IOGbE monitoring

Ariën Vijn arien.vijn@ams-ix.net

Demo Setup



- Traffic generator
 - Screen with Windows interface

Demo Setup



- Traffic generator
 - Screen with Windows interface.
- Ethernet switch called "lazerix"
 - One of the terminal window (look at prompt)

Demo Setup



- Traffic generator
 - Screen with Windows interface
- Ethernet switch called "lazerix"
 - One of the terminal window (look at prompt)
- We'll look at the traffic on the wire from 2:4:1 to 2/1
 - Monitorbox called "hysterix"

Let's have a look

IOGbE monitoring (some slides now) Ariën Vijn <u>arien.vijn@ams-ix.net</u>

Agenda

- An example.
- The problems with passive I0GbE monitoring.
 - Performance of generic computer hardware.
 - Issues with the mirror port feature.
- Overview of our solution.
 - Hardware architecture.
 - Homegrown firmware and software.
- A bit more detail.
- More to demo.

Agenda

• An example.

- The problems with passive I0GbE monitoring.
 - Performance of generic computer hardware.
 - Issues with the mirror port feature.
- Overview of our solution.
 - Hardware architecture.
 - Homegrown firmware and software.
- A bit more detail.
- More to demo.







- Normal situation
 - I0GbE member router



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park
- New switch



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park
- New switch
 - Connected to the switch park



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park
- New switch
 - Connected to the switch park



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park
- New switch
 - Connected to the switch park
 - Moved in production



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park
- New switch
 - Connected to the switch park
 - Moved in production
 - Member feels unhappy



- Normal situation
 - I0GbE member router
 - Photonic Cross Connect (PXC)
 - Redundant switch park
- New switch
 - Connected to the switch park
 - Moved in production
 - Member feels unhappy
 - Something must be wrong with the new switch



 Debugging the switch (the vendor's way)



- Debugging the switch (the vendor's way)
 - Layer-2 access list accounting



- Debugging the switch (the vendor's way)
 - Layer-2 access list accounting
 - Misses part of the egress path
 - "Cannot be the problem because the same chips serve the ingress path."
 - This took weeks



- Debugging the network (our way)
 - Passive monitoring
 - Watching the traffic as it passes by.



- Debugging the network (our way)
 - Passive monitoring
 - Watching the traffic as it passes by.
- What did we see?
 - A lot of incoming pause frames!



- Debugging the network (our way)
 - Passive monitoring
 - Watching the traffic as it passes by.
- What did we see?
 - A lot of incoming pause frames!
 - We had a problem in our switch park!
 - Issue found within two minutes after invoking the monitor system.



- Debugging the network (our way)
 - Passive monitoring
 - Watching the traffic as it passes by.
- What did we see?
 - A lot of incoming pause frames!
 - We had a problem in our switch park!
 - This took 2 minutes after the monitor system was installed.
 - Software upgrade solved the problem.

Agenda

- An example.
- The problems with passive I0GbE monitoring.
 - Performance of generic computer hardware.
 - Issues with the mirror port feature.
- Overview of our solution.
 - Hardware architecture.
 - Homegrown firmware and software.
- A bit more detail.
- More to demo.

Common setup



- Mirror port feature
- General purpose PC
 - Tools based on libpcap



- I0GbE frame rates
 - 14.8 million frames per second



- IOGbE frame rates
 - 14.8 million frames per second
- Test setup
 - Traffic generator to 'normal' Intel 10GbE NIC



- IOGbE frame rates
 - 14.8 million frames per second
- Test setup
 - Traffic generator to 'normal' Intel I0GbE NIC
- One out of one million.
 - Reduced data rate until all frames were captured
 - Note the more or less constant frame rate of 1.5Mfps
 - PCI-X bus speed harms larger sized frame.



- I0GbE frame rates
 - 14.8 million frames per second
- Test setup
 - Traffic generator to 'normal' Intel I0GbE NIC
- One out of one million frames
 - Reduced data rate until all frames were captured
 - Note the more or less constant frame rate of 1.5Mfps
 - PCI-X bus speed harms larger frame sizes.



- One out of one million frames
 - Reduced data rate until all frames were captured
 - Note the more or less constant frame rate of 1.5Mfps
 - PCI-X bus speed harms larger frame sizes.
- One out of 16 frames
 - Lots of context switching
 - Again a more or less constant frame rate.



- One out of 16 frames
 - Lot's of context switching
 - Again a more or less constant frame rate.
- Conclusion
 - General purpose hardware is not fast enough.

Mirror Port



- Device Under Test (DUT)
 - Foundry Networks RX series

Mirror Port



- Device Under Test (DUT)
 - Foundry Networks RX series


- Device Under Test (DUT)
 - Foundry Networks RX series
- Test I
 - Throughput without monitoring.
 - 100% at all frame rates



- Device Under Test (DUT)
 - Foundry Networks RX series
- Test 2
 - Ingress monitoring
 - Input traffic on port I/I is copied to port I/3
- Test 3
 - Egress monitoring
 - Output traffic to port I/I is copied to I/3.



- Device Under Test (DUT)
 - Foundry Networks RX series
- Test 2
 - Ingress monitoring
 - Input traffic on port I/I is copied to port I/3
- Test 3
 - Egress monitoring
 - Output traffic to port I/I is copied to I/3.



- Device Under Test (DUT)
 - Foundry Networks RX series
- Test 3
 - Egress monitoring:
 - Output traffic to port I/I is copied to I/3.
- Test 4
 - Both ingress and egress monitoring to one mirror port:
 - Bi-directional traffic at port 1/1 is copied to port 1/3.



- Device Under Test (DUT)
 - Foundry Networks RX series
- Test 4
 - Both ingress and egress monitoring to one mirror port:
 - Bi-directional traffic at port 1/1 is copied to port 1/3.
- Test 5
 - Ingress traffic at port I/I is copied to port I/3
 - Output traffic to port I/I is copied to I/4.



- Device Under Test (DUT)
 - Foundry Networks RX series
- Conclusion
 - No mirroring: no frame loss
 - Mirroring: frame loss above certain data rates.
 - Best case: 70% throughput
 - Worst case: 36% throughput



- Device Under Test (DUT)
 - Foundry Networks RX series
- Conclusion
 - No mirroring: no frame loss
 - Mirroring: frame loss above certain data rates.
 - Best case: 70% throughput
 - Worst case: 36% throughput
- Result of design decisions



- Device Under Test (DUT)
 - Foundry Networks RX series
- Conclusion
 - No mirroring: no frame loss
 - Mirroring: frame loss above certain data rates.
 - Best case: 70% throughput
 - Worst case: 36% throughput
- Result of design decisions
 - Replication takes place in packet processor.

Agenda

- An example.
- The problems with passive I0GbE monitoring.
 - Performance of generic computer hardware.
 - Issues with the mirror port feature.
- Overview of our solution.
 - Hardware architecture.
 - Homegrown firmware and software.
- A bit more detail.
- More to demo.

Solution



- We need to monitor in-line.
 - We can use our PXCs...

Solution



- We need to monitor in-line.
 - We can use our PXCs...

Solution



- We need to monitor in-line.
 - We can use our PXCs... to insert some kind of tap or repeater.



- We need to monitor in-line.
 - We can use our PXCs... to insert some kind of tap or repeater.
- Two full duplex transceivers
 - In-line, for signal regeneration



- We need to monitor in-line.
 - We can use our PXCs... to insert some kind of tap or repeater.
- Two full duplex transceivers
 - In-line, for signal regeneration
- Repeater and tap
 - Forwards traffic between transceivers and copies it for monitoring.



- Two full duplex transceivers
 - In-line, for signal regeneration
- Repeater and tap
 - Forwards traffic between transceivers and copies it for monitoring.
- User adjustable filter
 - Data rate reduction



- Two full duplex transceivers
- Repeater and tap
 - Forwards traffic between transceivers and copies it for monitoring.
- User adjustable filter
 - Data rate reduction
- Standard interface bus
 - PCI-X bus



- Two full duplex transceivers
- Repeater and tap
 - Forwards traffic between transceivers and copies it for monitoring.
- User adjustable filter
 - Data rate reduction
- Standard interface bus
 - PCI-X bus
- Software
 - Kernel module to make it look like a 'normal' ethernet interface.



- Two full duplex transceivers
- Repeater and tap
 - Forwards traffic between transceivers and copies it for monitoring.
- User adjustable filter
 - Data rate reduction
- Standard interface bus
 - PCI-X bus
- Software
 - Kernel module to make it 'look' like a normal ethernet interface.



- Summarizing
 - Full line rate hardware



- Summarizing
 - Full line rate hardware
 - Lower the data rate so that the PC can handle it by either filtering or sampling



- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PI0 IDS/IPS system
 - We wrote our own firmware and software for dynamic monitoring



- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PIO IDS/IPS system
 - We wrote our own firm and software for dynamic monitoring
- Front-end FPGA
 - Full line rate tap and repeater



- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PIO IDS/IPS system
 - We wrote our own firm and software for dynamic monitoring
- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - 128-bit wide words

0	preamble / SFD	dest. MAC addr.	src

- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PIO IDS/IPS system
 - We wrote our own firm and software for dynamic monitoring
- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - 128-bit wide words

	preamble / SFD			dest. MAC addr.			SrC	
							_	
	MAC addr.	type	4 ⊦ 4 L	Tos	total len.	ID	F/F	TTL pro to.
_		-	-				-	

- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PIO IDS/IPS system
 - We wrote our own firm and software for dynamic monitoring
- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - 128-bit wide words

preamble / SFD				dest	i. MAC a	SrC	
MAC addr. type		4 L TOS	total len.	ID F/F		TTL pro	
HCS src. IP addr.		1		IP payload			
HCS src. IP addr.							
HCS	src. IF	Paddr.	dst. IF	^D addr		P payloa	ad
CS	src. IF	Paddr.	dst. IF	^D addr		P payloa	ad
HCS	src. IF	Paddr.	dst. IF	P addr		P payloa	ad
HCS	src. IF	Paddr.	dst. IF	P addr		P payloa	ad
HCS	src. IF	Paddr.	dst. IF	P addr		P payloa	ad

- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PIO IDS/IPS system
 - We wrote our own firm and software for dynamic monitoring
- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - 128-bit wide words

I								
0	preamble / SFD			dest. MAC addr.			SrC	
ſ								
1	MAC addr. type		4 L TOS	total len.	ID	F/F	TTL pro to.	
ſ								
2	HCS src. IP addr.		dst. IF	9 addr	I	P payload		
1								1
i '								
n	IP payload				FC	S		

- ForceI0 MTP-I0G
 - Network Interface Card NIC
 - Part of PIO IDS/IPS system
 - We wrote our own firm and software for dynamic monitoring
- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - 128-bit wide words



- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - I 28-bit wide words
- Back-end FPGA
 - The heart of the system
 - MISD machine
 - Parallel processing on the 128-bit words.



- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - I 28-bit wide words
- Back-end FPGA
 - The heart of the system
 - MISD machine
 - Parallel processing on the 128-bit words.
- 128-bit word register (Single Data).



- Front-end FPGA
 - Full line rate tap and repeater
- Tapped traffic
 - I 28-bit wide words
- Back-end FPGA
 - The heart of the system
 - MISD machine
 - Parallel processing on the 128-bit words.
- 128-bit word register (Single Data).
- Multiple instruction registers.



- Back-end FPGA
 - The heart of the system
 - MISD machine
 - Parallel processing on the 128-bit words.
- 128-bit word register (Single Data).
- Multiple instruction registers.
- Multiple processors.



- Back-end FPGA
 - The heart of the system
 - MISD machine
 - Parallel processing on the 128-bit words.
- 128-bit word register (Single Data).
- Multiple instruction registers.
- Multiple processors.
 - Result is reduced to a match or no-match signal.



- Back-end FPGA
 - The heart of the system
 - MISD machine
 - Parallel processing on the 128-bit words.
- 128-bit word register (Single Data).
- Multiple instruction registers.
- Multiple processors.
 - Result is reduced to a match or no-match signal.
 - Frame is copied down to the PCI-X interface upon a match.



• Hardware



- Hardware
- Software
 - Two kernel modules.



- Hardware
- Software
 - Two kernel modules.
 - Network subsystem
 - libpcap


- Hardware
- Software
 - Two kernel modules.
 - Network subsystem
 - libpcap
 - Mapped memory
 - I/O



- Hardware
- Software
 - Two kernel modules.
 - Network subsystem
 - libpcap
 - Mapped memory
 - I/O
 - Counter readings



- Hardware
- Software
 - Two kernel modules.
 - Network subsystem
 - libpcap
 - Mapped memory
 - I/O
 - Counter readings
 - Dynamic filter settings



- Hardware
- Software part
 - Two kernel modules.
 - Network subsystem
 - libpcap
 - Mapped memory
 - I/O
 - Counter readings
 - Dynamic filter settings.
- Takes input from stdin.



- Takes input from stdin.
- Parses the input.
 - Chops addresses in patterns and offsets.
 - Results in an internal expression representation.

input	
not ether src 00:01:02:03:04:05	5
lexer/parser	
!(N00e0001&N0100203&N012040)5);
Quine-McCluskey	
!N00e0001I!N0100203I!N01204	05
register mapper	

- Takes input from stdin.
- Parses the input.
 - Chops addresses in patterns and offsets.
 - Results in an internal expression representation.
- Quine-McCluskey
 - Transforms the expression into its disjunctive normal form.



- Takes input from stdin.
- Parses the input.
 - Chops addresses in patterns and offsets.
 - Results in an internal expression representation.
- Quine-McCluskey
 - Transforms the expression into its disjunctive normal form.
- Maps each term into a instruction register in the back-end FPGA.

Agenda

- An example.
- The problems with passive I0GbE monitoring.
 - Performance of generic computer hardware.
 - Issues with the mirror port feature.
- Overview of our solution.
 - Hardware architecture.
 - Homegrown firmware and software.
- A bit more detail.
- More to demo.

not ether src 00:01:02:03:04:05

• Let's see how this filter expression is handled in detail.



- Let's see how this filter expression is handled in detail.
- The first I28bit word contains the first two bytes of the ethernet source address.



- Let's see how this filter expression is handled in detail.
- The first I28bit word contains the first two bytes of the ethernet source address.
- The second I28bit word contains the remaining address.
 - One "processor" can process two bytes.
 - Hence we need two here.



- So far so good, but what about this not operator?
 - Programmable network behind processors?
 - Limited in both speed and scalability.



- So far so good, but what about this not operator?
 - Programmable network behind processors?
 - Limited in both speed and scalability.
 - Do it is software by transforming the expression into sum of products
 - Disjunctive Normal Form.
 - Quine-McCluskey
 - Library



- So far so good, but what about this not operator?
 - Programmable network behind processors?
 - Limited in both speed and scalability.
 - Do it is software by transforming the expression into sum of products
 - Disjunctive Normal Form.
 - Quine-McCluskey
 - Library
- Registers

pattern (16 bits)	offset (11 bits)	instr. A L
	AND with previou las	s block

- Registers
 - Contain the pattern to match (16 bits)

pattern (16 bits)	offset (11 bits)	instr. A L
	AND with previous block —	

- Registers
 - Contain the pattern to match (16 bits)
 - The offset (II bits)
 - 2048 bytes
 - 1522 byte frames

	pattern (16 bits)	offset (11 bits)	instr. A L	
		AND with previou la	us block	
Instr.	meaning			
0.	Full 16 bit match.			
1.	Full 16 bit not match. The outcome is true if the data does not match the pattern.			
2.	Partial match, left half of pattern field contains the data to match, right half contains bit mask.			
3.	Partial not match, left half of pattern field contains data to match, right half contains bitmask.			
4.	Partial match, same as instruction 2 but left and right are reversed.			
5.	Partial not match, same as instruction 4 but left and right are reversed.			
6.	Reserved for future use.			
	Reserved for future use.			

- Registers
 - Contain the pattern to match (16 bits)
 - The offset (II bits)
 - 2048 bytes
 - 1522 byte frames
 - Instruction (3 bits)

	pattern (16 bits)	offset (11 bits)	instr. A	
		AND with previou la	us block ———	
Instr.	meaning			
0.	Full 16 bit match.			
1.	Full 16 bit not match. The outcome is true if the data does not match the pattern.			
2.	Partial match, left half of pattern field contains the data to match, right half contains bit mask.			
3.	Partial not match, left half of pattern field contains data to match, right half contains bitmask.			
4.	Partial match, same as instruction 2 but left and right are re- versed.			
5.	Partial not match, same as instruction 4 but left and right are reversed.			
6.	Reserved for future use.			
	Reserved for future use.			

- Registers
 - Contain the pattern to match (16 bits)
 - The offset (11 bits)
 - 2048 bytes
 - 1522 byte frames
 - Instruction (3 bits)
 - Last two bits to handle products that do not one processor.



- Registers
 - Contain the pattern to match (16 bits)
 - The offset (II bits)
 - 2048 bytes
 - 1522 byte frames
 - Instruction (3 bits)
 - Last two bits to handle products that do not one processor.

R	mask	pattern	offset	instr. A L
L	pattern	mask	offset	instr. A L
			AND with Ia	previous —
Instr.	meaning			
0.	Full 16 bit match.			
1.	Full 16 bit not match. The outcome is true if the data does not match the pattern.			
2.	Partial match, left half of pattern field contains the data to match, right half contains bit mask.			
3.	Partial not match, left half of pattern field contains data to match, right half contains bitmask.			
4.	Partial match, same as instruction 2 but left and right are reversed.			
5.	Partial not match, same as instruction 4 but left and right are reversed.			
6.	Reserved for future use.			
7.	Reserved for future use.			

- Registers
 - Contain the pattern to match (16 bits)
 - The offset (II bits)
 - 2048 bytes
 - 1522 byte frames
 - Instruction (3 bits)
 - Last two bits to handle products that do not one processor.
 - Partial matches
 - Via a bit mask.
 - Per 8 bit patterns.





Agenda

- An example.
- The problems with passive I0GbE monitoring.
 - Performance of generic computer hardware.
 - Issues with the mirror port feature.
- Overview of our solution.
 - Hardware architecture.
 - Homegrown firmware and software.
- A bit more detail.
- More to demo.

Demo



- Stream I
 - 15,000,000 frames.
 - Source MAC: 00:01:02:03:04:05.
- Stream 2
 - I frame from another source MAC.

Conclusion...

• System works.

Conclusion...

- System works.
 - It has been useful already...

... & future development

- System works.
 - It has been useful already, but it is still in development.
 - BPF to instructions for the MISD processor.
 - Integrate the system in libpcap.
 - Traffic generator.