

Monitorovanie sietí na rýchlosti 100 Gb/s

(Internet a Technologie 12)

Lukáš Kekely, Viktor Puš, Štěpán Friedl
(kekely, pus, friedl@cesnet.cz)



Praha, 24. 11. 2012

- 1 Motivácia
- 2 Ethernet 100Gb/s
- 3 Existujúce riešenie monitorovania
- 4 Návrh nového riešenia na 100Gb/s
- 5 Záver

- 1 Motivácia
- 2 Ethernet 100Gb/s
- 3 Existujúce riešenie monitorovania
- 4 Návrh nového riešenia na 100Gb/s
- 5 Záver

- rýchly nárast objemu dát tečúcich sieťou
 - stúpa počet služieb využívajúcich siete
 - rastie počet zariadení schopných komunikácie
 - zväčšuje sa objem prenášaných dát
- nutnosť zvýšenia prenosovej rýchlosti!

- rýchly nárast objemu dát tečúcich sieťou
 - stúpa počet služieb využívajúcich siete
 - rastie počet zariadení schopných komunikácie
 - zväčšuje sa objem prenášaných dát
- nutnosť zvýšenia prenosovej rýchlosti!

Ako rýchlo skutočne rastie vyťaženie sietí?

- rýchly nárast objemu dát tečúcich sieťou
 - stúpa počet služieb využívajúcich siete
 - rastie počet zariadení schopných komunikácie
 - zväčšuje sa objem prenášaných dát
- nutnosť zvýšenia prenosovej rýchlosti!

Ako rýchlo skutočne rastie vyťaženie sietí?

⇒ IEEE 802.3TM Industry Connections Ethernet Bandwidth Assessment (júl 2012)

http://www.ieee802.org/3/ad_hoc/bwa/BWA_Report.pdf

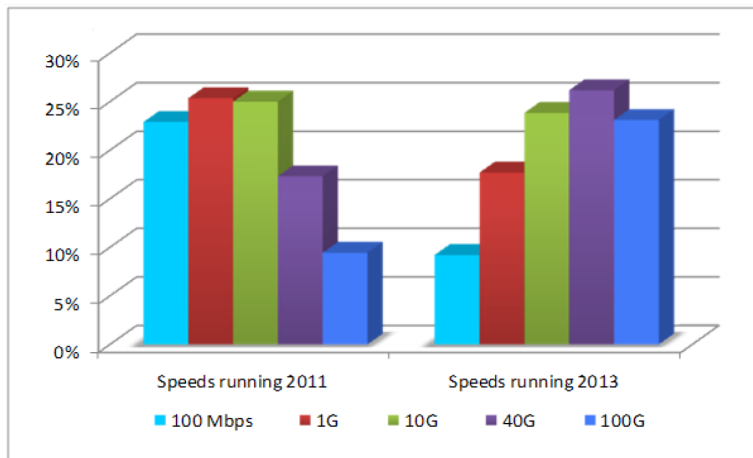


Figure 11—Data center study - percentage of links by speed

- v dátových centrách budú už budúci rok prevládať 40 a 100 Gb/s technológie

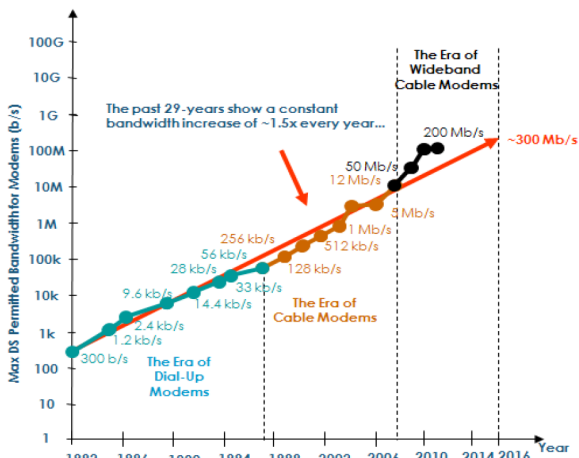


Figure 19—Maximum permitted downstream bandwidth trend

- *Nielsen's Law of Internet bandwidth*: rýchlosť pripojenia high-end užívateľov vzrastie o 50 % ročne

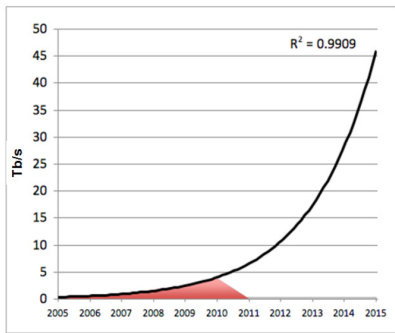
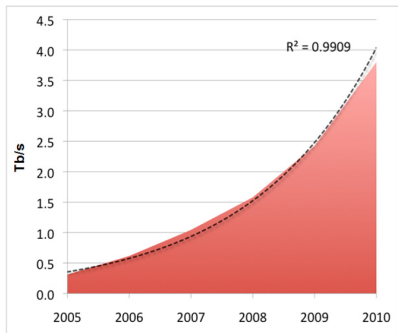


Figure 39—Five year peak European IXP traffic projection

- nárast prenosových rýchlostí priemerne o asi 57 % ročne

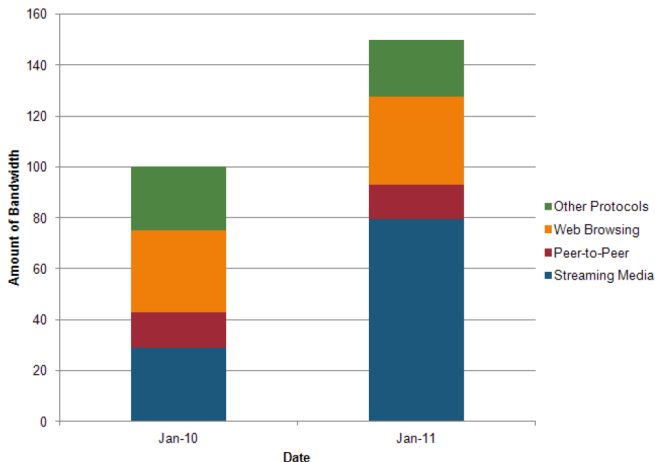


Figure 23—Mix of traffic type vs. time

- výrazne rastie pomer multimediálneho obsahu
- CESNET predviedol technológiu na prenos 8K videa

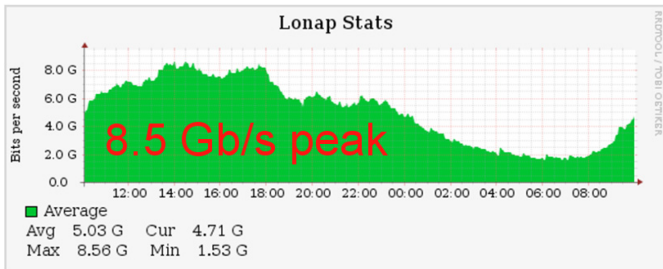


Figure 34—LONAP (London) traffic on a 'normal' weekday

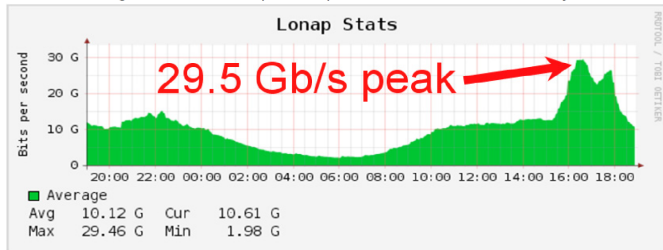
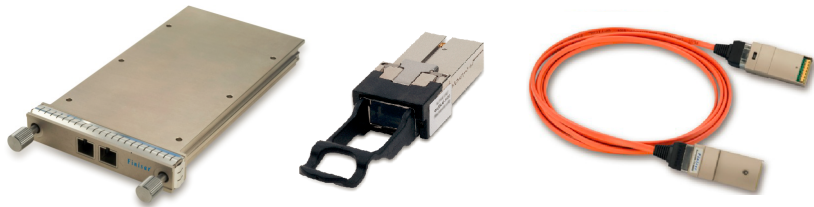


Figure 35—LONAP (London) traffic during the World Cup 2010 England vs. Slovenia

- 1 Motivácia
- 2 Ethernet 100Gb/s
- 3 Existujúce riešenie monitorovania
- 4 Návrh nového riešenia na 100Gb/s
- 5 Záver

- štandard 802.3ba prijatý IEEE v júni 2010
 - dve prenosové rýchlosti – 40 a 100 Gb/s
 - na spojovej vrstve (L2) zachováva väčšinu vlastností
 - enkapsulácia dát a formát rámca
 - minimálna a maximálna veľkosť rámca
 - adresovanie pomocou MAC adries
 - detekcia chýb pomocou CRC
 - na fyzickej vrstve (L1) nastali zmeny
 - niekoľko paralelných fyzických ciest
 - rámec je prenášaný naraz všetkými cestami
 - pridané značky na zarovnanie posunu ciest
- ⇒ potreba novej realizácie fyzickej vrstvy
- ⇒ 100 GE nie je len 10-krát rýchlejší 10 GE

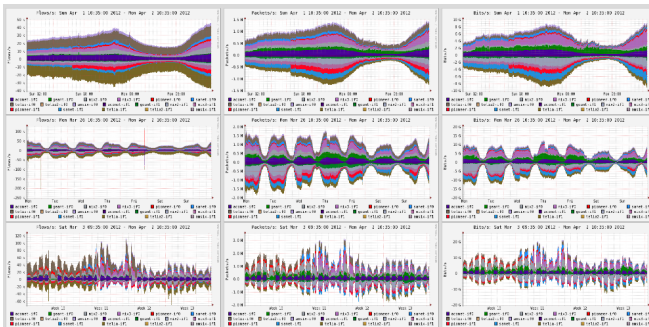
- dostupné optické moduly a optické médiá
- paralelný prenos rámca viacerými cestami
 - 1 vlákno, 4 vlnové dĺžky po 25 Gb/s (= 100 Gb/s)
 - 1 vlákno, 4 vlnové dĺžky po 10 Gb/s (= 40 Gb/s)
 - 10 vlákien po 10 Gb/s (= 100 Gb/s)
 - 4 vlákna po 10 Gb/s (= 40 Gb/s)
 - 1 vlákno, 10 vlnových dĺžok po 10 Gb/s (= 100 Gb/s)
- rôzne realizácie optických transceiverov v 2 generáciách

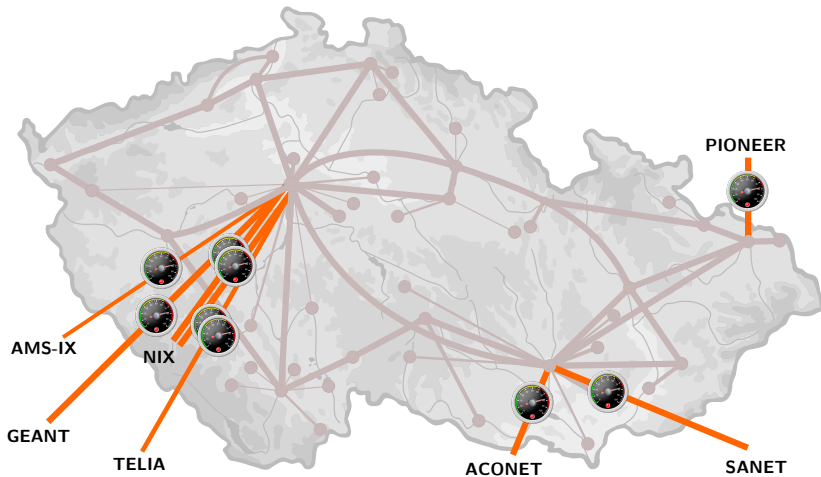


* obrázky sú zo stránok firmy Finisar

- 1 Motivácia
- 2 Ethernet 100Gb/s
- 3 Existujúce riešenie monitorovania**
- 4 Návrh nového riešenia na 100Gb/s
- 5 Záver

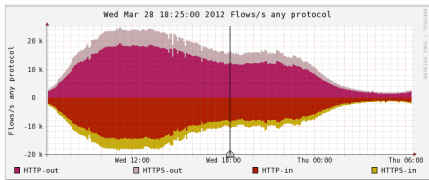
- Oddělení nástrojů pro monitoring a konfiguraci
- monitorovacia infraštruktúra na zber informácií o IP tokoch
- celkovo 9 meracích bodov na zber dát
- dáta spracovávané na kolektore
- pokrytie všetkých hraničných liniek siete CESNET2
- podpora monitorovania na plnej rýchlosti liniek
- okrem zberu dát aj vyhodnotenie, reprezentácia a ďalšie

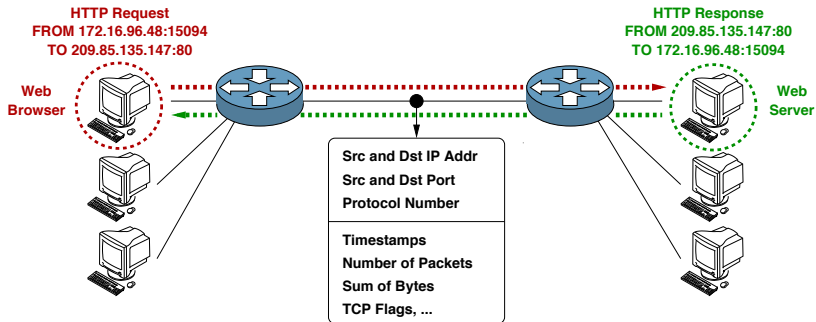




- kto komunikuje s kým, ako dlho, akým protokolom atď
- založené na CISCO NetFlow v5/v9 a IETF IPFIX
 - IPFIX prináša flexibilitu záznamov \Rightarrow možnosť rozšírenia o pokročilejšie informácie (napr. L7)
- monitoring a analýza sieťových tokov v reálnom čase
 - detekcia anomálií, útokov, výpadkov, nakazených PC atď
- ukladanie a prezentácia dlhodobých štatistík
 - agregácia na 5 minútové intervaly (raw cca 140 GB/deň)
 - sledovanie dlhodobého vývoja siete
 - dohľadanie podrobností k riešeniu incidentov

Duration	Proto	Src IP Addr:Port	Dst IP Addr:Port	Flags
2.096	TCP	108.7.1.50:80	108.7.1.50:80	.AP.S.
0.094	TCP	59.173.182.61:49442	59.173.182.61:49442	.AP.S.
0.368	TCP	108.7.1.50:80	59.173.182.61:49440	.AP.S.
0.737	TCP	108.7.1.50:80	59.173.182.61:49434	.AP.S.
0.379	TCP	59.173.182.61:49438	59.173.182.61:49438	.AP.S.
0.296	TCP	108.7.1.50:80	108.7.1.50:80	.AP.S.
0.575	TCP	108.7.1.50:80	108.7.1.50:80	.AP.S.
0.574	TCP	108.7.1.50:80	108.7.1.50:80	.AP.S.
0.451	TCP	108.7.1.50:80	108.7.1.50:80	.AP..
1.281	TCP	108.7.1.50:80	108.7.1.50:80	.AP.SF
1.280	TCP	213.173.41667	213.173.41667	.AP.SF
5.886	TCP	108.7.1.50:80	108.129.1687	.AP..
6.051	TCP	108.7.1.50:80	108.7.1.50:80	.AP..
2.800	TCP	210.242.141.183:1324	210.242.141.183:1324	.AP.S.
2.980	TCP	210.242.141.183:1324	210.242.141.183:1324	.AP.S.
1.693	TCP	108.7.1.50:80	157.242.141.183:1324	.AP.S.
1.778	TCP	108.7.1.50:80	157.242.141.183:1324	.AP.S.
0.604	TCP	157.242.141.183:1325	108.7.1.50:80	.AP.S.
1.990	TCP	157.242.141.183:1324	108.7.1.50:80	.AP.S.





Flow start	Duration	Proto	Src IP Addr:Port		Dst IP Addr:Port	Flags	Packets	Bytes
09:41:21.763	0.101	TCP	172.16.96.48:15094	->	209.85.135.147:80	.AP.SF	4	715
09:41:21.893	0.031	TCP	209.85.135.147:80	->	172.16.96.48:15094	.AP.SF	4	1594

- komoditný server pod Linuxom
- vlastný viacvláknový SW exportér (NetFlow/IPFIX)
- hardvérová sonda COMBOv2 vlastnej výroby
 - PCI Express karta s FPGA Virtex-5

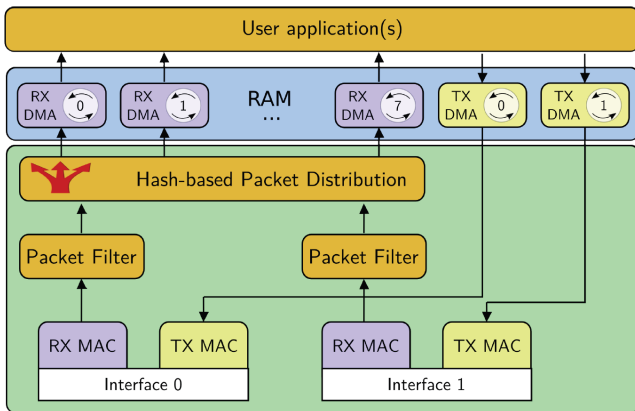


COMBOI-10G2



COMBOI-10G4

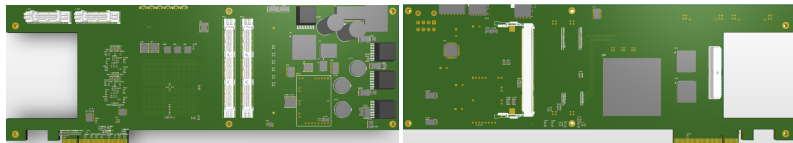
- firmvér pre FPGA na Combov2, akcelerovaná NIC s možnosťou doplnenia vlastného predspracovania paketov
- nastaviteľná distribúcia dát na jadrá
 - základom je hash vybraných políček z hlavičiek
 - distribúcia zachováva toky



- 1 Motivácia
- 2 Ethernet 100Gb/s
- 3 Existujúce riešenie monitorovania
- 4 Návrh nového riešenia na 100Gb/s**
- 5 Záver

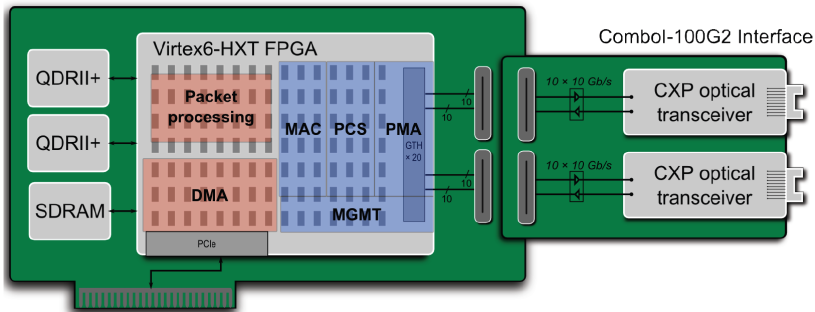
- zachovanie systému meracích bodov a kolektoru
- monitorovanie založené na IETF *IPFIX*
- nutnosť vytvorenia novej karty do meracieho bodu
- rozšírenie monitorovania o ďalšie informácie
- nový model HW/SW kodizajnu v meracom bode

- vynútená zmenou fyzickej vrstvy a nárastom rýchlosti
- základňová PCI Express karta s FPGA Virtex6-HXT
 - 24x rýchly GTH transceiver (11,3 Gb/s každý)
 - 2x PCIe x8 rozhranie generácie 2 (64 Gb/s)
 - 2x QDRII+ pamäť (80 Gb/s)
 - SODIMM slot na DDR3 pamäť
- možnosť pripojiť rôzne interface karty



- jednoduchá zmena fyzických portov karty
- plánované rôzne varianty pre 40 aj 100 Gb/s
 - 1x 40 GbE, 2x 40 GbE
 - 2x 100 GbE
 - rozdeľovací kábel na 6x 40 GbE alebo 24x 10 GbE

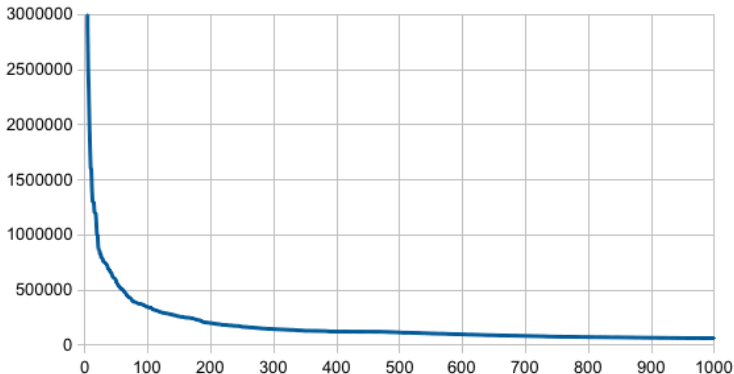
Combo-HXT Main Board



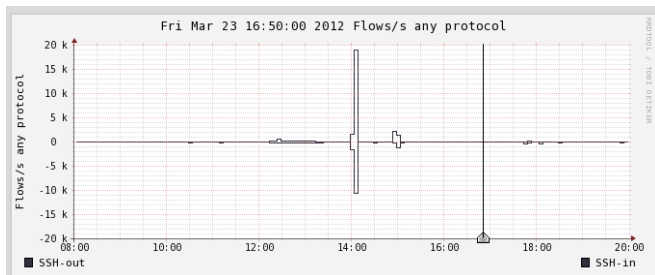
- komoditný server pod Linuxom s Combo-HXT
- hlavné problémy vynucujúce zmenu
 - PCIe kritickým bodom na priepustnosť (len 64 Gb/s)
 - podpora pokročilejšieho monitorovania (detekcia)
 - snaha odľahčiť SW spracovanie paketov
- využiteľné vlastnosti sieťových tokov
 - *heavy-tail* rozloženie tokov
 - nezaujímavosť dátovej časti väčšiny paketov
- riešenie inšpirované *Software Defined Networking*
 - užšie previazanie SW riadenia a HW akcelerácie
 - SW sa stará o riadenie a využíva HW na odľahčenie



- pozorované správanie veľkosti sieťových tokov na linke
- málo najväčších tokov tvorí veľkú časť z toku linkou
- meranie v sieti CESNET2
 - 1 000 najväčších tokov obsahuje cca 30 % paketov
 - 10 000 najväčších tokov obsahuje cca 60 % paketov

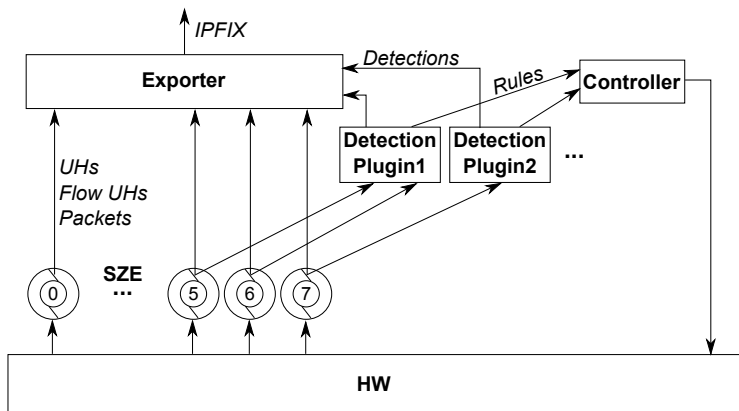


- priemerná veľkosť paketu je cca 600-800 B
- základné monitorovanie potrebuje len metadáta
 - adresy, porty, protokol, časová značka, veľkosť
 - veľkosť metadát je podstatne menšia ako väčšina paketov
- pokročilé monitorovanie nemusí mať všetky pakety
 - nespracovateľné protokoly (napr. šifrované) tvoria asi 10%
 - detekcia SPAMu potrebuje len SMTP (menej ako 1%)
 - detekcia DNS útokov potrebuje len DNS (menej ako 1%)

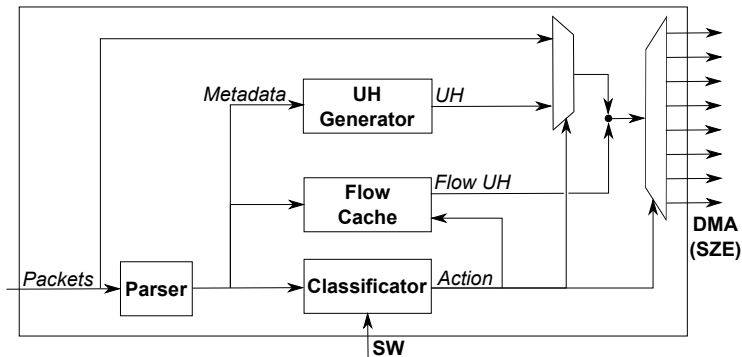


SSH sken z Prahy do sveta

- počítanie a export celkového monitoringu tokov
- jednoduchá rozšíriteľnosť merania (pluginy)
- konfigurácia klasifikácie a akcií v HW
 - dátovo nezaujímavé toky (stačí UH)
 - *heavy-tail* nezaujímavé toky (stačí Flow UH)



- rozšírenie architektúry Hanic
- nastaviteľná klasifikácia paketov s akciou
- schopnosť extrahovať a posilať len metadáta do SW (UH)
- počítanie čiastkového monitoringu tokov (Flow UH)



- 1 Motivácia
- 2 Ethernet 100Gb/s
- 3 Existujúce riešenie monitorovania
- 4 Návrh nového riešenia na 100Gb/s
- 5 Záver

- monitoring nesmie zaostávať za vývojom infraštruktúry
- poskytovanie ešte lepších služieb aj na vyššej rýchlosti

- monitoring nesmie zaostávať za vývojom infraštruktúry
- poskytovanie ešte lepších služieb aj na vyššej rýchlosti

“Optická síť CESNET2 je plně připravena na nasazení přenosové technologie 100 Gb/s, kterou poskytují jen nejpokročilejší infrastruktury.”

Ďakujem za pozornosť.