `export P4_PROJECT_NAME=minip4` instead of `export P4_PROJECT_NAME=switch_calc`. Sample code for installation:

```
mkdir -p ~/projects
git clone git@github.com:NetFPGA/P4-NetFPGA-live.git P4-NetFPGA
sed -i 's/\(P4_PROJECT_NAME=\)\.*/\lminip4/' ~/projects/P4-NetFPGA/tools/settings.sh
```

Version **v1.3.1-46-g97d3aaa** of the P4-NetFPGA repository was used for creating the bitfiles of this project.

```
nico@nsg-System:~/projects/P4-NetFPGA$ git describe --always
v1.3.1-46-g97d3aaa
```

A.5 P4/NetFPGA compilation process

After having setup the compile host as described above, the script `bin/do-all-steps.sh` that is included in the thesis' git repository. With a NetFPGA card installed in the host, this script will compile the P4 source code to PX and in a second step to HDL and then upload the resulting bitstream to the NetFPGA. The compilation process will log its output to the directory `~/master-thesis/netpfga/log/`.

A.6 P4/NetFPGA Tests

In the following sections we describe functionality tests of our code on the NetFPGA.

A.6.1 Test 1: IPv4 Egress

In this test we test whether setting the output port based on the IPv4 address. First we get the integer values of the IPv4 addresses in python:

```
>>> int(ipaddress.IPv4Address(u"10.0.0.42"))
167772202
>>> int(ipaddress.IPv4Address(u"10.0.0.4"))
167772164
>>>
```

After that we set the table entries for the NetFPGA.

```
>> table_cam_add_entry realmain_v4_networks_0 realmain.set_egress_port 167772202 => 16 0 0 0
fields = [(u'hit', 1), (u'action_run', 3), (u'out_port', 8), (u'out_port', 8), (u'mac_addr', 48), (u'task', 16), (u'table_id', 16)]
action_name = TopPipe.realmain.set_egress_port
field_vals = [1, '16', '0', '0', '0', '0']
CAM_Init_ValidateContext() - done
WROTE 0x44020250 = 0xa00002a
WROTE 0x44020280 = 0x0000
WROTE 0x44020284 = 0x0000
WROTE 0x44020288 = 0x10000000
WROTE 0x4402028c = 0x0001
READ 0x44020244 = 0x0001
WROTE 0x44020240 = 0x0001
READ 0x44020244 = 0x0001
READ 0x44020244 = 0x0001
success
>> table_cam_add_entry realmain_v4_networks_0 realmain.set_egress_port 167772164 => 16 0 0 0
fields = [(u'hit', 1), (u'action_run', 3), (u'out_port', 8), (u'out_port', 8), (u'mac_addr', 48), (u'task', 16), (u'table_id', 16)]
action_name = TopPipe.realmain.set_egress_port
field_vals = [1, '16', '0', '0', '0', '0']
CAM_Init_ValidateContext() - done
WROTE 0x44020250 = 0xa000004
WROTE 0x44020280 = 0x0000
WROTE 0x44020284 = 0x0000
WROTE 0x44020288 = 0x10000000
WROTE 0x4402028c = 0x0001
READ 0x44020244 = 0x0001
WROTE 0x44020240 = 0x0001
READ 0x44020244 = 0x0001
READ 0x44020244 = 0x0001
success
>>
```

On the host we setup the ARP entries:

```
root@ESPRIMO-P956:~# ip neigh add 10.0.0.6 lladdr f8:f2:1e:09:62:d1 dev enp2s0f0
root@ESPRIMO-P956:~# ip neigh add 10.0.0.4 lladdr f8:f2:1e:09:62:d1 dev enp2s0f0
```

And then we generate test packets and expect 4 packets to show up on `enp2s0f0`. The following `tcpdump` output shows the expected packets arriving on `enp2s0f0`:

```
nico@ESPRIMO-P956:~$ sudo tcpdump -ni enp2s0f0
tcpdump: verbose output suppressed, use -v or -vv for full protocol decode
listening on enp2s0f0, link-type EN10MB (Ethernet), capture size 262144 bytes
10:49:28.200407 IP 10.0.0.42 > 10.0.0.4: ICMP echo request, id 4440, seq 1, length 64
10:49:28.200445 IP 10.0.0.42 > 10.0.0.4: ICMP echo request, id 4440, seq 1, length 64
10:49:29.222340 IP 10.0.0.42 > 10.0.0.4: ICMP echo request, id 4440, seq 2, length 64
10:49:29.222418 IP 10.0.0.42 > 10.0.0.4: ICMP echo request, id 4440, seq 2, length 64
```

A.6.2 Test 2: IPv6 egress

This test shows how setting the egress port based on the IPv6 address works with the NetPFGA. Similar to the previous test, we first the the Integer values of the IPv6 addresses:

```
>>> int (ipaddress.IPv6Address (u"2001:db8:42::4" )
42540766411362381960998550477184434180L
>>> int (ipaddress.IPv6Address (u"2001:db8:42::6" )
42540766411362381960998550477184434182L
>>> int (ipaddress.IPv6Address (u"2001:db8:42::42" )
42540766411362381960998550477184434242L
```

After that we set the table entries:

```
>> table_cam_add_entry realmain_v6_networks_0 realmain.set_egress_port 42540766411362381960998550477184434182 => 64 0 0 0 0
fields = [(u'hit', 1), (u'action_run', 3), (u'out_port', 8), (u'out_port', 8), (u'mac_addr', 48), (u'task', 16), (u'table_id', 16)]
action_name = TopPipe.realmain.set_egress_port
field_vals = [1, '64', '0', '0', '0', '0']
CAM_Init_ValidateContext() - done
WROTE 0x44020350 = 0x0006
WROTE 0x44020354 = 0x0000
WROTE 0x44020358 = 0x420000
WROTE 0x4402035c = 0x20010db8
WROTE 0x44020380 = 0x0000
WROTE 0x44020384 = 0x0000
WROTE 0x44020388 = 0x40000000
WROTE 0x4402038c = 0x0001
READ 0x44020344 = 0x0001
WROTE 0x44020340 = 0x0001
READ 0x44020344 = 0x0001
READ 0x44020344 = 0x0001
success
>> table_cam_add_entry realmain_v6_networks_0 realmain.set_egress_port 42540766411362381960998550477184434242 => 64 0 0 0 0
fields = [(u'hit', 1), (u'action_run', 3), (u'out_port', 8), (u'out_port', 8), (u'mac_addr', 48), (u'task', 16), (u'table_id', 16)]
action_name = TopPipe.realmain.set_egress_port
field_vals = [1, '64', '0', '0', '0', '0']
CAM_Init_ValidateContext() - done
WROTE 0x44020350 = 0x0042
WROTE 0x44020354 = 0x0000
WROTE 0x44020358 = 0x420000
WROTE 0x4402035c = 0x20010db8
WROTE 0x44020380 = 0x0000
WROTE 0x44020384 = 0x0000
WROTE 0x44020388 = 0x40000000
WROTE 0x4402038c = 0x0001
READ 0x44020344 = 0x0001
WROTE 0x44020340 = 0x0001
READ 0x44020344 = 0x0001
READ 0x44020344 = 0x0001
success
>>
```

On the host we set the IPv6 neighbor entries:

```
nico@ESPRIMO-P956:~$ sudo ip -6 neigh add 2001:db8:42::6 lladdr f8:f2:1e:09:62:d0 dev enp2s0f1
nico@ESPRIMO-P956:~$ sudo ip -6 neigh add 2001:db8:42::4 lladdr f8:f2:1e:09:62:d0 dev enp2s0f1
```

And generate the test packets:

```
nico@ESPRIMO-P956:~$ ping6 -c2 2001:db8:42::6
PING 2001:db8:42::6 (2001:db8:42::6) 56 data bytes
```

```
nico@ESPRIMO-P956:~$ sudo tcpdump -ni enp2s0f1
tcpdump: verbose output suppressed, use -v or -vv for full protocol decode
listening on enp2s0f1, link-type EN10MB (Ethernet), capture size 262144 bytes
11:30:17.287577 IP6 2001:db8:42::42 > 2001:db8:42::6: ICMP6, echo request, seq 1, length 64
11:30:17.287599 IP6 2001:db8:42::42 > 2001:db8:42::6: ICMP6, echo request, seq 1, length 64
11:30:18.310178 IP6 2001:db8:42::42 > 2001:db8:42::6: ICMP6, echo request, seq 2, length 64
11:30:18.310258 IP6 2001:db8:42::42 > 2001:db8:42::6: ICMP6, echo request, seq 2, length 64
```

The packets are successfully seen by tcpdump.

A.7 P4/BMV2 environment and tests

All BMV2 based compilations were made with the following compiler:

```
p4@ubuntu:~$ p4c --version
p4c 0.5 (SHA: 5ae30ee)
```

The installation is based on the vagrant files that were provided in the “Advanced Topics in Communication Networks Fall 2018” course of ETHZ (<https://adv-net.ethz.ch/2018/>) and contains p4tools as well as all utilities that came with the vagrant installation. For running the diff based checksum code, the following steps are necessary: First compile the p4 code and then start the switch, both with p4run.

```
cd ~/master-thesis/p4app
sudo p4run --config nat64-diff.json
```

Then with starting the controller the required table entries will

```
cd ~/master-thesis/p4app  
sudo python ./controller.py --mode range_router
```

Appendix B

NetFPGA Logs

The log files of the NetFPGA compilations are stored inside the source code directory stored at `netpfga/logs`. It follows a selection of excerpts of log files that might be relevant for reproducing the work.

B.1 NetFPGA Flash Errors

Sometimes flashing bitfiles to the NetFPGA will fail. A random amount of reboots (1 to 3) and a random amount of reflashing will fix this problem. Below can be found the log output from the flashing process.

```
nico@nsg-System:~/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/bitfiles$
sudo bash -c ". $HOME/master-thesis/netpfga/bashinit && $(pwd -P)/program_switch.sh"
++ which vivado
+ xilinx_tool_path=/opt/Xilinx/Vivado/2018.2/bin/vivado
+ bitimage=minip4.bit
+ configWrites=config_writes.sh
+ '[' -z minip4.bit ']'
+ '[' -z config_writes.sh ']'
+ '[' /opt/Xilinx/Vivado/2018.2/bin/vivado == '' ']'
+ rmdir sume_riffa
+ xsct /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/tools/run_xsct.tcl -tclargs minip4.bit
rlwrap: warning: your $TERM is 'screen' but rlwrap couldn't find it in the terminfo database. Expect some problems.
RUN loading image file.
minip4.bit
100% 19MB 1.7MB/s 00:11
fpga configuration failed. DONE PIN is not HIGH
  invoked from within
  "::tcf::eval -progress ::xsdb::print_progress (: :tcf::cache_enter tcfchan#0 (tcf_cache_eval {process_tcf_actions_cache_client ::tcfclient#0::arg}))"
    (procedure "::tcf::cache_eval_with_progress" line 2)
  invoked from within
  "::tcf::cache_eval_with_progress [dict get $arg chan] [list process_tcf_actions_cache_client $argvar] $progress"
    (procedure "process_tcf_actions" line 1)
  invoked from within
  "process_tcf_actions $arg ::xsdb::print_progress"
    (procedure "fpga" line 430)
  invoked from within
  "fpga -f $bitimage"
    (file "/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/tools/run_xsct.tcl" line 33)

+ bash /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/tools/pci_rescan_run.sh
Check programming FPGA or Reboot machine !
+ rmdir sume_riffa
rmdir: ERROR: Module sume_riffa is not currently loaded
+ modprobe sume_riffa
+ ifconfig nf0 up
nf0: ERROR while getting interface flags: No such device
+ ifconfig nf1 up
nf1: ERROR while getting interface flags: No such device
+ ifconfig nf2 up
nf2: ERROR while getting interface flags: No such device
+ ifconfig nf3 up
nf3: ERROR while getting interface flags: No such device
+ bash config_writes.sh
```

B.2 NetFPGA Flash Success

A successful flashing process also emits a couple of errors, however the message “fpga configuration failed. DONE PIN is not HIGH” and its succeeding lines are missing, as seen below. After that in all cases a reboot is required; the PCI rescan in none of our test cases re enabled the nf devices.

```
nico@nsg-System:~$ cd $NF_DESIGN_DIR/bitfiles/
nico@nsg-System:~/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/bitfiles$
sudo bash -c ". $HOME/master-thesis/netpfga/bashinit && $(pwd -P)/program_switch.sh"
++ which vivado
+ xilinx_tool_path=/opt/Xilinx/Vivado/2018.2/bin/vivado
+ bitimage=minip4.bit
```

```

+ configWrites=config_writes.sh
+ '[' -z minip4.bit ']'
+ '[' -z config_writes.sh ']'
+ '[' /opt/Xilinx/Vivado/2018.2/bin/vivado == '' ']'
+ rmmmod sume_riffa
+ xsct /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/tools/run_xsct.tcl -tclargs minip4.bit
rlwrap: warning: your $TERM is 'xterm-256color' but rlwrap couldn't find it in the terminfo database. Expect some problems.
RUN loading image file.
minip4.bit
attempting to launch hw_server

***** Xilinx hw_server v2018.2
**** Build date : Jun 14 2018-20:18:37
** Copyright 1986-2018 Xilinx, Inc. All Rights Reserved.

INFO: hw_server application started
INFO: Use Ctrl-C to exit hw_server application

INFO: To connect to this hw_server instance use url: TCP:127.0.0.1:3121

100% 19MB 1.7MB/s 00:11
+ bash /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/tools/pci_rescan_run.sh
Check programming FPGA or Reboot machine !
+ rmmmod sume_riffa
rmmmod: ERROR: Module sume_riffa is not currently loaded
+ modprobe sume_riffa
+ ifconfig nf0 up
nf0: ERROR while getting interface flags: No such device
+ ifconfig nf1 up
nf1: ERROR while getting interface flags: No such device
+ ifconfig nf2 up
nf2: ERROR while getting interface flags: No such device
+ ifconfig nf3 up
nf3: ERROR while getting interface flags: No such device
+ bash config_writes.sh
nico@nsg-System:~/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/bitfiles$

```

B.3 NetFPGA Kernel module

After a successful flash, loading the kernel module will enable nf devices to appear in the operating system.

```

nico@nsg-System:~$ ip l
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN mode DEFAULT group default qlen 1000
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
2: eth0: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc pfifo_fast state UP mode DEFAULT group default qlen 1000
    link/ether 74:d0:2b:98:38:f6 brd ff:ff:ff:ff:ff:ff
3: eth1: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether f8:f2:1e:41:44:9c brd ff:ff:ff:ff:ff:ff
4: eth2: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether f8:f2:1e:41:44:9d brd ff:ff:ff:ff:ff:ff
5: wg0: <POINTOPOINT,NOARP,UP,LOWER_UP> mtu 1420 qdisc noqueue state UNKNOWN mode DEFAULT group default qlen 1000
    link/none
nico@nsg-System:~$ ~/master-thesis/bin/build-load-drivers.sh
+ cd /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0
+ sudo modprobe -r sume_riffa
+ make clean
make -C /lib/modules/4.15.0-55-generic/build M=/home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0 clean
make[1]: Entering directory '/usr/src/linux-headers-4.15.0-55-generic'
  CLEAN /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0/.tmp_versions
  CLEAN /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0/Module.symvers
make[1]: Leaving directory '/usr/src/linux-headers-4.15.0-55-generic'
+ make all
make -C /lib/modules/4.15.0-55-generic/build M=/home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0 modules
make[1]: Entering directory '/usr/src/linux-headers-4.15.0-55-generic'
  CC [M] /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0/sume_riffa.o
Building modules, stage 2.
MODPOST 1 modules
CC /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0/sume_riffa.mod.o
LD [M] /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0/sume_riffa.ko
make[1]: Leaving directory '/usr/src/linux-headers-4.15.0-55-generic'
+ sudo make install
make -C /lib/modules/4.15.0-55-generic/build M=/home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0 modules
make[1]: Entering directory '/usr/src/linux-headers-4.15.0-55-generic'
  Building modules, stage 2.
MODPOST 1 modules
make[1]: Leaving directory '/usr/src/linux-headers-4.15.0-55-generic'
install -o root -g root -m 0755 -d /lib/modules/4.15.0-55-generic/extra/sume_riffa/
install -o root -g root -m 0755 sume_riffa.ko /lib/modules/4.15.0-55-generic/extra/sume_riffa/
depmod -a 4.15.0-55-generic
+ sudo modprobe sume_riffa
+ grep sume_riffa
+ lsmod
sume_riffa                28672  0
nico@nsg-System:~$
nico@nsg-System:~$ ip l
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN mode DEFAULT group default qlen 1000
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
2: eth0: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc pfifo_fast state UP mode DEFAULT group default qlen 1000
    link/ether 74:d0:2b:98:38:f6 brd ff:ff:ff:ff:ff:ff
3: eth1: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether f8:f2:1e:41:44:9c brd ff:ff:ff:ff:ff:ff
4: eth2: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether f8:f2:1e:41:44:9d brd ff:ff:ff:ff:ff:ff
5: wg0: <POINTOPOINT,NOARP,UP,LOWER_UP> mtu 1420 qdisc noqueue state UNKNOWN mode DEFAULT group default qlen 1000
    link/none
6: nf0: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether 02:53:55:4d:45:00 brd ff:ff:ff:ff:ff:ff
7: nf1: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether 02:53:55:4d:45:01 brd ff:ff:ff:ff:ff:ff
8: nf2: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether 02:53:55:4d:45:02 brd ff:ff:ff:ff:ff:ff
9: nf3: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000
    link/ether 02:53:55:4d:45:03 brd ff:ff:ff:ff:ff:ff
nico@nsg-System:~$

```

B.4 NetFPGA compile logs

This section shows a compilation of of NetFPGA compile output and errors. Unfound tbl files that are not correctly generated fail the compilation:

```
# Fix introduced for SDNet 2017.4
sed -i 's/$(XSIM).dir/$(XSC)/$(DPI).so/$(DPI).so/g' nf_sume_sdnet_ip/SimpleSumeSwitch/vivado_sim.bash
sed -i 's/$(XSIM).dir/$(XSC)/$(DPI).so/$(DPI).so/g' nf_sume_sdnet_ip/SimpleSumeSwitch/vivado_sim_waveform.bash
# Fix introduced for SDNet 2018.2
sed -i 's/glbl_sim/glbl/g' nf_sume_sdnet_ip/SimpleSumeSwitch/vivado_sim_waveform.bash
cp src/*.tbl nf_sume_sdnet_ip/SimpleSumeSwitch/
cp: cannot stat 'src/*.tbl': No such file or directory
make: *** [Makefile:23: cpp_test] Error 1
[23:12] loch:minip4%
```

Failure to generate an intermediate file:

```
ERROR: [XSIM 43-3409] Failed to compile generated C file xsim.dir/work.SimpleSumeSwitch_tb#work.glbl/obj/xsim_3.c.
ERROR: [XSIM 43-3915] Encountered a fatal error. Cannot continue. Exiting...
```

Failure to compile because libncurses.so.5 is missing:

```
/opt/Xilinx/Vivado/2018.2/data/./tps/llvm/3.1/linux64.o/bin/clang -fPIC -c -std=gnu89 -nobuiltinc -nostdinc++ -w -Wl,--unres
olved-symbols=ignore-in-object-files -fbracket-depth=1048576 -I/opt/Xilinx/Vivado/2018.2/data/./tps/llvm/3.1/linux64.o/bin/./li
b/clang/3.1/include -fPIC -m64 -I"/opt/Xilinx/Vivado/2018.2/data/xsim/include" "xsim.dir/work.SimpleSumeSwitch_tb#work.glbl/ob
j/xsim_3.c" -O0 -sim -o "xsim.dir/work.SimpleSumeSwitch_tb#work.glbl/obj/xsim_3.lnx64.o" -DXILINX_SIMULATOR
/opt/Xilinx/Vivado/2018.2/data/./tps/llvm/3.1/linux64.o/bin/clang: error while loading shared libraries: libncurses.so.5: cannot
open shared object file: No such file or directory
ERROR: [XSIM 43-3409] Failed to compile generated C file xsim.dir/work.SimpleSumeSwitch_tb#work.glbl/obj/xsim_3.c.
ERROR: [XSIM 43-3915] Encountered a fatal error. Cannot continue. Exiting...
[20:00] rainbow:SimpleSumeSwitch%
```

Failure to access txt files that were not correctly generated in a different compilation step:

```
# Fix introduced for SDNet 2018.2
sed -i 's/glbl_sim/glbl/g' nf_sume_sdnet_ip/SimpleSumeSwitch/vivado_sim_waveform.bash
sed -i 's/SimpleSumeSwitch_tb_sim#work.glbl/SimpleSumeSwitch_tb/g' nf_sume_sdnet_ip/SimpleSumeSwitch/vivado_sim_waveform.bash
cp src/*.tbl nf_sume_sdnet_ip/SimpleSumeSwitch/
cp testdata/*.txt nf_sume_sdnet_ip/SimpleSumeSwitch/
cp: cannot stat 'testdata/*.txt': No such file or directory
make: *** [Makefile:17: all] Error 1
[15:46] rainbow:minip4%
```

Missing pcap files of non generated testdata causing compile abortion:

```
make -C testdata/
make[1]: Entering directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/testdata'
./gen_testdata.py
/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/bin/pcap2axi --output Packet_in.axi --bus_width 256 src.pcap
Traceback (most recent call last):
  File "/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/bin/pcap2axi", line 108, in <module>
    write_to_file(args.file_pcap, args.output)
  File "/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/bin/pcap2axi", line 88, in write_to_file
    for pkt in rdpcap(file_in):
  File "/usr/lib/python2.7/dist-packages/scapy/utils.py", line 728, in rdpcap
    with PcapReader(filename) as fdesc:
  File "/usr/lib/python2.7/dist-packages/scapy/utils.py", line 751, in __call__
    filename, fdesc, magic = cls.open(filename)
  File "/usr/lib/python2.7/dist-packages/scapy/utils.py", line 778, in open
    fdesc = open(filename, "rb")
IOError: [Errno 2] No such file or directory: 'src.pcap'
make[1]: *** [Makefile:5: all] Error 1
make[1]: Leaving directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/testdata'
make: *** [Makefile:32: frontend] Error 2
[15:47] rainbow:minip4%
```

Syntax errors due to incorrect generation of a python script:

```
update_compile_order: Time (s): cpu = 00:00:17 ; elapsed = 00:00:09 . Memory (MB): peak = 1995.594 ; gain = 0.016 ; free physic
al = 21975 ; free virtual = 33161
loading libsume..
Traceback (most recent call last):
  File "/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/test/sim_switch_de
fault/run.py", line 42, in <module>
    import config_writes
  File "/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/test/sim_switch_de
fault/config_writes.py", line 7
    ^
IndentationError: expected an indented block
while executing
"exec python $::env(NF_DESIGN_DIR)/test/${test_name}/run.py"
  invoked from within
"set output [exec python $::env(NF_DESIGN_DIR)/test/${test_name}/run.py]"
(file "/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/hw/tcl/simple_s
ume_switch_sim.tcl" line 177)
INFO: [Common 17-206] Exiting Vivado at Sat May 18 15:21:21 2019...
```

Missing axi files don't abort the compilation process: (shortened for formatting)

```
cp: cannot stat 'simple_sume_switch/test/nf_interface_0_log.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_0_stim.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_0_expected.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_1_log.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_1_stim.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_1_expected.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_2_log.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_2_stim.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_2_expected.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_3_log.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_3_stim.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/nf_interface_3_expected.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/dma_0_log.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/dma_0_expected.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/reg_stim.log': No such file or directory
cp: cannot stat 'simple_sume_switch/test/reg_expect.axi': No such file or directory
cp: cannot stat 'simple_sume_switch/test/reg_stim.axi': No such file or directory
=== Running test /tmp/nico/test/simple_sume_switch/sim_switch_default
... using cmd
['/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/test/sim_switch_default/run.py',
 '--sim', 'xsim']
```



```
[7:05] rainbow:netpfga% bash build-load-drivers.sh
+ cd /home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0
+ make all
make -C /lib/modules/5.0.0-16-generic/build M=/home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0 modules
make[1]: Entering directory '/usr/src/linux-headers-5.0.0-16-generic'
Building modules, stage 2.
MODPOST 1 modules
make[1]: Leaving directory '/usr/src/linux-headers-5.0.0-16-generic'
+ sudo make install
make -C /lib/modules/5.0.0-16-generic/build M=/home/nico/projects/P4-NetFPGA/lib/sw/std/driver/sume_riffa_v1_0_0 modules
make[1]: Entering directory '/usr/src/linux-headers-5.0.0-16-generic'
Building modules, stage 2.
MODPOST 1 modules
make[1]: Leaving directory '/usr/src/linux-headers-5.0.0-16-generic'
install -o root -g root -m 0755 -d /lib/modules/5.0.0-16-generic/extra/sume_riffa/
install -o root -g root -m 0755 sume_riffa.ko /lib/modules/5.0.0-16-generic/extra/sume_riffa/
depmod -a 5.0.0-16-generic
+ sudo modprobe sume_riffa
modprobe: ERROR: could not insert 'sume_riffa': Exec format error
[7:06] rainbow:netpfga%
```

Java traceback when trying to install SDNET: (reason was a hidden window)

```
Exception in thread "AWT-EventQueue-0" java.lang.IllegalArgumentException: Window must not be zero
at java.desktop/sun.awt.X11.XAtom.checkWindow(Unknown Source)
at java.desktop/sun.awt.X11.XAtom.getAtomData(Unknown Source)
at java.desktop/sun.awt.X11.XToolkit.getWorkArea(Unknown Source)
at java.desktop/sun.awt.X11.XToolkit.getInsets(Unknown Source)
at java.desktop/sun.awt.X11.XToolkit.getScreenInsets(Unknown Source)
at java.desktop/java.awt.Window.init(Unknown Source)
at java.desktop/java.awt.Window.<init>(Unknown Source)
at java.desktop/java.awt.Window.<init>(Unknown Source)
at java.desktop/java.awt.Dialog.<init>(Unknown Source)
at java.desktop/java.awt.Dialog.<init>(Unknown Source)
at java.desktop/java.awt.Dialog.<init>(Unknown Source)
at java.desktop/javaw.swing.JDialog.<init>(Unknown Source)
at java.desktop/javaw.swing.JOptionPane.createDialog(Unknown Source)
at java.desktop/javaw.swing.JOptionPane.createDialog(Unknown Source)
at j.a.c(Unknown Source)
at j.a.a(Unknown Source)
at j.a.a(Unknown Source)
at j.a.c(Unknown Source)
at com.xilinx.installer.gui.panel.destination.b.a(Unknown Source)
at com.xilinx.installer.gui.panel.destination.DestinationPanel.z(Unknown Source)
at com.xilinx.installer.gui.E.a(Unknown Source)
at com.xilinx.installer.gui.InstallerGUI.l(Unknown Source)
at com.xilinx.installer.gui.i.actionPerformed(Unknown Source)
at java.desktop/javaw.swing.AbstractButton.fireActionPerformed(Unknown Source)
at java.desktop/javaw.swing.AbstractButton$Handler.actionPerformed(Unknown Source)
at java.desktop/javaw.swing.DefaultButtonModel.fireActionPerformed(Unknown Source)
at java.desktop/javaw.swing.DefaultButtonModel.setPressed(Unknown Source)
at java.desktop/javaw.swing.plaf.basic.BasicButtonListener.mouseReleased(Unknown Source)
at java.desktop/java.awt.Component.processMouseEvent(Unknown Source)
at java.desktop/javaw.swing.JComponent.processMouseEvent(Unknown Source)
at java.desktop/java.awt.Component.processEvent(Unknown Source)
at java.desktop/java.awt.Container.processEvent(Unknown Source)
at java.desktop/java.awt.Component.dispatchEventImpl(Unknown Source)
at java.desktop/java.awt.Container.dispatchEventImpl(Unknown Source)
at java.desktop/java.awt.Component.dispatchEvent(Unknown Source)
at java.desktop/java.awt.LightweightDispatcher.retargetMouseEvent(Unknown Source)
at java.desktop/java.awt.LightweightDispatcher.processMouseEvent(Unknown Source)
at java.desktop/java.awt.LightweightDispatcher.dispatchEvent(Unknown Source)
at java.desktop/java.awt.Container.dispatchEventImpl(Unknown Source)
at java.desktop/java.awt.Window.dispatchEventImpl(Unknown Source)
at java.desktop/java.awt.Component.dispatchEvent(Unknown Source)
at java.desktop/java.awt.EventQueue.dispatchEventImpl(Unknown Source)
at java.desktop/java.awt.EventQueue.access$500(Unknown Source)
at java.desktop/java.awt.EventQueue$3.run(Unknown Source)
at java.desktop/java.awt.EventQueue$3.run(Unknown Source)
at java.base/java.security.AccessController.doPrivileged(Native Method)
at java.base/java.security.ProtectionDomain$JavaSecurityAccessImpl.doIntersectionPrivilege(Unknown Source)
at java.base/java.security.ProtectionDomain$JavaSecurityAccessImpl.doIntersectionPrivilege(Unknown Source)
at java.desktop/java.awt.EventQueue$4.run(Unknown Source)
at java.desktop/java.awt.EventQueue$4.run(Unknown Source)
at java.base/java.security.AccessController.doPrivileged(Native Method)
at java.base/java.security.ProtectionDomain$JavaSecurityAccessImpl.doIntersectionPrivilege(Unknown Source)
at java.desktop/java.awt.EventQueue.dispatchEvent(Unknown Source)
at java.desktop/java.awt.EventDispatchThread.pumpOneEventForFilters(Unknown Source)
at java.desktop/java.awt.EventDispatchThread.pumpEventsForFilter(Unknown Source)
at java.desktop/java.awt.EventDispatchThread.pumpEventsForHierarchy(Unknown Source)
at java.desktop/java.awt.EventDispatchThread.pumpEvents(Unknown Source)
at java.desktop/java.awt.EventDispatchThread.pumpEvents(Unknown Source)
at java.desktop/java.awt.EventDispatchThread.run(Unknown Source)
```

Failures when testing the first NetFPGA card

```
-----
[ddr3B]: Running Auto Test
-----
Traceback (most recent call last):
File "/usr/lib/python2.7/dist-packages/wx-3.0-gtk2/wx/_core.py", line 16765, in <lambda>
lambda event: event.callable(*event.args, **event.kw)
File "sw/host/script/NFSumeTest.py", line 848, in UpdateProgress
self.progressDlg.Update(self.curProgress, str(localLine))
File "/usr/lib/python2.7/dist-packages/wx-3.0-gtk2/wx/_core.py", line 16710, in __getattr__
raise PyDeadObjectError(self.attrStr % self._name)
wx._core.PyDeadObjectError: The C++ part of the NFSumeProgress object has been deleted, attribute access no longer allowed.
Exception in thread Thread-18:
Traceback (most recent call last):
File "/usr/lib/python2.7/threading.py", line 801, in __bootstrap_inner
self.run()
File "sw/host/script/NFSumeTest.py", line 947, in run
self.target(*self.data)
File "sw/host/script/NFSumeTest.py", line 355, in StartAutoTest
self.TestInterface(testName)
File "sw/host/script/NFSumeTest.py", line 465, in TestInterface
self.ProgramFpga('.../bitfiles/' + self.nfsumeTestConfiguration[testName]['bitstream'])
File "sw/host/script/NFSumeTest.py", line 586, in ProgramFpga
self.getFpgaIndex()
File "sw/host/script/NFSumeTest.py", line 574, in getFpgaIndex
p = Popen(['djtgcfg', 'init', '-d', 'NetSUME'], stdout=PIPE, bufsize = 1)
File "/usr/lib/python2.7/subprocess.py", line 711, in __init__
errread, errwrite)
File "/usr/lib/python2.7/subprocess.py", line 1343, in _execute_child
raise child_exception
OSError: [Errno 2] No such file or directory
```

More failures when testing the first NetFPGA card

```
-----
[pcie]: Running Auto Test
-----
Traceback (most recent call last):
  File "/usr/lib/python2.7/dist-packages/wx-3.0-gtk2/wx/_core.py", line 16765, in <lambda>
    lambda event: event.callable(*event.args, **event.kw) )
  File "sw/host/script/NfSumeTest.py", line 848, in UpdateProgress
    self.progressDlg.Update(self.curProgress, str(localLine))
  File "/usr/lib/python2.7/dist-packages/wx-3.0-gtk2/wx/_core.py", line 16710, in __getattr__
    raise PyDeadObjectError(self.attrStr % self._name)
wx._core.PyDeadObjectError: The C++ part of the NfSumeProgress object has been deleted, attribute access no longer allowed.
Exception in thread Thread-21:
Traceback (most recent call last):
  File "/usr/lib/python2.7/threading.py", line 801, in __bootstrap_inner
    self.run()
  File "sw/host/script/NfSumeTest.py", line 947, in run
    self.target(*self.data)
  File "sw/host/script/NfSumeTest.py", line 466, in TestInterface
    self.serialCon.readlines()
  File "/usr/lib/python2.7/dist-packages/serial/serialposix.py", line 495, in read
    raise SerialException('device reports readiness to read but returned no data (device disconnected or multiple access on port?)')
SerialException: device reports readiness to read but returned no data (device disconnected or multiple access on port?)
```

Unexpected EOF during compilation:

```
ERROR: [VRFC 10-1491] unexpected EOF
[/home/nico/master-thesis/netpfga/minip4/nf_sume_sdnet_ip/
SimpleSumeSwitch/S_CONTROLLERs.HDL/S_CONTROLLER_SimpleSumeSwitch.vp:37]
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_0_ErrorCheck
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_1_ErrorCheck
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_2_ErrorCheck
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_3_ErrorCheck
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_4_ErrorCheck
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_5_ErrorCheck
INFO: [VRFC 10-311] analyzing module TopDeparser_t_EngineStage_6_ErrorCheck
```

The function syntax is not supported by p4/netfpga:

```
make[1]: Entering directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/src'
p4c-sdnet -o minip4.sdnet --sdnet_info .sdnet_switch_info.dat minip4_solution.p4
headers.p4(246):syntax error, unexpected IDENTIFIER, expecting (
bit<16> ones_complement_sum
^^^^^^^^^^^^^^^^^^^^^^^^
error: 1 errors encountered, aborting compilation
Makefile:34: recipe for target 'all' failed
make[1]: *** [all] Error 1
make[1]: Leaving directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/src'
Makefile:31: recipe for target 'frontend' failed
make: *** [frontend] Error 2
nico@nsg-System:~/master-thesis/netpfga$
```

The config_writes.py is missing due to a previous, non critical compilation error:

```
nico@nsg-System:~/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/test/sim_switch_default$
cd $NF_DESIGN_DIR/test/sim_switch_default && make 2>&1 | tee ~/master-thesis/netpfga/log/step8-8(date +%F-%H%M%S)
rm -f config_writes.py*
rm -f *.pyc
cp /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/testdata/config_writes.py ./
cp: cannot stat '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/testdata/config_writes.py': No such file or directory
Makefile:36: recipe for target 'all' failed
make: *** [all] Error 1
```

Failed to synthesizing module errors:

```
WARNING: [Synth 8-689] width (12) of port connection 'control_S_AXI_ARADDR' does not match port width (8) of module 'SimpleSumeSwitch'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/hw/project/
simple_sume_switch.srscs/sources_1/ip/nf_sume_sdnet_ip/nf_sume_sdnet_ip_wrapper/nf_sume_sdnet.v:199]
ERROR: [Synth 8-448] named port connection 'tuple_out_sume_metadata_VALID' does not exist for instance 'SimpleSumeSwitch_inst' of module 'SimpleSumeSwitch'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/hw/project/
simple_sume_switch.srscs/sources_1/ip/nf_sume_sdnet_ip/nf_sume_sdnet_ip_wrapper/nf_sume_sdnet.v:218]
ERROR: [Synth 8-448] named port connection 'tuple_out_sume_metadata_DATA' does not exist for instance 'SimpleSumeSwitch_inst' of module 'SimpleSumeSwitch'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/hw/project/
simple_sume_switch.srscs/sources_1/ip/nf_sume_sdnet_ip/nf_sume_sdnet_ip_wrapper/nf_sume_sdnet.v:219]
ERROR: [Synth 8-6156] failed synthesizing module 'nf_sume_sdnet'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/hw/project/
simple_sume_switch.srscs/sources_1/ip/nf_sume_sdnet_ip/nf_sume_sdnet_ip_wrapper/nf_sume_sdnet.v:44]
ERROR: [Synth 8-6156] failed synthesizing module 'nf_sume_sdnet_ip'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch/hw/project/
simple_sume_switch.srscs/sources_1/ip/nf_sume_sdnet_ip/synth/nf_sume_sdnet_ip.v:57]
ERROR: [Synth 8-6156] failed synthesizing module 'nf_datapath'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/
simple_sume_switch/hw/hdl/nf_datapath.v:44]
ERROR: [Synth 8-6156] failed synthesizing module 'top'
[/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/
simple_sume_switch/hw/hdl/top.v:43]
```

Missing "source" files abort CLI compilation errors:

```
cc -c -fPIC /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/sw/API/CAM.c
-I/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/sw/API
cc -std=c99 -Wall -Werror -fPIC -c libcam.c
-I/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/sw/sume
-I/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/sw/API
cc -L/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/sw/sume
-shared -o libcam.so libcam.o CAM.o -lsumereg
/usr/bin/ld: cannot find -lsumereg
collect2: error: ld returned 1 exit status
Makefile:52: recipe for target 'libcam' failed
make[1]: *** [libcam] Error 1
make[1]: Leaving directory '/home/nico/master-thesis/netpfga/minip4/sw/CLI'
ERROR: could not compile libcam source files
```

Generated axi files not found at a different stage:

```

cp: cannot stat '/home/nico/projects/P4-NetFPGA/contrib-projects/
sume-sdnet-switch/projects/minip4/simple_sume_switch/test/dma_0_expected.axi': No such file or directory
cp: cannot stat '/home/nico/projects/P4-NetFPGA/contrib-projects/
sume-sdnet-switch/projects/minip4/simple_sume_switch/test/Makefile': No such file or directory
cp: cannot stat '/home/nico/projects/P4-NetFPGA/contrib-projects/
sume-sdnet-switch/projects/minip4/simple_sume_switch/test/reg_stim.log': No such file or directory
cp: cannot stat '/home/nico/projects/P4-NetFPGA/contrib-projects/
sume-sdnet-switch/projects/minip4/simple_sume_switch/test/reg_expect.axi': No such file or directory
cp: cannot stat '/home/nico/projects/P4-NetFPGA/contrib-projects/
sume-sdnet-switch/projects/minip4/simple_sume_switch/test/reg_stim.axi': No such file or directory
NetFPGA environment:
  Root dir:      /home/nico/projects/P4-NetFPGA
  Project name:  simple_sume_switch
  Project dir:   /tmp/nico/test/simple_sume_switch
  Work dir:     /tmp/nico
512
=== Work directory is /tmp/nico/test/simple_sume_switch
=== Setting up test in /tmp/nico/test/simple_sume_switch/sim_switch_default
=== Running test /tmp/nico/test/simple_sume_switch/sim_switch_default ... using
cmd ['/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/
minip4/simple_sume_switch/test/sim_switch_default/run.py', '--sim', 'xsim']+ date
Die Jul 23 13:34:54 CEST 2019
+ [ = no ]
+ cd /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/simple_sume_switch
+ make
make: *** No targets specified and no makefile found.  Stop.

```

Renaming variables as follows breaks the compile process

```

@Xilinx_MaxPacketRegion(1024)
control TopDeparser (
-   packet_out b,
-   in Parsed_packet p,
+   packet_out packet,
+   in Parsed_packet hdr,
    in user_metadata_t user_metadata,
    inout digest_data_t digest_data,
    inout sume_metadata_t sume_metadata) {

    apply {
-       b.emit(p.ethernet);
+       packet.emit(hdr.ethernet);
    }

+
+
}

```

In NetPFGA the LPM table size must be != 64:

```

minip4_solution.p4(38): [--Wwarn=uninitialized_out_param] warning: out parameter meta may be uninitialized when RealParser terminates
    out metadata meta,
        ^^^^^
minip4_solution.p4(35)
parser RealParser(
    ^^^^^^^^^^^
error: LPM table size should be 2^n - 1
actions_nat64_generic.p4(169): error: could not not map table size size
    size = 64;
    ^^^^^
error: table match_types are not the same
actions_arp.p4(35): error: could not map table key(s) KeyElement
    hdr.arp.dst_ipv4_addr: lpm;
    ^^^^^^^^^^^^^^^^^^^^^
error: LPM table size should be 2^n - 1
actions_arp.p4(55): error: could not not map table size size
    size = 64;
    ^^^^^
Makefile:34: recipe for target 'all' failed
make[1]: *** [all] Error 1
make[1]: Leaving directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/src'
Makefile:31: recipe for target 'frontend' failed
make: *** [frontend] Error 2
nico@nsg-System:~/master-thesis/netpfga/log$

```

Cannot mix the key table types with P4/NetFPGA:

```

make[1]: Entering directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/src'
p4c-sdnet -o minip4.sdnet --sdnet_info .sdnet_switch_info.dat minip4_solution.p4
actions_egress.p4(52): warning: Table v6_networks is not used; removing
table v6_networks {
    ^^^^^^^^^^^
actions_egress.p4(69): warning: Table v4_networks is not used; removing
table v4_networks {
    ^^^^^^^^^^^
actions_nat64_generic.p4(174): warning: Table nat46 is not used; removing
table nat46 {
    ^^^^^
minip4_solution.p4(38): [--Wwarn=uninitialized_out_param] warning: out parameter meta may be uninitialized when RealParser terminates
    out metadata meta,
        ^^^^^
minip4_solution.p4(35)
parser RealParser(
    ^^^^^^^^^^^

```

```

error: table match_types are not the same
actions_arp.p4(35): error: could not map table key(s) KeyElement
      hdr.arp.dst_ipv4_addr: lpm;
      ^^^^^^^^^^^^^^^^^^^^^^^
Makefile:34: recipe for target 'all' failed
make[1]: *** [all] Error 1

```

```

table v4_arp {
  key = {
    hdr.ethernet.dst_addr: exact;
    hdr.arp.opcode: exact;
    hdr.arp.dst_ipv4_addr: lpm;
  }
  actions = {
    controller_debug_table_id;
    arp_reply;
    NoAction;
  }
  size = ICMP6_TABLE_SIZE;
  default_action = controller_debug_table_id(TABLE_ARP);
}

```

Implicit error saying that LPM tables don't work in P4/NetFPGA:

```

s/sume-sdnet-switch/projects/minip4/nf_sume_sdnet_ip/SimpleSumeSwitch/realmain_lookup_table_0_t.HDL/xpm_memory.sv
[SW] LPM_Init() - start
[SW] LPM_Init() - done
[SW] LPM_LoadDataset() - start
[SW] LPM_LoadDataset() failed with error code = 12
FATAL_ERROR: Vivado Simulator kernel has encountered an exception from DPI C function: LPM_VerifyDataset(). Please correct.
Time: 2016466 ps Iteration: 0 Process: /SimpleSumeSwitch_tb/LPM_VerifyDataset
File: /home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/nf_sume_sdnet_ip/SimpleSumeSwitch/Testbench/SimpleSumeSwitch_tb.sv

```

The table for exact matches must be at least 64 in P4/NetFPGA:

```

minip4_solution.p4(35)
parser RealParser(
  ^^^^^^^^^^^
actions_nat64_generic.p4(173): error: table size too small for match_type(EM): 63 < 64
      size = 63;
      ^^
actions_nat64_generic.p4(173): error: could not not map table size size
      size = 63;
      ^^^^

```

Unsupported default parameters in P4/NetFPGA:

```

actions_egress.p4(89): error: data-plane arguments in default_actions are currently unsupported: realmain_controller_debug_table_id_0
      default_action = controller_debug_table_id(TABLE_V4_NETWORKS);
      ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
terminate called after throwing an instance of 'Util::CompilerBug'
  what():  In file: /wrk/hdscratch/staff/mohan/p4c_sdnet/build/p4c/extensions/sdnet/translate/core/lookupEngine.cpp:137
Compiler Bug: actions_egress.p4(89): unhandled expression realmain_controller_debug_table_id/realmain_controller_debug_table_id_0(5);
      default_action = controller_debug_table_id(TABLE_V4_NETWORKS);
      ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^

```

Causing compiler bug by using an if statement at a wrong place in P4/NetFPGA:

```

minip4_solution.p4(39)
parser RealParser(
  ^^^^^^^^^^^
terminate called after throwing an instance of 'Util::CompilerBug'
  what():  In file: /wrk/hdscratch/staff/mohan/p4c_sdnet/build/p4c/extensions/sdnet/writers/pxWriter.h:20
Compiler Bug: unhandled node: <IfStatement>(471564)

Makefile:34: recipe for target 'all' failed
make[1]: *** [all] Error 134
make[1]: Leaving directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/src'
Makefile:31: recipe for target 'frontend' failed

```

Applying table "twice" in different branches is impossible in P4/NetFPGA causes a different compiler bug:

```

make -C src/
make[1]: Entering directory '/home/nico/projects/P4-NetFPGA/contrib-projects/sume-sdnet-switch/projects/minip4/src'
p4c-sdnet -o minip4.sdnet --sdnet_info .sdnet_switch_info.dat minip4_solution.p4
minip4_solution.p4(19): [--Wwarn=uninitialized_out_param] warning: out parameter meta may be uninitialized when RealParser terminates
      out metadata meta,
      ^^^^
minip4_solution.p4(16)
parser RealParser(
  ^^^^^^^^^^^
terminate called after throwing an instance of 'Util::CompilerBug'
  what():  In file: /wrk/hdscratch/staff/mohan/p4c_sdnet/build/p4c/extensions/sdnet/translate/core/tupleEngine.cpp:324
Compiler Bug: overwrite

Makefile:34: recipe for target 'all' failed

```

Adding table entries requires setting parameters for all possible actions that are registered in a table:

```

>> table_cam_add_entry realmain_v6_networks_0 realmain.set_egress_port
42540766411362381960998550477184434178 => 1 ERROR: not enough fields provided to complete _hexify()

```

The table handling scripts do not handle conversion for long integers for P4/NetFPGA:

```

>> table_cam_delete_entry realmain_v6_networks_0 42540766411362381960998550477184434179
ERROR: failed to convert 42540766411362381960998550477184434179 of type <type 'long'> to an integer
nico@nsg-System:~/master-thesis/netpfa/minip4/sw/CLI$

```

A P4/BMV2 compiler error:


```
        ret = ret + 1;
    }
    return ret[15:0];
}'''
And p4c version:
'''p4@ubuntu:~/master-thesis/p4app$ p4c --version
p4c 0.5 (SHA: 5ae30ee)'''
```

Appendix C

Benchmark Logs

C.1 Enabling hardware offloading

The following commands enable hardware offloading even though error messages are printed:

```
root@ESPRIMO-P956:~# ethtool -K enp2s0f1 tx on
Cannot get device udp-fragmentation-offload settings: Operation not supported
Cannot get device udp-fragmentation-offload settings: Operation not supported
Actual changes:
tx-checksumming: on
    tx-checksum-ip-generic: on
    tx-checksum-sctp: on
tcp-segmentation-offload: on
    tx-tcp-segmentation: on
    tx-tcp6-segmentation: on
root@ESPRIMO-P956:~# ethtool -K enp2s0f1 rx on
Cannot get device udp-fragmentation-offload settings: Operation not supported
Cannot get device udp-fragmentation-offload settings: Operation not supported
root@ESPRIMO-P956:~#
```

This results into the following:

```
root@ESPRIMO-P956:~# ethtool -k enp2s0f1
Features for enp2s0f1:
Cannot get device udp-fragmentation-offload settings: Operation not supported
rx-checksumming: on
tx-checksumming: on
    tx-checksum-ipv4: off [fixed]
    tx-checksum-ip-generic: on
    tx-checksum-ipv6: off [fixed]
    tx-checksum-fcoe-crc: on [fixed]
    tx-checksum-sctp: on
scatter-gather: on
    tx-scatter-gather: on
    tx-scatter-gather-fraglist: off [fixed]
tcp-segmentation-offload: on
    tx-tcp-segmentation: on
    tx-tcp-ecn-segmentation: off [fixed]
    tx-tcp-mangleid-segmentation: off
    tx-tcp6-segmentation: on
udp-fragmentation-offload: off
generic-segmentation-offload: on
generic-receive-offload: on
large-receive-offload: off
rx-vlan-offload: on
```

```

tx-vlan-offload: on
ntuple-filters: off
receive-hashing: on
highdma: on [fixed]
rx-vlan-filter: on
vlan-challenged: off [fixed]
tx-lockless: off [fixed]
netns-local: off [fixed]
tx-gso-robust: off [fixed]
tx-fcoe-segmentation: on [fixed]
tx-gre-segmentation: on
tx-gre-csum-segmentation: on
tx-ipxip4-segmentation: on
tx-ipxip6-segmentation: on
tx-udp_tnl-segmentation: on
tx-udp_tnl-csum-segmentation: on
tx-gso-partial: on
tx-sctp-segmentation: off [fixed]
tx-esp-segmentation: off [fixed]
fcoe-mtu: off [fixed]
tx-nocache-copy: off
loopback: off [fixed]
rx-fcs: off [fixed]
rx-all: off
tx-vlan-stag-hw-insert: off [fixed]
rx-vlan-stag-hw-parse: off [fixed]
rx-vlan-stag-filter: off [fixed]
l2-fwd-offload: off
hw-tc-offload: off
esp-hw-offload: off [fixed]
esp-tx-csum-hw-offload: off [fixed]
rx-udp_tunnel-port-offload: off
root@ESPRIMO-P956:~#

```

C.2 Tayga

Tayga is installed from the regular package database:

```
ii  tayga                                0.9.2-6
```

We prepare the networking as follows:

```

[15:12] nsg-System:~# ip addr add 10.0.0.77/24 dev eth1
[15:12] nsg-System:~# ip l s eth1 up

nico@ESPRIMO-P956:~$ ~/master-thesis/bin/init_ipv4_esprimo.sh
nico@ESPRIMO-P956:~$ cat ~/master-thesis/bin/init_ipv4_esprimo.sh
#!/bin/sh

sudo ip addr add 10.0.0.42/24 dev enp2s0f0
sudo ip link set enp2s0f0 up

nico@ESPRIMO-P956:~$ sudo ip route add 10.0.1.0/24 via 10.0.0.77

```

And verify that networking works:

```

[15:12] nsg-System:~# ping 10.0.0.42
PING 10.0.0.42 (10.0.0.42) 56(84) bytes of data.
64 bytes from 10.0.0.42: icmp_seq=1 ttl=64 time=0.304 ms

```



```
64 bytes from 10.0.0.42: icmp_seq=2 ttl=64 time=0.097 ms
^C
--- 10.0.0.42 ping statistics ---
2 packets transmitted, 2 received, 0% packet loss, time 1011ms
rtt min/avg/max/mdev = 0.097/0.200/0.304/0.104 ms
[15:12] nsg-System:~#
```

We also setup the IPv6 networking:

```
nico@ESPRIMO-P956:~$ ip addr show dev enp2s0f1
13: enp2s0f1: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc mq state UP group default qlen 1000
    link/ether f8:f2:1e:09:62:d1 brd ff:ff:ff:ff:ff:ff
    inet6 2001:db8:42::42/64 scope global
        valid_lft forever preferred_lft forever
    inet6 fe80::faf2:1eff:fe09:62d1/64 scope link
        valid_lft forever preferred_lft forever
nico@ESPRIMO-P956:~$ sudo ip route add 2001:db8:23::/96 via 2001:db8:42::77

[15:12] nsg-System:~# ip addr add 2001:db8:42::77/64 dev eth2
[15:15] nsg-System:~# ip link set eth2 up
```

And verify that IPv6 networking works:

```
nico@ESPRIMO-P956:~$ ping6 -c2 2001:db8:42::77
PING 2001:db8:42::77 (2001:db8:42::77) 56 data bytes
64 bytes from 2001:db8:42::77: icmp_seq=1 ttl=64 time=0.169 ms
64 bytes from 2001:db8:42::77: icmp_seq=2 ttl=64 time=0.153 ms

--- 2001:db8:42::77 ping statistics ---
2 packets transmitted, 2 received, 0% packet loss, time 1010ms
rtt min/avg/max/mdev = 0.153/0.161/0.169/0.008 ms
nico@ESPRIMO-P956:~$
```

We enable IPv6 and IPv4 forwarding:

```
[15:16] nsg-System:~# sysctl -w net.ipv6.conf.all.forwarding=1
net.ipv6.conf.all.forwarding = 1
```

```
[15:20] nsg-System:~# sysctl -w net.ipv4.ip_forward=1
net.ipv4.ip_forward = 1
```

And we test NAT64 with tayga:

```
nico@ESPRIMO-P956:~$ ping -c2 10.0.1.42
PING 10.0.1.42 (10.0.1.42) 56(84) bytes of data.
64 bytes from 10.0.1.42: icmp_seq=1 ttl=61 time=0.356 ms
64 bytes from 10.0.1.42: icmp_seq=2 ttl=61 time=0.410 ms

--- 10.0.1.42 ping statistics ---
2 packets transmitted, 2 received, 0% packet loss, time 1019ms
rtt min/avg/max/mdev = 0.356/0.383/0.410/0.027 ms
nico@ESPRIMO-P956:~$

nico@ESPRIMO-P956:~$ sudo tcpdump -ni enp2s0f1
tcpdump: verbose output suppressed, use -v or -vv for full protocol decode
listening on enp2s0f1, link-type EN10MB (Ethernet), capture size 262144 bytes
15:21:39.851057 IP6 2001:db8:23::a00:2a > 2001:db8:42::42: ICMP6, echo request, seq 1, length 64
15:21:39.851124 IP6 2001:db8:42::42 > 2001:db8:23::a00:2a: ICMP6, echo reply, seq 1, length 64
15:21:40.870448 IP6 2001:db8:23::a00:2a > 2001:db8:42::42: ICMP6, echo request, seq 2, length 64
15:21:40.870507 IP6 2001:db8:42::42 > 2001:db8:23::a00:2a: ICMP6, echo reply, seq 2, length 64
^C
4 packets captured
4 packets received by filter
0 packets dropped by kernel
nico@ESPRIMO-P956:~$
```

And test NAT64 from IPv6 to IPv4:

```
nico@ESPRIMO-P956:~$ ping6 -c2 2001:db8:23::a00:2a
PING 2001:db8:23::a00:2a (2001:db8:23::a00:2a) 56 data bytes
64 bytes from 2001:db8:23::a00:2a: icmp_seq=1 ttl=61 time=0.240 ms
64 bytes from 2001:db8:23::a00:2a: icmp_seq=2 ttl=61 time=0.400 ms

--- 2001:db8:23::a00:2a ping statistics ---
2 packets transmitted, 2 received, 0% packet loss, time 1003ms
rtt min/avg/max/mdev = 0.240/0.320/0.400/0.080 ms
nico@ESPRIMO-P956:~$
```

C.3 Jool

We install Jool 4.0.1 from source from <https://www.jool.mx/en/download.html> as follows:

```
nico@nsg-System:~$ wget https://github.com/NICMx/Jool/releases/download/v4.0.1/jool_4.0.1.tar.gz
nico@nsg-System:~$ tar xvfz jool_4.0.1.tar.gz
nico@nsg-System:~$ cd jool-4.0.1/
nico@nsg-System:~/jool-4.0.1$ sudo apt install linux-headers-$(uname -r)
nico@nsg-System:~/jool-4.0.1$ sudo apt install libnl-genl-3-dev
nico@nsg-System:~/jool-4.0.1$ sudo apt install iptables-dev
nico@nsg-System:~/jool-4.0.1$ sudo make install
```

We enable forwarding:

```
sysctl -w net.ipv4.conf.all.forwarding=1
sysctl -w net.ipv6.conf.all.forwarding=1
```

We configure jool to map the network prefixes and setup iptables to redirect the traffic into the jool instance:

```
[16:53] nsg-System:~# modprobe jool_siit
[16:54] nsg-System:~# jool_siit instance add "example" --iptables
[16:54] nsg-System:~# jool_siit -i example eamt add 2001:db8:42::/120 10.0.1.0/24
[16:55] nsg-System:~# jool_siit -i example eamt add 2001:db8:23::/120 10.0.0.0/24
[16:57] nsg-System:~# iptables -t mangle -A PREROUTING -s 2001:db8:42::/120 -d 2001:db8:23::/120 -j JOOL_SIIT --instance example
[16:57] nsg-System:~# iptables -t mangle -A PREROUTING -s 10.0.0.0/24 -d 10.0.1.0/24 -j JOOL_SIIT --instance example
```

Afterwards we test NAT64:

```
nico@ESPRIMO-P956:~/master-thesis/iperf$ ping6 2001:db8:23::2a
PING 2001:db8:23::2a(2001:db8:23::2a) 56 data bytes
64 bytes from 2001:db8:23::2a: icmp_seq=1 ttl=63 time=0.199 ms
64 bytes from 2001:db8:23::2a: icmp_seq=2 ttl=63 time=0.282 ms
64 bytes from 2001:db8:23::2a: icmp_seq=3 ttl=63 time=0.186 ms
^C
--- 2001:db8:23::2a ping statistics ---
3 packets transmitted, 3 received, 0% packet loss, time 2040ms
rtt min/avg/max/mdev = 0.186/0.222/0.282/0.044 ms
nico@ESPRIMO-P956:~/master-thesis/iperf$ ping 10.0.1.66
PING 10.0.1.66 (10.0.1.66) 56(84) bytes of data.
64 bytes from 10.0.1.66: icmp_seq=1 ttl=63 time=0.218 ms
64 bytes from 10.0.1.66: icmp_seq=2 ttl=63 time=0.281 ms
64 bytes from 10.0.1.66: icmp_seq=3 ttl=63 time=0.280 ms
^C
--- 10.0.1.66 ping statistics ---
3 packets transmitted, 3 received, 0% packet loss, time 2051ms
rtt min/avg/max/mdev = 0.218/0.259/0.281/0.034 ms
nico@ESPRIMO-P956:~/master-thesis/iperf$
```

List of Abbreviations

ARP	Address resolution protocol
ASIC	Application-specific integrated circuit
DAC	Direct attach cable
FGPA	Field-programmable gate array
LPM	Longest prefix matching
MTU	Maximum transfer unit
NAT	Network Address Translation
NAT64	Network Address Translation from / to IPv6 to / from IPv4
NDP	Neighbor Discovery Protocol
RIR	Regional Internet Registry
RTT	Round Trip Time

Bibliography

- [1] AFRINIC. Afrinic ipv4 exhaustion. <https://afrinic.net/exhaustion>.
- [2] Akamai. Ipv6 adoption visualization. <https://www.akamai.com/us/en/resources/our-thinking/state-of-the-internet-report/state-of-the-internet-ipv6-adoption-visualization.jsp#countries>.
- [3] T. Anderson and A. L. Popper. Explicit Address Mappings for Stateless IP/ICMP Translation. RFC 7757 (Proposed Standard), Feb. 2016.
- [4] APNIC. Apnic's ipv4 pool status. <https://www.apnic.net/community/ipv4-exhaustion/graphical-information/>.
- [5] ARIN. Ipv4 addressing options. <https://www.arin.net/resources/guide/ipv4/>.
- [6] M. Bagnulo, P. Matthews, and I. van Beijnum. Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers. RFC 6146 (Proposed Standard), Apr. 2011.
- [7] M. Bagnulo, A. Sullivan, P. Matthews, and I. van Beijnum. DNS64: DNS Extensions for Network Address Translation from IPv6 Clients to IPv4 Servers. RFC 6147 (Proposed Standard), Apr. 2011.
- [8] C. Bao, C. Huitema, M. Bagnulo, M. Boucadair, and X. Li. IPv6 Addressing of IPv4/IPv6 Translators. RFC 6052 (Proposed Standard), Oct. 2010.
- [9] S. Blake-Wilson, M. Nystrom, D. Hopwood, J. Mikkelsen, and T. Wright. Transport Layer Security (TLS) Extensions. RFC 4366 (Proposed Standard), Apr. 2006. Obsoleted by RFCs 5246, 6066, updated by RFC 5746.
- [10] BMV2. Implementing your switch target with bmv2. <http://www.bmv2.org/>.
- [11] R. Braden, D. Borman, and C. Partridge. Computing the Internet checksum. RFC 1071 (Informational), Sept. 1988. Updated by RFC 1141.
- [12] S. Bradner. Key words for use in RFCs to Indicate Requirement Levels. RFC 2119 (Best Current Practice), Mar. 1997. Updated by RFC 8174.
- [13] G. Camarillo, S. Srinivasan, R. Even, and J. Urpalainen. Conference Event Package Data Format Extension for Centralized Conferencing (XCON). RFC 6502 (Proposed Standard), Mar. 2012.
- [14] CISCO. 6lab - the place to monitor ipv6 adoption. <https://6lab.cisco.com/stats/>.
- [15] A. Conta, S. Deering, and M. Gupta (Ed.). Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification. RFC 4443 (Internet Standard), Mar. 2006. Updated by RFC 4884.
- [16] M. Crawford. Transmission of IPv6 Packets over Ethernet Networks. RFC 2464 (Proposed Standard), Dec. 1998. Updated by RFCs 6085, 8064.
- [17] S. Deering, W. Fenner, and B. Haberman. Multicast Listener Discovery (MLD) for IPv6. RFC 2710 (Proposed Standard), Oct. 1999. Updated by RFCs 3590, 3810.

- [18] S. Deering and R. Hinden. Internet Protocol, Version 6 (IPv6) Specification. RFC 2460 (Draft Standard), Dec. 1998. Obsoleted by RFC 8200, updated by RFCs 5095, 5722, 5871, 6437, 6564, 6935, 6946, 7045, 7112.
- [19] J. Dugan, S. Elliott, B. A. Mah, J. Poskanzer, and K. Prabhu. iperf - the ultimate speed test tool for tcp, udp and sctp. <https://iperf.fr/>. Requested on 2019-08-19.
- [20] R. Fielding, J. Gettys, J. Mogul, H. Frystyk, L. Masinter, P. Leach, and T. Berners-Lee. Hypertext Transfer Protocol – HTTP/1.1. RFC 2616 (Draft Standard), June 1999. Obsoleted by RFCs 7230, 7231, 7232, 7233, 7234, 7235, updated by RFCs 2817, 5785, 6266, 6585.
- [21] L. Foundation. Open vswitch. <https://www.openvswitch.org/>.
- [22] Google. Ipv6 - google. <https://www.google.com/intl/en/ipv6/statistics.html>.
- [23] S. P. D. L. V. T. T. B. Hendrik Züllig. P4-programming on an fpga, semester thesis sa-2019-02. https://gitlab.ethz.ch/nsg/student-projects/sa-2019-02_p4_programming_sume_netfpga/blob/master/SA-2019-02.pdf.
- [24] R. Hinden and S. Deering. IP Version 6 Addressing Architecture. RFC 4291 (Draft Standard), Feb. 2006. Updated by RFCs 5952, 6052, 7136, 7346, 7371, 8064.
- [25] C. Hornig. A Standard for the Transmission of IP Datagrams over Ethernet Networks. RFC 894 (Internet Standard), Apr. 1984.
- [26] G. Huston. Ipv4 address report. <https://ipv4.potaroo.net/>. Requested on 2019-08-18.
- [27] G. Huston, A. Lord, and P. Smith. IPv6 Address Prefix Reserved for Documentation. RFC 3849 (Informational), July 2004.
- [28] E. Juskevicius. Definition of IETF Working Group Document States. RFC 6174 (Informational), Mar. 2011.
- [29] LACNIC. Ipv4 depletion phases. <https://www.lacnic.net/1039/1/lacnic/ipv4-depletion-phases>.
- [30] X. Li, C. Bao, and F. Baker. IP/ICMP Translation Algorithm. RFC 6145 (Proposed Standard), Apr. 2011. Obsoleted by RFC 7915, updated by RFCs 6791, 7757.
- [31] N. Lutchansky. Tayga - simple, no-fuss nat64 for linux. <http://www.litech.org/tayga/>.
- [32] J. McCann, S. Deering, J. Mogul, and R. Hinden (Ed.). Path MTU Discovery for IP version 6. RFC 8201 (Internet Standard), July 2017.
- [33] N. Mexico. Jool an open source siit and nat64 for linux. <https://www.jool.mx/en/index.html>.
- [34] J. Mogul and S. Deering. Path MTU discovery. RFC 1191 (Draft Standard), Nov. 1990.
- [35] T. Narten, E. Nordmark, W. Simpson, and H. Soliman. Neighbor Discovery for IP version 6 (IPv6). RFC 4861 (Draft Standard), Sept. 2007. Updated by RFCs 5942, 6980, 7048, 7527, 7559, 8028, 8319, 8425.
- [36] NetFPGA. P4-netfpga-public repository at github. <https://github.com/NetFPGA/P4-NetFPGA-public>.
- [37] A. Networks. Arista 7170 series. <https://www.arista.com/en/products/7170-series>. Requested on 2019-08-21.
- [38] B. Networks. Barefoot tofino. <https://www.barefootnetworks.com/products/brief-tofino/>.

- [39] B. Networks. Barefoot tofino2. <https://barefootnetworks.com/products/brief-tofino-2/>.
- [40] NGINX. Nginx | high performance load balancer, web server, & reverse proxy. <https://www.nginx.com/>.
- [41] E. Nordmark and R. Gilligan. Basic Transition Mechanisms for IPv6 Hosts and Routers. RFC 4213 (Proposed Standard), Oct. 2005.
- [42] D. Plummer. An Ethernet Address Resolution Protocol: Or Converting Network Protocol Addresses to 48.bit Ethernet Address for Transmission on Ethernet Hardware. RFC 826 (Internet Standard), Nov. 1982. Updated by RFCs 5227, 5494.
- [43] J. Postel. User Datagram Protocol. RFC 768 (Internet Standard), Aug. 1980.
- [44] J. Postel. Internet Control Message Protocol. RFC 792 (Internet Standard), Sept. 1981. Updated by RFCs 950, 4884, 6633, 6918.
- [45] J. Postel. Internet Protocol. RFC 791 (Internet Standard), Sept. 1981. Updated by RFCs 1349, 2474, 6864.
- [46] J. Postel. Transmission Control Protocol. RFC 793 (Internet Standard), Sept. 1981. Updated by RFCs 1122, 3168, 6093, 6528.
- [47] Y. Rekhter, B. Moskowitz, D. Karrenberg, G. J. de Groot, and E. Lear. Address Allocation for Private Internets. RFC 1918 (Best Current Practice), Feb. 1996. Updated by RFC 6761.
- [48] E. Rescorla. The Transport Layer Security (TLS) Protocol Version 1.3. RFC 8446 (Proposed Standard), Aug. 2018.
- [49] G. Rieger. socat - multipurpose relay. <http://www.dest-unreach.org/socat/>. Requested on 2019-08-19.
- [50] RIPE. ipv4 exhaustion. <https://www.ripe.net/publications/ipv6-info-centre/about-ipv6/ipv4-exhaustion>.
- [51] N. Schottelius. Add access to table keys. <https://github.com/p4lang/p4-spec/issues/745>.
- [52] N. Schottelius. Casting bit<16> to bit<32> in checksum causes incorrect json to be generated. <https://github.com/p4lang/p4c/issues/1765>.
- [53] N. Schottelius. Extern for checksum'ing payload (p4-netpfga-public). <https://github.com/NetFPGA/P4-NetFPGA-public/issues/13>.
- [54] N. Schottelius. High speed nat64 in p4 (git repository). https://gitlab.ethz.ch/nsg/student-projects/ma-2019-19_high_speed_nat64_with_p4.
- [55] N. Schottelius and S. Plocher. Implementation of a layer 7 ipv4 to ipv6 reverse proxy. Protected git repository <https://gitlab.ethz.ch/nicosc/sdn-nat64/>, part of the Advanced topics in communication networks course fall 2019, <https://adv-net.ethz.ch/>, 2018.
- [56] P. Srisuresh and K. Egevang. Traditional IP Network Address Translator (Traditional NAT). RFC 3022 (Informational), Jan. 2001.
- [57] theojeppen. Get size of header. <https://github.com/p4lang/p4-spec/issues/660>.
- [58] ungleich. Die ipv4, die! <https://ungleich.ch/en-us/cms/blog/2019/01/09/die-ipv4-die/>.
- [59] ungleich. The ungleich network infrastructure. https://redmine.ungleich.ch/projects/open-infrastructure/wiki/The_ungleich_network_infrastructure. Requested on 2019-08-18.

- [60] L. Vanbever. Programming network data planes. https://github.com/nsg-ethz/p4-learning/blob/master/slides/02_p4_env.pdf.
- [61] Various. Should i use python 2 or python 3 for my development activity? <https://wiki.python.org/moin/Python2orPython3>. Requested on 2019-08-19.
- [62] E. Vyncke. Ipv6 deployment aggregated status. <https://www.vyncke.org/ipv6status/>.
- [63] Wikipedia. Ipv4 header checksum. https://en.wikipedia.org/wiki/IPv4_header_checksum. Requested on 2019-08-12.
- [64] Wikipedia. Ipv6 transition mechanism. https://en.wikipedia.org/wiki/IPv6_transition_mechanism. As requested on 2019-08-08.
- [65] Wikipedia. Jumbo frame. https://en.wikipedia.org/wiki/Jumbo_frame. Requested on 2019-08-15.
- [66] Wikipedia. Solicited-node multicast address. https://en.wikipedia.org/wiki/Solicited-node_multicast_address. Requested on 2019-08-13.
- [67] N. Zilberman, Y. Audzevich, G. A. Covington, and A. W. Moore. Netfpga sume: Toward 100 gbps as research commodity. *IEEE Micro*, vol.34, no.5, pp.32-41, Sept.-Oct. 2014, doi: 10.1109/MM.2014.61.