

MonIPÊ: Serviço de Medições de Desempenho de Redes

Rede Nacional de Ensino e Pesquisa

Alex Soares de Moura – alex@rnp.br



Ministério da
Cultura

Ministério da
Saúde

Ministério da
Educação

Ministério da
**Ciência, Tecnologia
e Inovação**

GT-ER 39

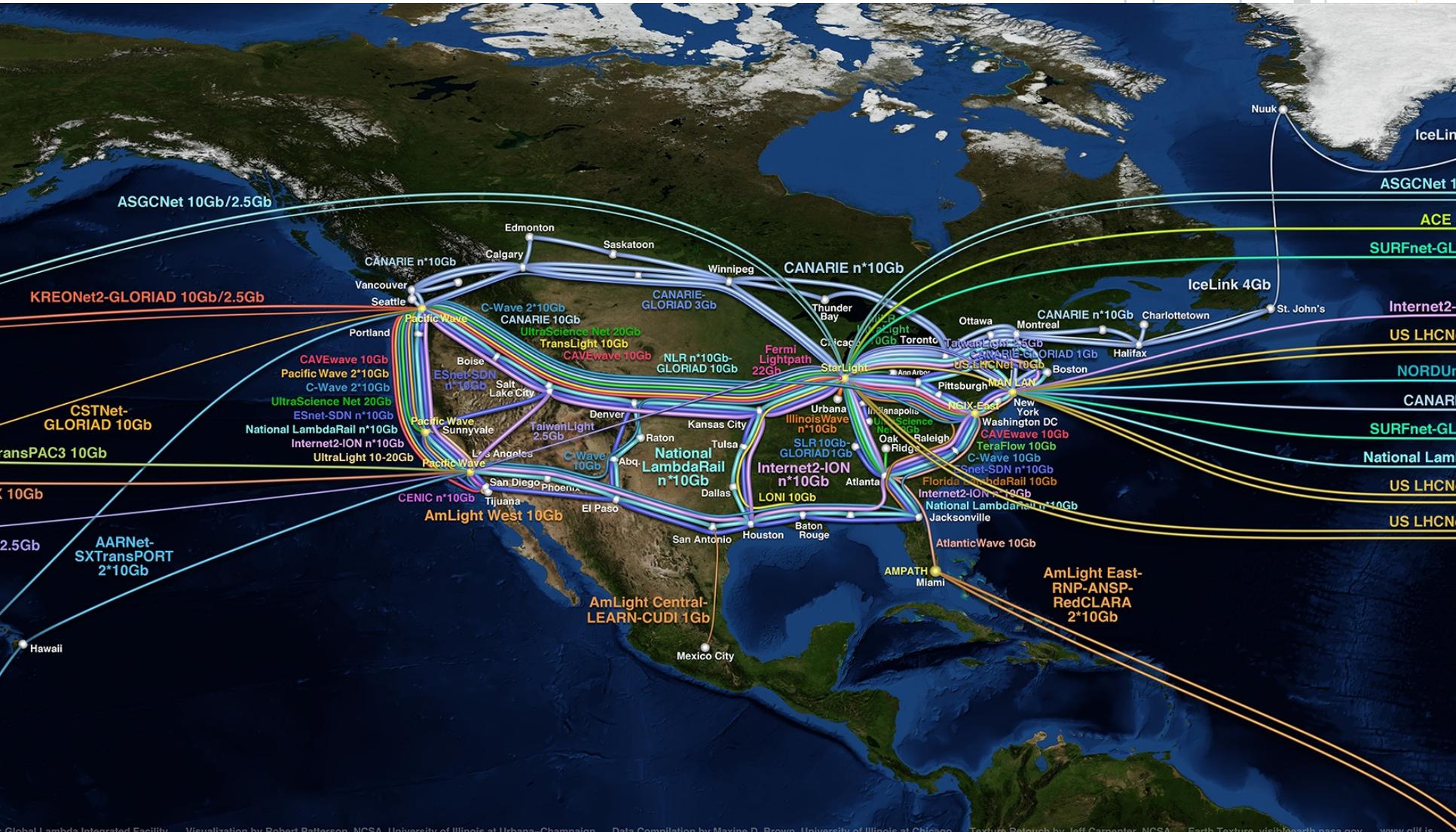
28 a 30 de Maio de 2015, Rio de Janeiro - RJ

Monitorar o desempenho das redes é fundamental para a eficiência das comunicações e aplicações

Com o Serviço MonIPÊ – compatível com o padrão perfSONAR – é possível realizar, com alta precisão, o monitoramento fim a fim do desempenho de redes

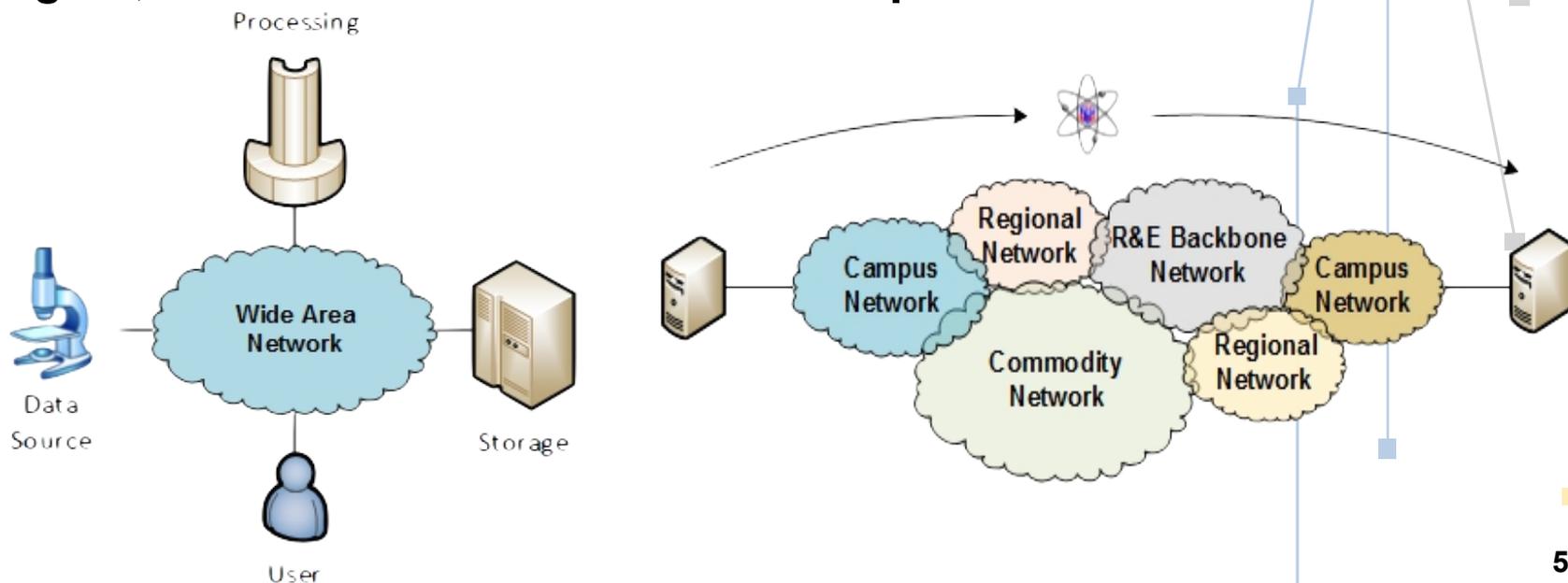
Monitoramento melhor = Redes melhores

O Problema



O Problema

- Este complexo conjunto de redes heterogêneas precisa ser operada de forma integrada, “fim a fim” para suportar colaborações de pesquisa científica distribuídas globalmente
- Na prática, problemas de desempenho são distribuídos
- Quando uma rede tem baixo desempenho, é difícil identificar a origem, e testes de rede local nem sempre são suficientes



Desempenho das redes: expectativas

- Expectativas sobre as transferências de dados
 - Em condições ideais, qual desempenho você **espera** da sua rede?
 - Quanto tempo demora transferir 1TB em diferentes velocidades¹ ? 
- e-Ciência requer **pesquisa colaborativa**, transferências de **volumes de dados crescentes** e tem como requisito **baixo tempo de transferência**²
- Pesquisadores realizam trabalhos científicos colaborativos e **compartilham dados e recursos computacionais**²
- Com essa colaboração, **é necessário movimentar grandes volumes de dados**, da ordem de **gigabytes** ou até mesmo **terabytes** por dia ²

TRANSFERÊNCIA DE 1TB

Rede (máx.)	Duração
10 Mbps	300h (12,5 dias)
100 Mbps	30h
1 Gbps	3h
10 Gbps	20min

(valores aproximados baseados na vazão máxima teórica de cada rede)

1. Fonte: ESnet Fasterdata - <http://fasterdata.es.net/fasterdata-home/requirements-and-expectations/>

2. Fonte: RNP Projeto Science DMZ (WRNP 2014): <http://indico.rnp.br/getFile.py/access?contribId=31&resId=0&materialId=slides&confId=188>

Desempenho das redes: expectativas (2)

Vazão de Dados no Tempo

- Quanto tempo demora para transferir N bytes (em bits/s)

Dados	1min	5min	20min	1h	8h	24h	7d	30d
1XB					277,78Tbps	92,59Tbps	13,23Tbps	3,09Tbps
100PB					27,78Tbps	9,26Tbps	1,32Tbps	308,64Gbps
10PB	1.333,33Tbps	266,67Tbps	66,67Tbps	22,22Tbps	2,78Tbps	925,93Gbps	132,28Gbps	30,86Gbps
1PB	133,33Tbps	26,67Tbps	6,67Tbps	2,22Tbps	277,78Gbps	92,59Gbps	13,23Gbps	3,09Gbps
100TB	13,33Tbps	2,67Tbps	666,67Gbps	222,22Gbps	27,78Gbps	9,26Gbps	1,32Gbps	308,64Mbps
10TB	1,33Tbps	266,67Gbps	66,67Gbps	22,22Gbps	2,78Gbps	925,93Mbps	132,28Mbps	30,86Mbps
1TB	133,33Gbps	26,67Gbps	6,67Gbps	2,22Gbps	277,78Mbps	92,59Mbps	13,23Mbps	3,09Mbps
100GB	13,33Gbps	2,67Gbps	666,67Mbps	222,22Mbps	27,78Mbps	9,26Mbps	1,32Mbps	0,31Mbps
10GB	1,33Gbps	266,67Mbps	66,67Mbps	22,22Mbps	2,78Mbps	0,93Mbps	0,13Mbps	0,03Mbps
1GB	133,33Mbps	26,67Mbps	6,67Mbps	2,22Mbps				
100MB	13,33Mbps	2,67Mbps	0,67Mbps	0,22Mbps				

LEGENDA

Requer vazão abaixo de 100Mbps

Requer vazão entre de 100Mbps e 10Gbps

Requer vazão entre de 10Gbps e 100Gbps

Requer vazão acima de 100Gbps

Nota: Kilo, Mega etc. estão em unidades SI (ex.: 1KB = 1000 bytes)

Protocolo TCP: ubíquo e frágil

As redes fornecem conectividade entre *hosts*

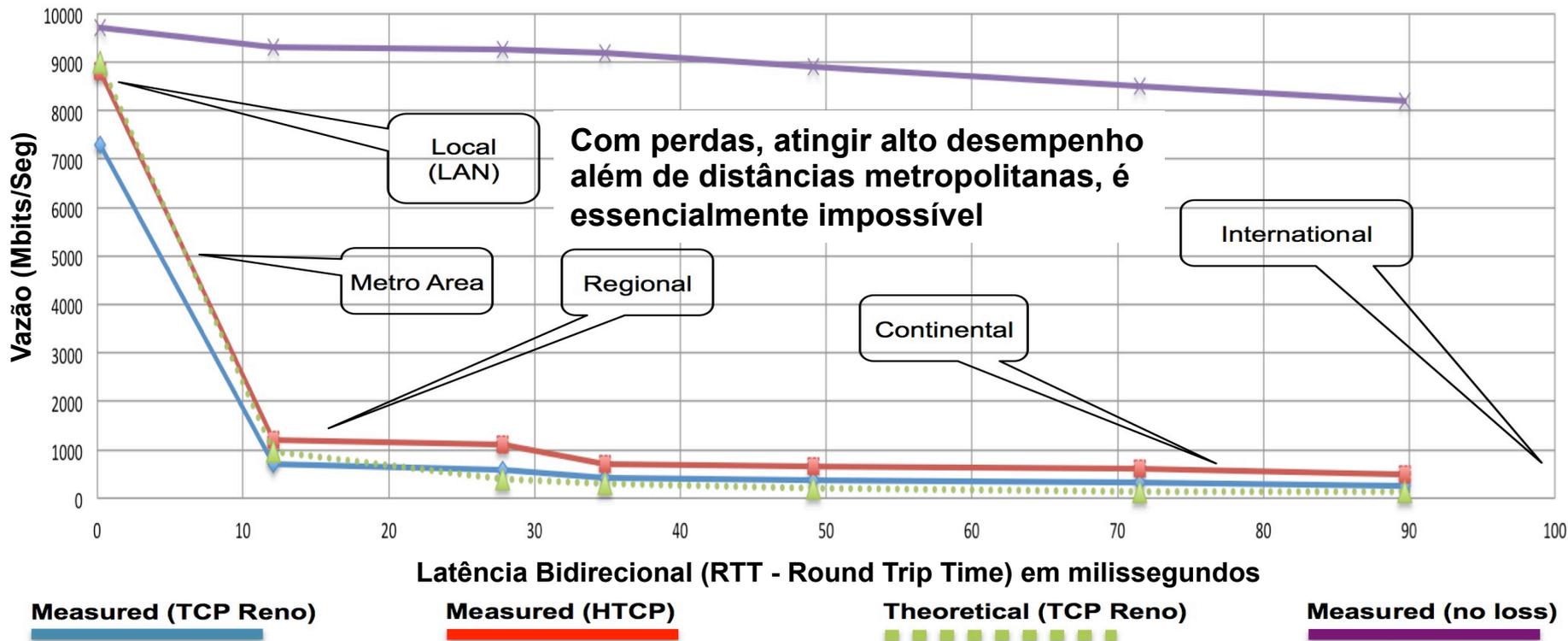
- Como aplicações e *hosts* enxergam a rede?
- Para uma aplicação, a interface para a outra extremidade é um *socket*
- A comunicação é entre aplicações, na maioria sobre TCP
- Protocolo TCP (*Transmission Control Protocol*)
 - Para o TCP, perdas de pacotes são interpretadas como congestionamentos e significam: “reduzir a taxa de transmissão”
 - Para o TCP, as perdas de pacotes, somadas à alta latência (alto RTT) causam enormes impactos no desempenho das redes
 - O TCP é o protocolo mais usado pela grande maioria das aplicações de transferências de dados (HTTP, FTP, SMTP etc)
 - Ex.: na RNP ~85% do tráfego é TCP; na Esnet, acima de 95%

Falhas completas x Falhas Parciais

- “Falhas completas” são o tipo de problema que todos entendem
 - Cortes de fibra
 - Falha de energia
 - Hardware que para de funcionar
- Sistemas de monitoramento são bons em alertar falhas
 - Ex.: NOC percebe um alerta vermelho
 - Engenheiros são acionados
- “Falhas parciais” são diferentes e, frequentemente passam despercebidas
- Conectividade básica funciona
- Desempenho é sofrível
- Quanto devemos nos importar com falhas parciais?

Desempenho do TCP: quanto mais distante, pior. Falhas parciais causam perdas de pacotes, que o afetam

Vazão com **latência crescente** e com **perdas de pacotes de 0,0046%**



9 – ESnet Science Engagement (engage@es.net) - 4/21/14

Lawrence Berkeley National Laboratory

U.S. Department of Energy | Office of Science

Onde estão os problemas?



Onde estão os problemas?

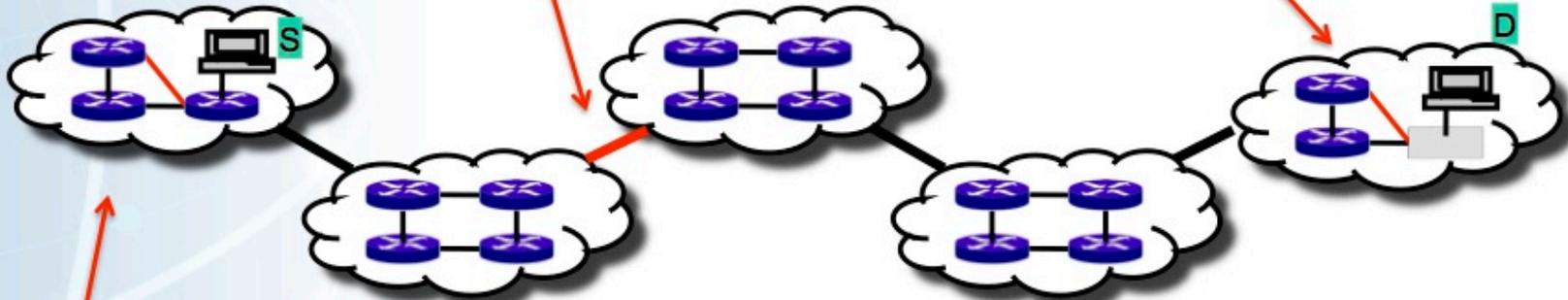
Links congestionados
ou com falhas entre
domínios

Problemas dependentes
da latência dentro de domínios
com baixo RTT

Campus
Origem

Backbone

Campus
Destino



Rede acadêmica

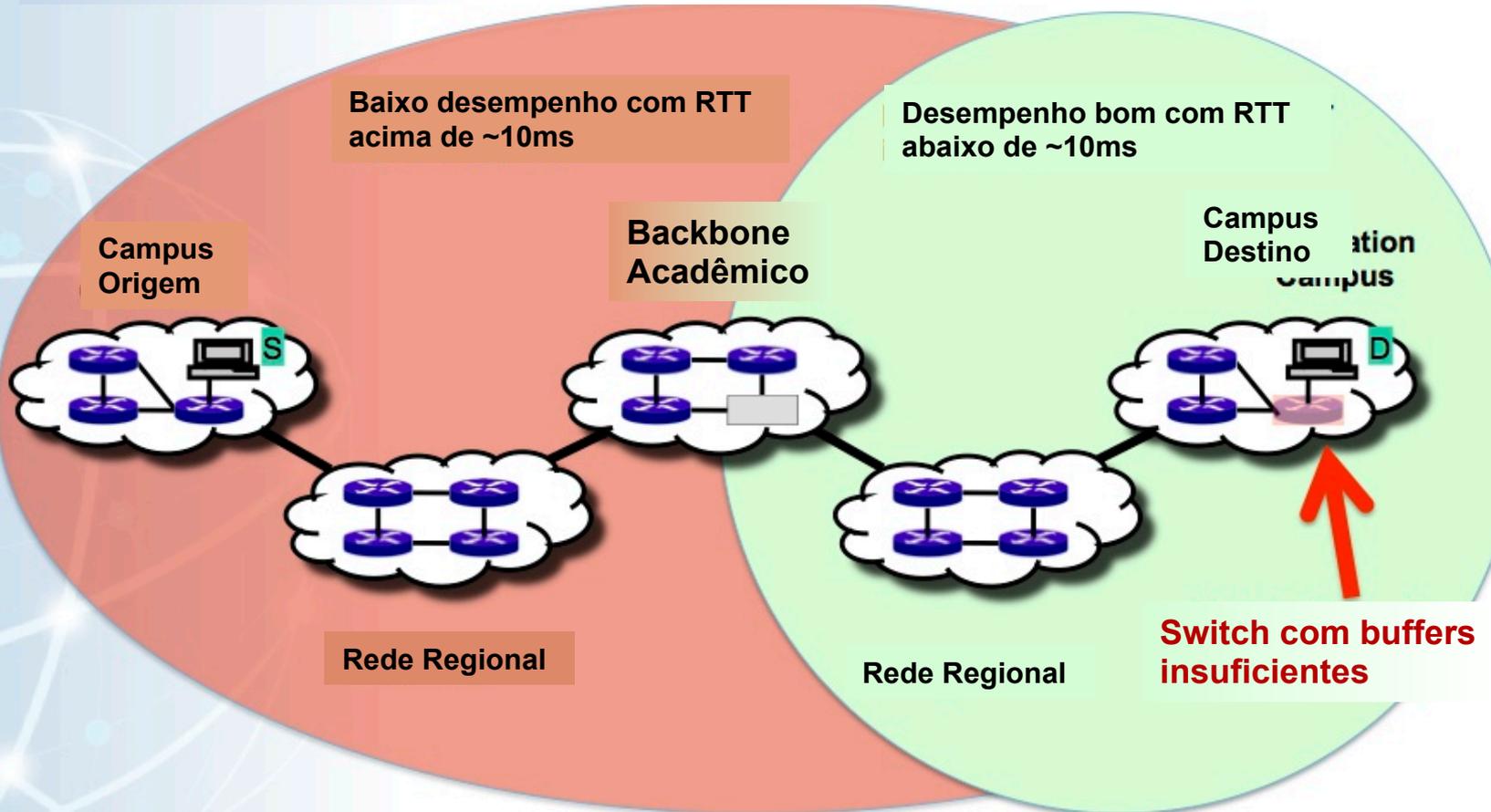
Rede metropolitana

Links intra campus
congestionados

Testes locais não encontram todas as causas



Testes Locais Não Encontram Todas as Causas



Causas de Perdas de Pacotes

- Rede Congestionada
 - Fácil de confirmar via SNMP, “simples” de consertar com \$\$
 - Isto não é uma ‘falha parcial’; apenas um prob. de capacidade
 - Com frequência, pessoas assumem que congestionamento é a causa, quando na verdade não é.
- Switches com *buffers* insuficientes descartando pacotes
 - Difícil confirmar
- Firewall com processamento insuficiente descartando pacotes
 - Difícil confirmar
- Conectores ou fibras sujas, níveis ópticos falhando
 - Às vezes fácil de confirmar verificando contadores de erros
- Host sobrecarregado ou com baixo processamento perdendo pacotes
 - Fácil de confirmar olhando a carga de CPU

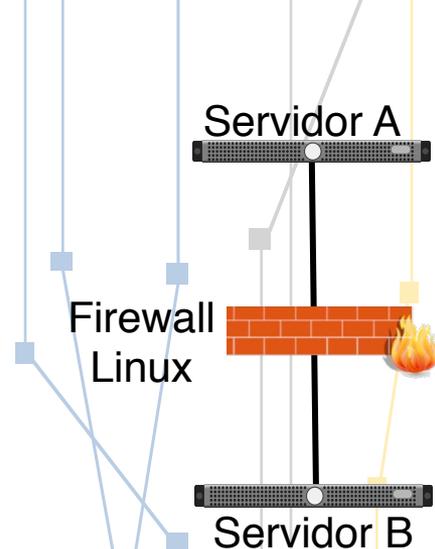
Falhas Parciais (*soft failures*)

- **Falhas totais** (*hard failures*) são fáceis de detectar e consertar
- **Falhas parciais** (*soft network failures*) são falhas onde a conectividade e comunicação básica funciona, mas não é possível atingir alto desempenho.
- O TCP foi intencionalmente projetado para esconder os erros transmissões do usuário:
 - “*Enquanto o TCP funcionar adequadamente e o sistema internet não sofrer interrupções, nenhum erro de transmissão afetará os usuários.*”
(IEN 129, RFC 716)
- Algumas falhas parciais afetam somente fluxos de alta vazão e fluxos de alta latência (RTT acima de 100ms).
- Falhas parciais podem durar por anos sem serem detectadas

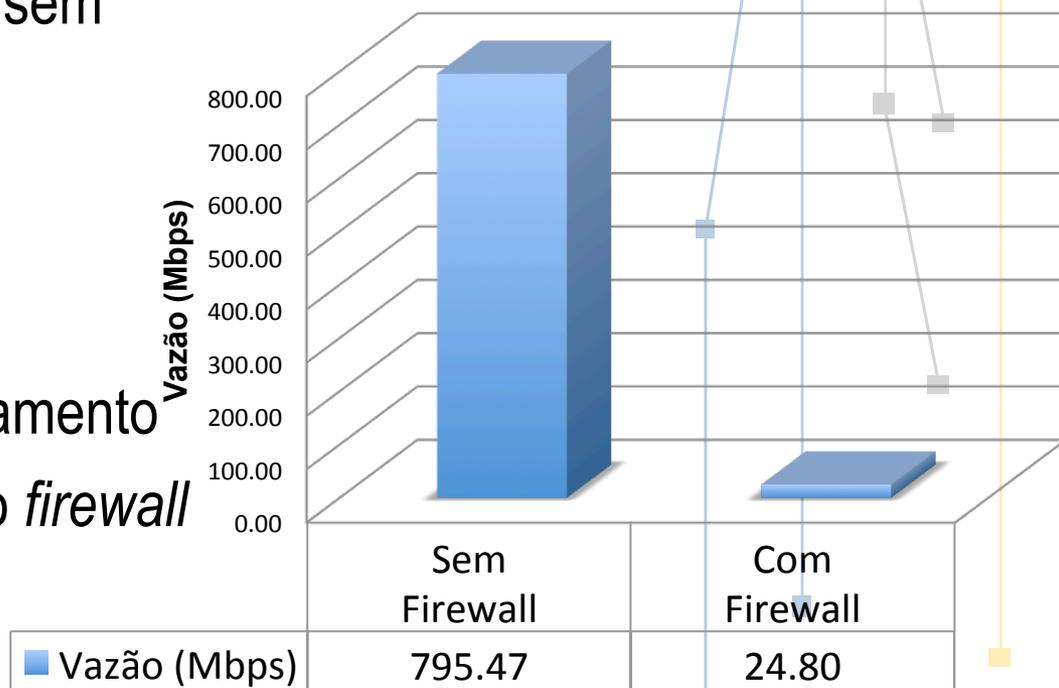
Falhas Parciais: Firewalls (1)

e-Ciência e transferências de dados

- **Segurança**
- Transferência de arquivos com e sem *firewall*
- Servidores físicos dedicados sem otimização
- Transferência disco-a-disco:
 - 1 tamanho de arquivo: 1GB
 - 1 ferramenta: xrootd
- Enlaces de 1Gbps, com roteamento
- Testes com ou sem regras no *firewall*
 - 6324 regras *iptables*



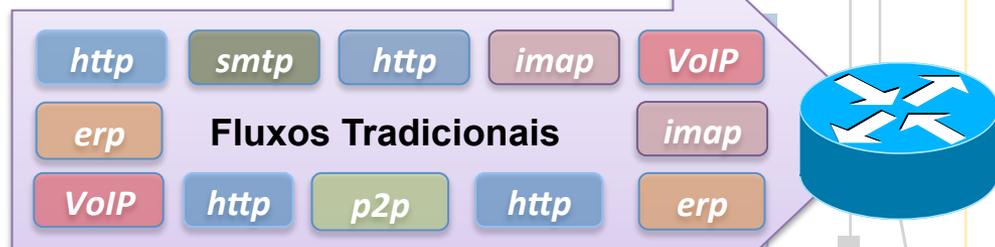
Vazão Xrootd - Arquivo 1G (Mbps)



Falhas Parciais: Firewalls (2)

- **Fluxos Tradicionais**

- Grande número de fluxos consumindo pouca banda
- Pequena taxa de perda de pacotes não afeta desempenho de forma significativa.
- Filtragem complexa



- **Fluxos Científicos**

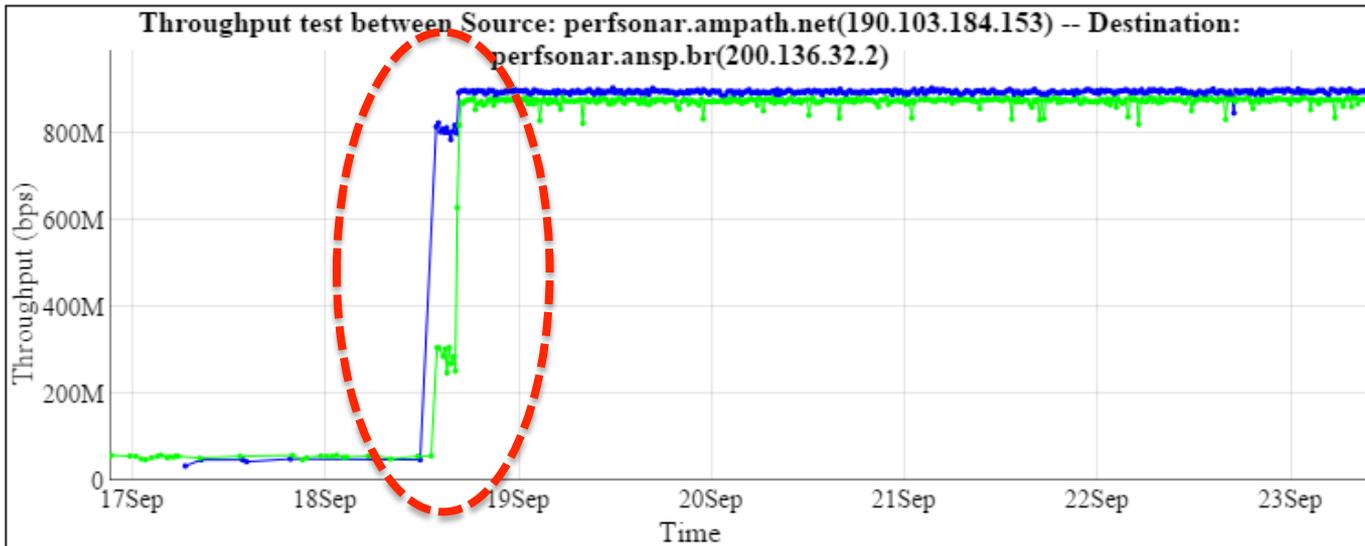
- Pequeno número de fluxos consumindo muita banda
- Pequena taxa de perda de pacotes afeta desempenho de forma significativa
- Controle simples



Falhas Parciais: Ajustes em Hosts (1)

Desempenho

- **Vazão:** ajustes nos hosts envolvidos em comunicações e transferências de dados são necessários em caminhos com alto RTT



Graph Key

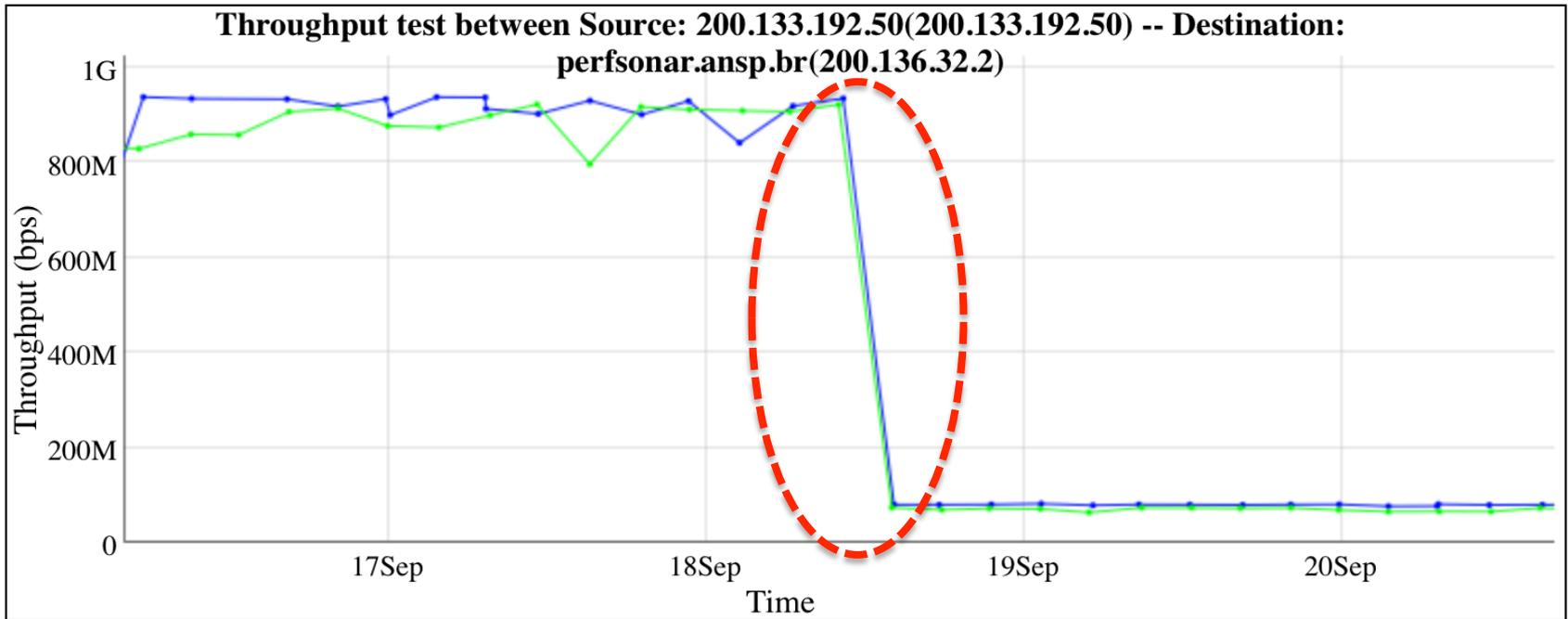
- Src-Dst throughput
- Dst-Src throughput

Direction	Max throughput(bps)	Mean throughput(bps)	Min throughput(bps)
Src > Dst	906.72M	884.68M	33.37M
Dst > Src	885.58M	814.93M	47.62M

Falhas Parciais: Ajustes em Hosts (2)

Desempenho

- Ajustes nos hosts podem ser necessários em caminhos **com baixo** RTT (~1ms)



Direction	Max throughput(bps)	Mean throughput(bps)	Min throughput(bps)
Src > Dst	940.22M	605.26M	75.66M
Dst > Src	938.07M	590.89M	63.7M

Métricas de Interesse

Use a ferramenta correta para o trabalho

- Para determinar a ferramenta correta, é preciso iniciar pelo que se deseja fazer...
- Que é importante medir?
 - **Perdas ou pacotes** fora de sequência, Duplicação (camada de transporte)
 - **Banda Alcançável** (p.ex.: vazão ou “*throughput*”)
 - **Latência unidirecional e bidirecional** (RTT e *One Way Delay*)
 - **Variação do atraso** (*jitter, delay variation*)
 - **Utilização da interface, descartes, erros** (camada de rede)
 - **Rotas dos fluxos** (*traceroute*)
 - **MTU** (*Maximum Transmission Unit*)

Monitoramento de Rede

- Todas as redes possuem algum monitoramento.
- Atende necessidades da equipe local para entender o estado da rede
- Essa informação poderia ser útil para usuários externos?
- Essas ferramentas podem funcionar em múltiplos domínios?

Além dos métodos passivos, há ferramentas para **medições ativas**.

- Frequentemente precisamos de um valor de vazão (*'throughput'*). (É possível automatizar essa idéia?)
- Seria bom ter estatísticas de desempenho (Por dia? Por semana? Por ano? De múltiplos pontos?)
- Onde está o *middleware* de medições? (Algo que permita a fácil troca de métricas coletadas localmente, em escala global?)

perfSONAR

- Atualmente mais de 1200 pontos de medições



<http://stats.es.net/ServicesDirectory/>

Projeto perfSONAR

- Projeto iniciado há 10 anos nos EUA, c/ participação da GÉANT (Europa) e RNP (Brasil)
- **Motivação**
 - As redes são parte essencial das e-Ciências
 - O desempenho é fator crítico
 - Dificuldades no uso efetivo das redes WAN por cientistas
- **O perfSONAR possui:**
 - API padrão aberto
 - Camada de Web Services para comunicação entre MPs
- **O perfSONAR é uma ferramenta para:**
 - Normalizar ou compatibilizar as

expectativas sobre o desempenho das redes

- Encontrar problemas (“*soft failures*”)
- Ajudar a consertar os problemas em múltiplos domínios de rede
- Os problemas são mais difíceis quando múltiplas redes são envolvidas
- perfSONAR oferece um padrão para publicação e intercâmbio de dados de monitoramentos ativos e passivos
- Esses dados são de interesse para pesquisadores e operadores de redes

Serviço MonIPÊ

Serviço de medições

de desempenho fim-a-fim entre a RNP, seus clientes, e com outras redes

Escopo e Objetivos (2015)

- Atender melhor instituições clientes da RNP;
- Estender a cobertura do monitoramento até a rede da instituição cliente;
- Colaborações em projetos de e-Ciência — p.ex: física, astronomia — e outras comunidades, colaborando com outras redes (p.ex.: ESnet, GÉANT, Internet2, NTT e RedCLARA);
- Usuários gestores de TI e comunidades de pesquisa;
- Equipamentos mais baratos;
- Simplificar a instalação;

- Melhorar a usabilidade;
- Ambiente de verificação e acompanhamento do desempenho

MEDIÇÕES

- sob demanda; temporárias e periódicas;
- de alta precisão do desempenho da rede;
- da última milha;
- **Backbone:** medições entre PoPs
- **PoPs:** medições entre PoPs e clientes diretamente conectados ao PoP
- **Internacional:** testes a outras redes acadêmicas e instituições (ex.: CERN)

MÉTRICAS

- Perdas de pacotes,
- Atraso {bi,uni}direcional
- Vazão (banda alcançável) em TCP e UDP

MonIPÊ: Componentes

Interface (GUI)

- Portal de Medições

Infraestrutura

- Pontos de Medição (*Measurement Points - MP*)
 - VMs e Kits Baixo Custo: até 1Gbps
 - Servidores dedicados: até 10Gbps

Virtualização

• MPs de baixo custo

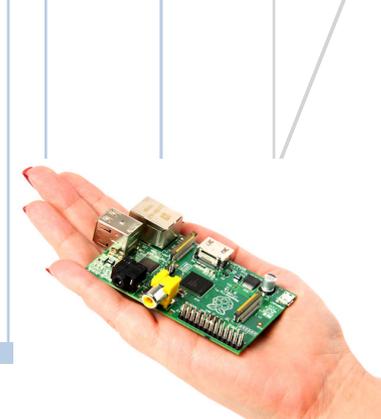
(kit: 1 Mini PC ou 2 SBCs¹ + GPS Adafruit)

- 1a G. **Raspberry Pi e CuBox**
CPU ARM, RAM 512MB, NIC 1GbE
- 2a G. Mini PC **Blue Appliance 847**
Intel Dual Core 847, 2G RAM, (2x) NICs 1GbE

▪ MP 10G

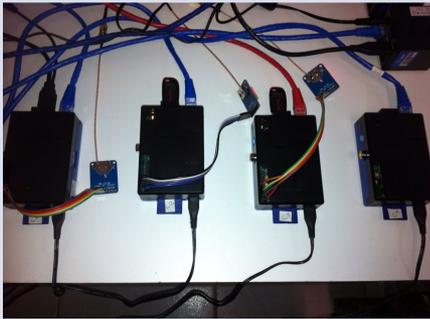
- 1a. G. **Dell R620 - CPU:** Intel Xeon 2GHz
RAM: 16GB, **HDD:** 2x 500GB (raid 1), **NIC:**
2x 10GbE + 2x 1GbE (BCM57800)

1. Single Board Computer



Kits de Baixo Custo: 1ª e 2ª geração

Hardware de 1a Geração (2013) - Custo total aproximado ~R\$ 2.000,00



Raspberry Pi
Testes latência unidirecional



Adafruit GPS + antena
Sincronização dos relógios



CuBox Pro
Testes de vazão c/ TCP e UDP

Hardware de 2a Geração (2014) - Custo total aproximado ~R\$ 1.200,00



Blue Appliance 847
Testes vazão e latência

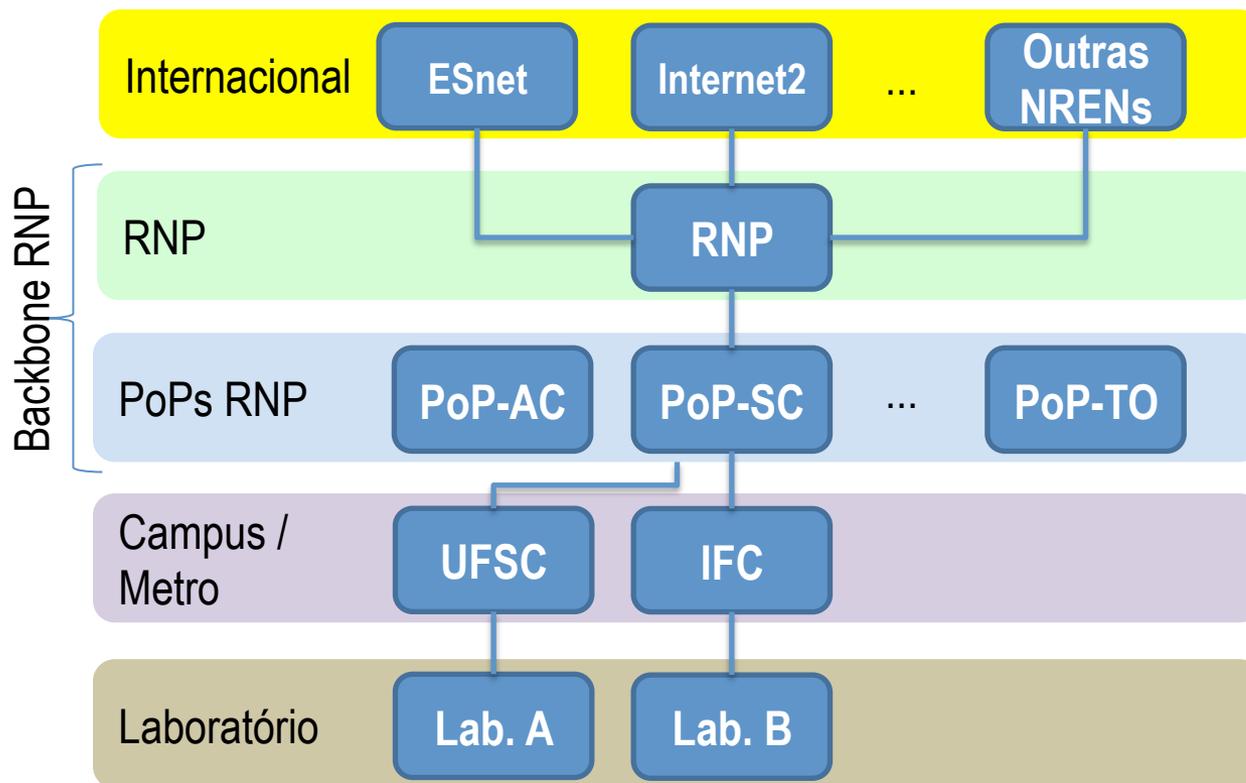


Adafruit GPS + antena
Sincronização dos relógios

Redução de 40% no custo total entre os kits da 1ª e 2ª versão

1. Raspberry Pi <http://www.raspberrypi.org/>
2. CuBox <http://http://cubox-i.com/>
3. Componentes MonIPÊ <http://goo.gl/rNEFWO>

MonIPÊ: arquitetura lógica em domínios

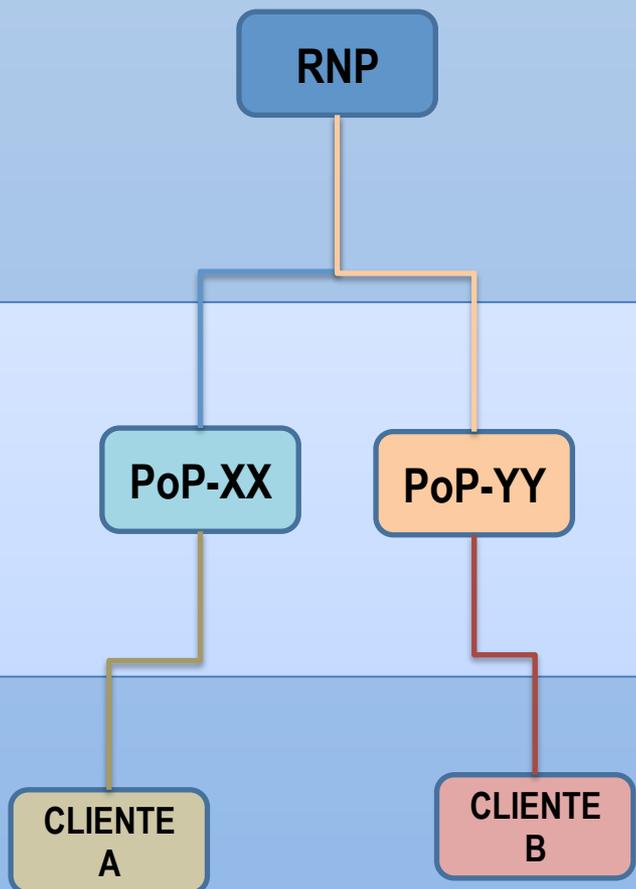


Domínio	Medições (em ou entre)
RNP	PoPs
PoP	PoP e instituição cliente diretamente conectada
Metro	PoP e Rede Metro
Campus	Campus e laboratórios de pesquisa ou outros campi
Laboratório	Dentro do campus
Internacional	Backbone RNP e outras NRENs

MonIPÊ: Componentes da Arquitetura

HIERARQUIA

COMPONENTES



Domínio RNP

- Portal de Medições
- Módulo Agendamento
- Módulo Sob-demanda
- **perfSONAR PHP-MA-SQL**
- **perfSONAR Lookup Service**
- **perfSONAR CLMP**

Domínio PoP

- Portal de Medições
- Módulo Agendamento
- Módulo Sob-demanda
- **Antena GPS**
- **perfSONAR PHP-MA-SQL**
- **perfSONAR CLMP**

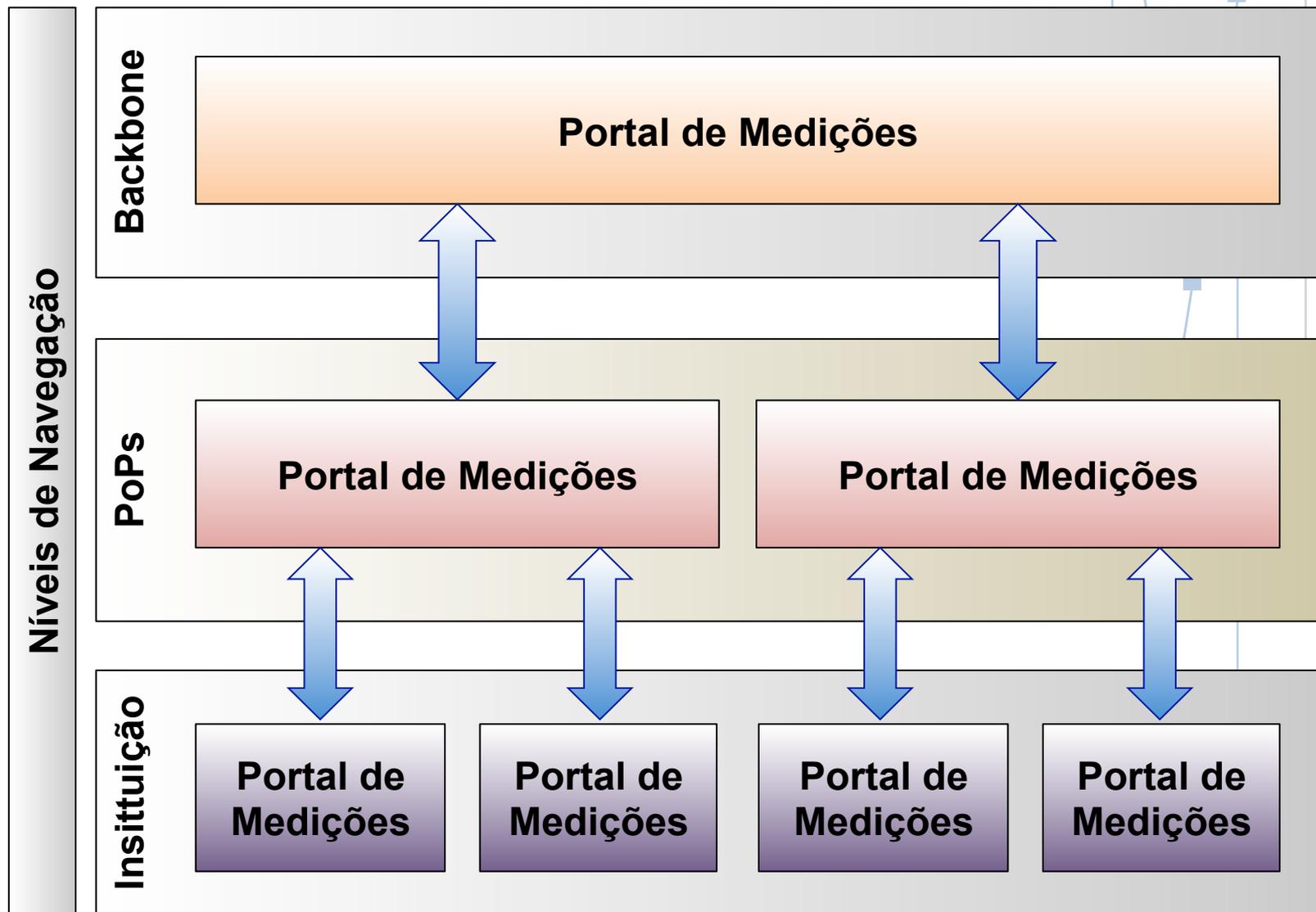
Domínio Instituição

- Portal de Medições
- Módulo Sob-demanda
- **Antena GPS de baixo custo**
- **perfSONAR CLMP**

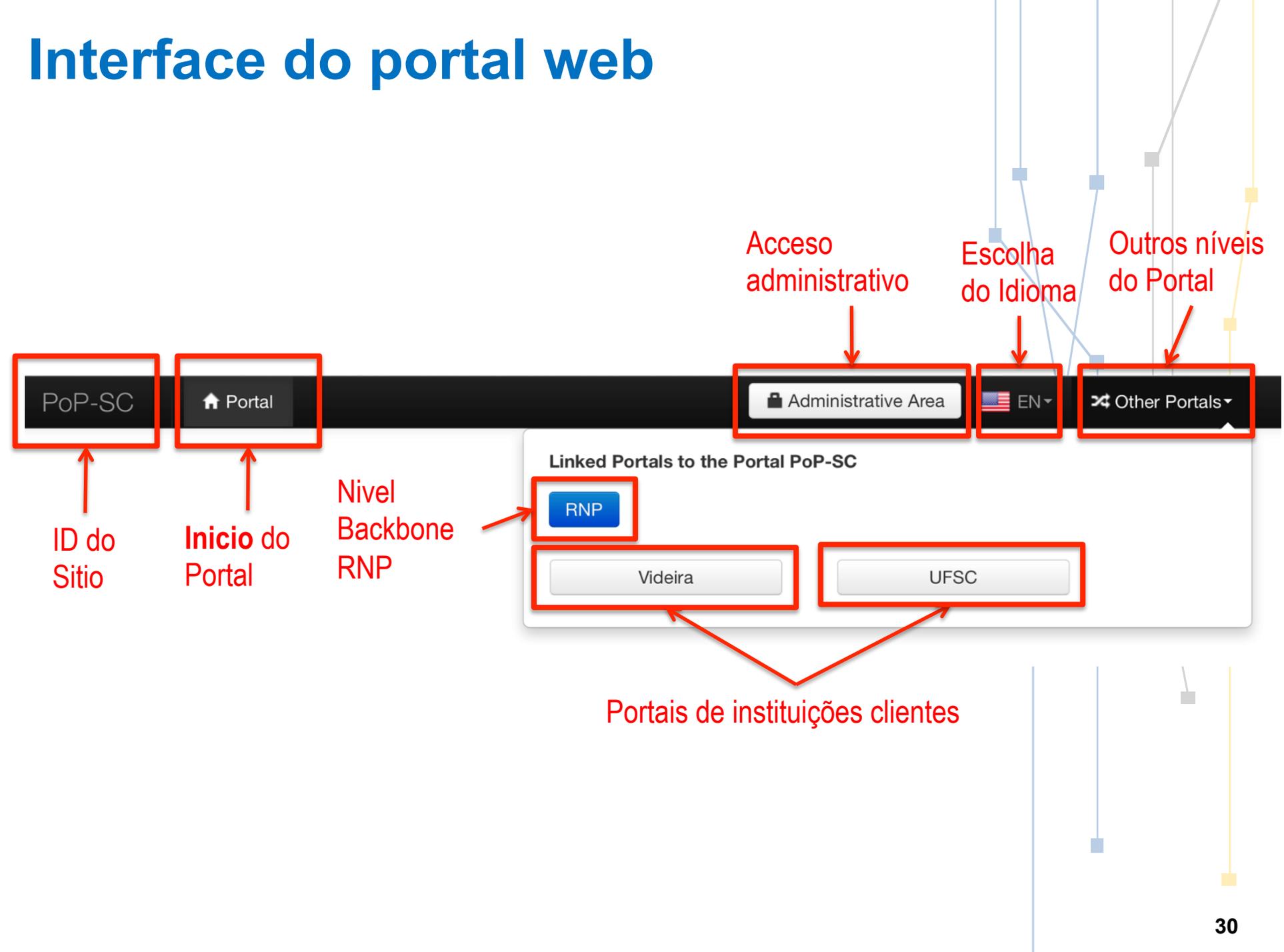
MonIPÊ: Portal

- **Interface gráfica web**
 - Configurações, gerenciamento, navegação multinível (domínios e *hosts*) e *feedback* de usuários (falhas/*bugs*, sugestões, críticas)
- **Agendamentos de testes:** sob demanda, temporários e periódicos
- **Armazenamento e recuperação** de medições
- **Visões:** pública e acesso restrito por senha (em desenv.: federado)
- **Integração** ao projeto perfSONAR (*Global Lookup Service - GLS*)
 - **Ferramentas:** ping, owamp, bwctl, traceroute* e ndt*
- **Testes** para MPs registrados no GLS
- **Atualizações:** manual ou automática

Portal web: navegação hierárquica entre os vários pontos de medições



Interface do portal web



Interface do portal web: teste de vazão

The screenshot displays the perfSONAR web portal interface for a bandwidth test. The browser address bar shows the URL `200.237.196.145/painel/testere`. The page header includes the MONIPE logo, navigation links for Portal, Painel de Monitoração, and Configurações, and user information for admin in Brazil (BR).

TESTE SUA REDE

- Vazão (selected)
- Atraso Bidirecional
- Atraso Unidirecional

AGENDAMENTO

- Ponto-a-Ponto
- Ponto-a-Multiponto

VISUALIZAR

- Meus Testes
- Testes Agendados

Teste de Vazão

Modelo de Teste: Vazão TCP - Vazão TCP Padrão

Host A: 200.237.196.145

Sentido: [Direção]

Host B: PoP-MG BWCTL Server - tc...

Intervalo: 1

Duração: 10

Executar Teste

Resposta

Gráfico | Tabela | Texto Simples | Requisição perfSONAR | Resposta perfSONAR

Teste de Vazão - 200.237.196.145 -> 200.131.0.165

Intervalo (s)	Enviado (Mbps)	Recebido (Mbps)
1.0	~45	~45
2.0	~55	~55
3.0	~55	~55
4.0	~55	~55
5.0	53.48	53.19
6.0	~55	~55
7.0	~55	~55
8.0	~55	~55
9.0	~55	~55
10.0	~55	~55

Salvar Resultados Evento: Nenhum evento associado. Tornar resultado público

Versões dos componentes:
CORE: 1.0.59 | WEB: 1.0.19
CLMP: 1.0.29 | NMWG: 1.0.5 | LIBS: 1.0.3

MONIPE Serviço de monitoramento da rede IqE | PoP-SC | RNP | perfSONAR powered | [Enviar feedback](#)

Display a menu

Interface do portal web: teste de latência

Portal perfSONAR

MONIPE Portal Painel de Monitoração Configurações admin BR Outros Portais

TESTE SUA REDE
Vazão
Atraso Bidirecional
Atraso Unidirecional

AGENDAMENTO
Ponto-a-Ponto
Ponto-a-Multiponto

VISUALIZAR
Meus Testes
Testes Agendados

Teste de Atraso Bidirecional

Modelo de Teste: Atraso Bidirecional - Atraso Bidirecional Padrão

Host A: 200.237.196.145

Host B: PoP-MG PING Server - 200.1...

Intervalo: Contagem: 10

Executar Teste

Resposta

Gráfico Tabela Texto Simples Requisição perfSONAR Resposta perfSONAR

Teste de Atraso Bidirecional - 200.237.196.145 -> 200.131.0.164

Iteração	RTT (ms)
1	22.05
2	22.05
3	22.05
4	22.05
5	22.05
6	22.05
7	21.9
8	22.05
9	21.9
10	21.9

Salvar Resultados Evento: Nenhum evento associado. Tornar resultado público

Versões dos componentes:
CORE: 1.0.59 | WEB: 1.0.19
CLMP: 1.0.29 | NMWG: 1.0.5 | LIBS: 1.0.3

MONIPÊ Serviço de monitoramento da rede Ipt

POP-SC

RNP

perfSONAR powered

Enviar feedback

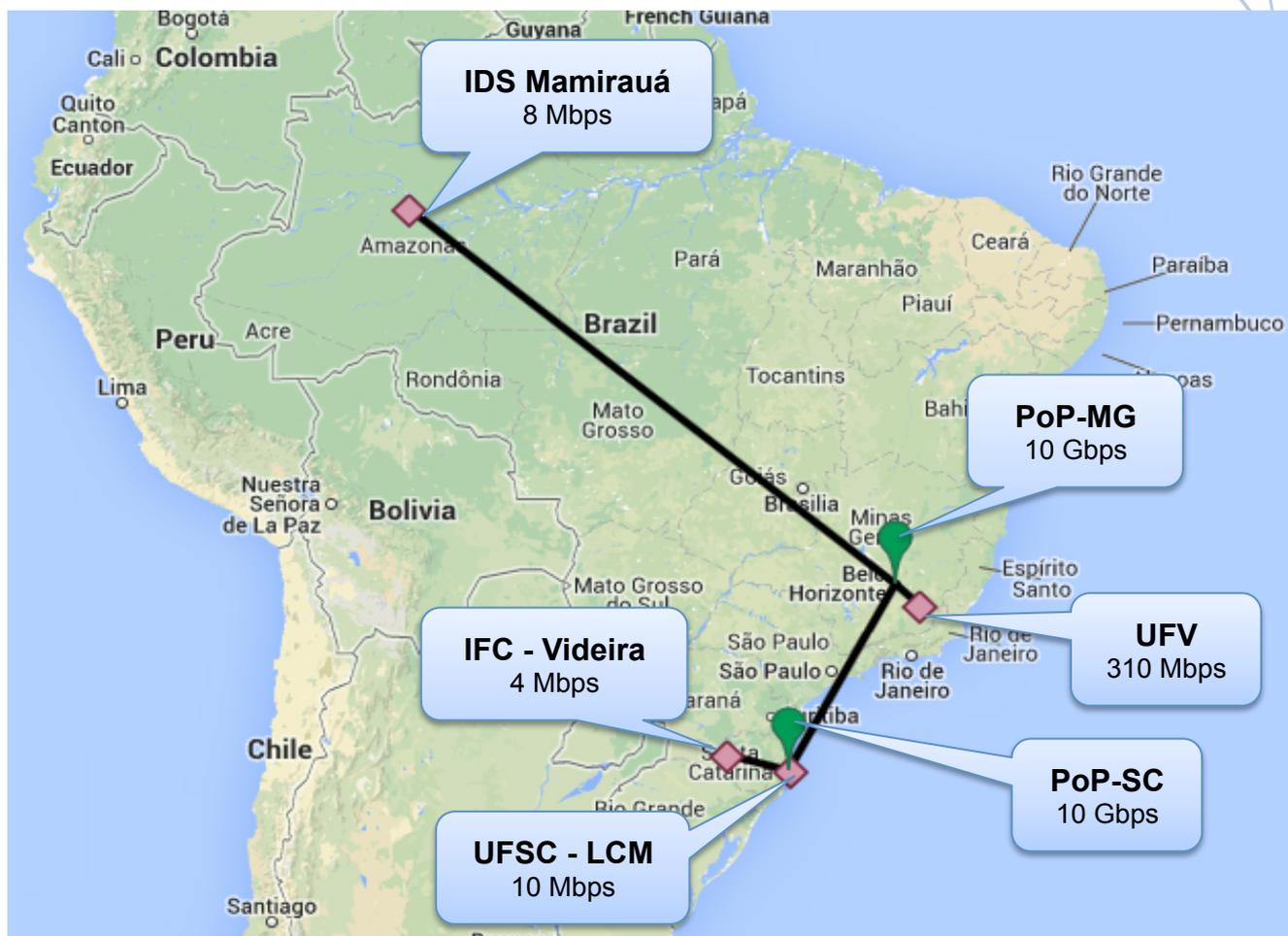
Display a menu

MonIPÊ: Realizações

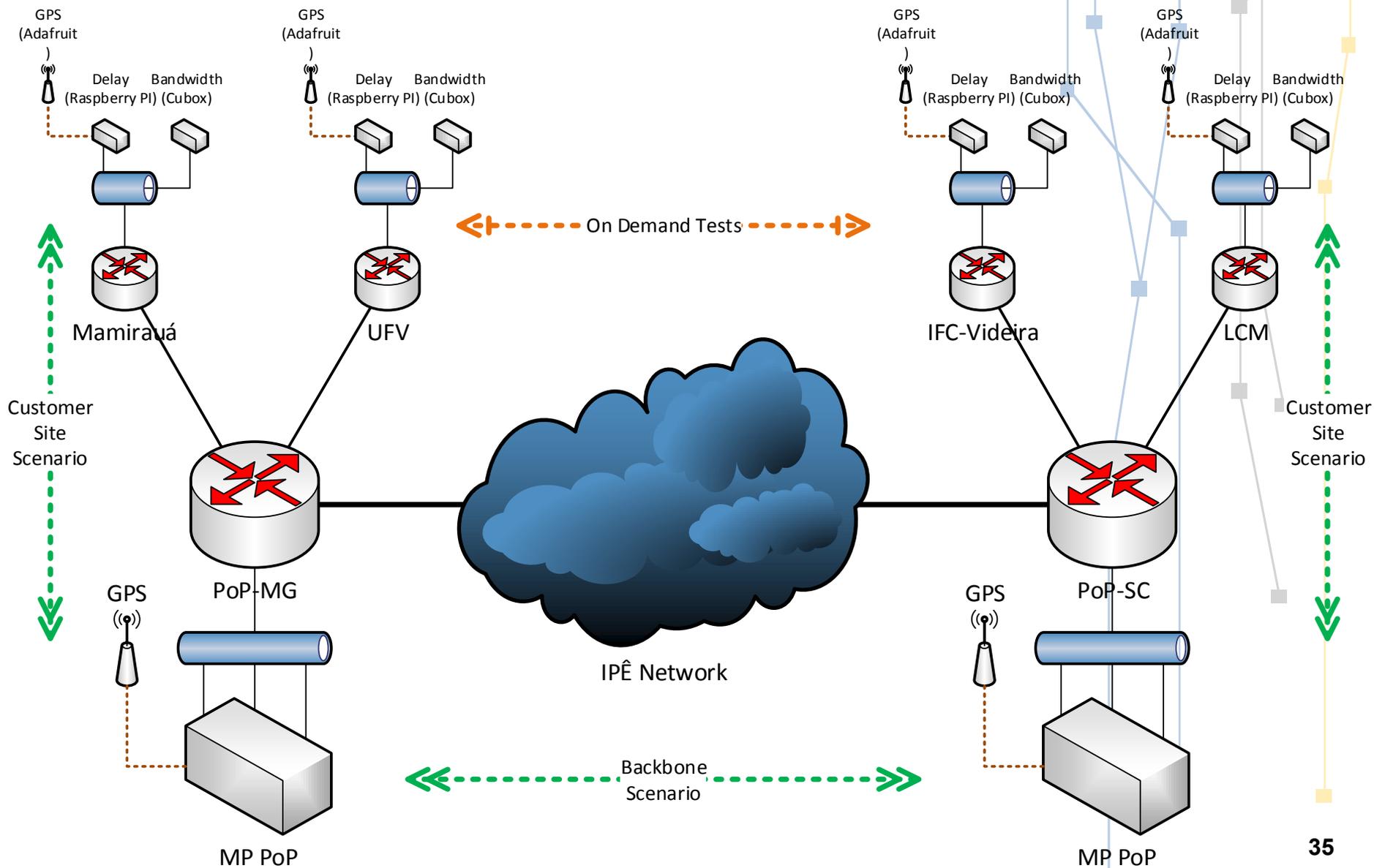
- Protótipos construídos e testados:
 - MPs para medições até 1Gbps e até 10Gbps (2013-2014)
- Piloto realizado com 4 instituições (2013)
- Serviço em implantação no backbone RNP
- Prevista instalação em:
 - MPs VMs (em todos os PoPs)
 - MPs 10G (em 11 PoPs)
 - Kits baixo custo em ~100 instituições clientes da RNP
- Transição do Serviço Experimental MonIPÊ para produção, atendendo clientes da RNP
- Interoperabilidade com perfSONAR de outras redes

Piloto MonIPÊ (Nov/2013)

- **PoPs:** SC e MG (●)
- **Clientes:** Mamirauá, UFG, IFC-Videira e LCM/UFSC (◆)



MonIPÊ: infraestrutura do Piloto (Nov/2013)



Evolução e Futuro: apoio do CT-Mon

CT-Mon: Comitê Técnico em Monitoramento de Redes

EQUIPE

- Coordenador Artur Ziviani (LNCC)
- Secretário: Alex Moura (RNP)
- Membros da comunidade acadêmica
- Acompanha a evolução do perfSONAR e em medições
- Apoia a RNP na evolução do Serviço MonIPÊ
- Colabora com o esforço de padronização do perfSONAR e em nível nacional e internacional
- Reuniões presenciais e remotas
- Apresentações de novos desenvolvimentos
- Tarefas específicas atribuídas a membros do comitê e/ou grupos de estudo/testes
- Chamadas temáticas de P&D de curta duração

TEMAS

Soluções e tecnologias para monitoramento:

- Medições ativas e passivas: 1G, 10G, 100Gbps+
- Redes sem fio
- Metodologias, técnicas e ferramentas
- Grandes eventos: Copa do Mundo, Olimpíadas
- Transferências de grandes volumes de dados
- Monitoração nas camadas abaixo da camada 3
- Geração de alertas (previsão de problemas)
- Armazenamento, compactação recuperação e compartilhamento de dados históricos
- Evolução do tráfego na RNP, QoE, DPI etc...

Equipe



COORDENAÇÃO

RNP - Diretoria de Pesquisa e Desenvolvimento

Direção	Michael Stanton
Direção Adjunta	Iara Machado
Gerência	Alex Soares de Moura
Coordenação	Marcos Schwarz
Coordenação	Fausto Vetter

Desenvolvimento

PoP-SC - Ponto de Presença de Santa Catarina (UFSC)

Coordenação

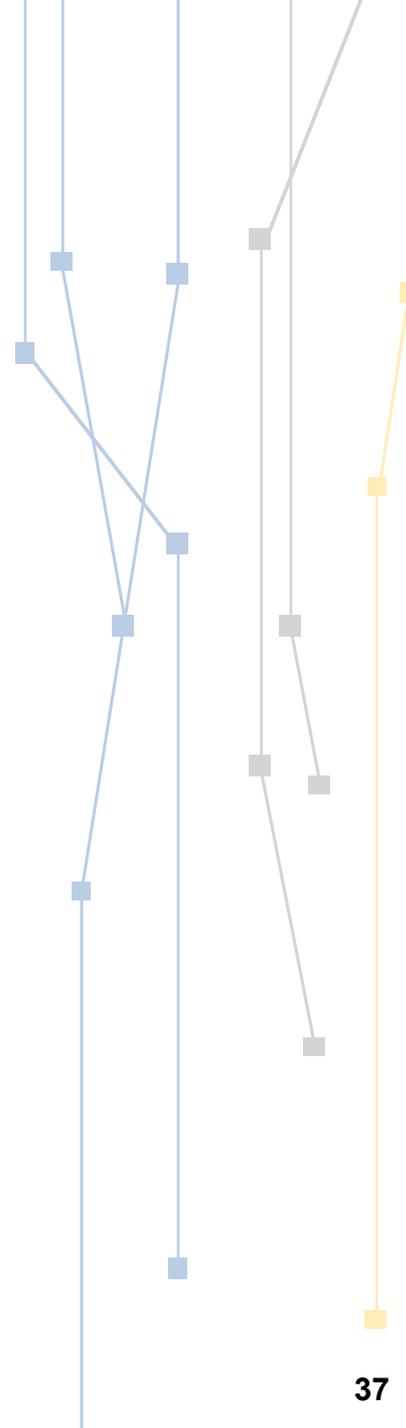
Coordenação Administrativa do Desenvolvimento	Edison Melo
Coordenação Técnica do Desenvolvimento	Murilo Vetter

Infraestrutura - Desenvolvimento de Hardware

Desenvolvimento de Hardware e Infraestrutura	Rodrigo Pescador
Desenvolvimento de Hardware e Infraestrutura	Guilherme Rhoden
Desenvolvimento de Hardware e Infraestrutura	Rodrigo Gonçalves
Desenvolvimento de Software Web	Paulo Brandtner
Desenvolvimento de Software Web	Luis Fernando Cordeiro

BOLSISTAS

Desenvolvimento de Software	Leonardo Schlüter Leite
Desenvolvimento de Software	Kádio Francisco Miguel Colzani



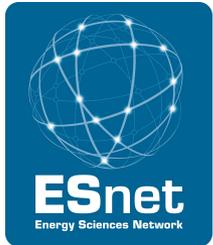
Agradecimentos



PoP-SC



perfSONAR



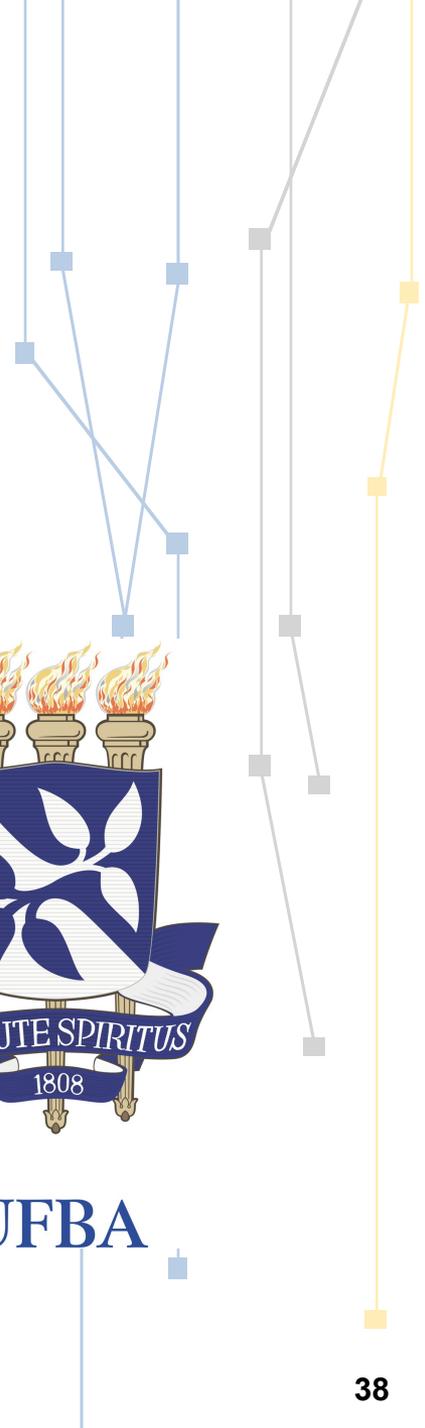
INTERNET²

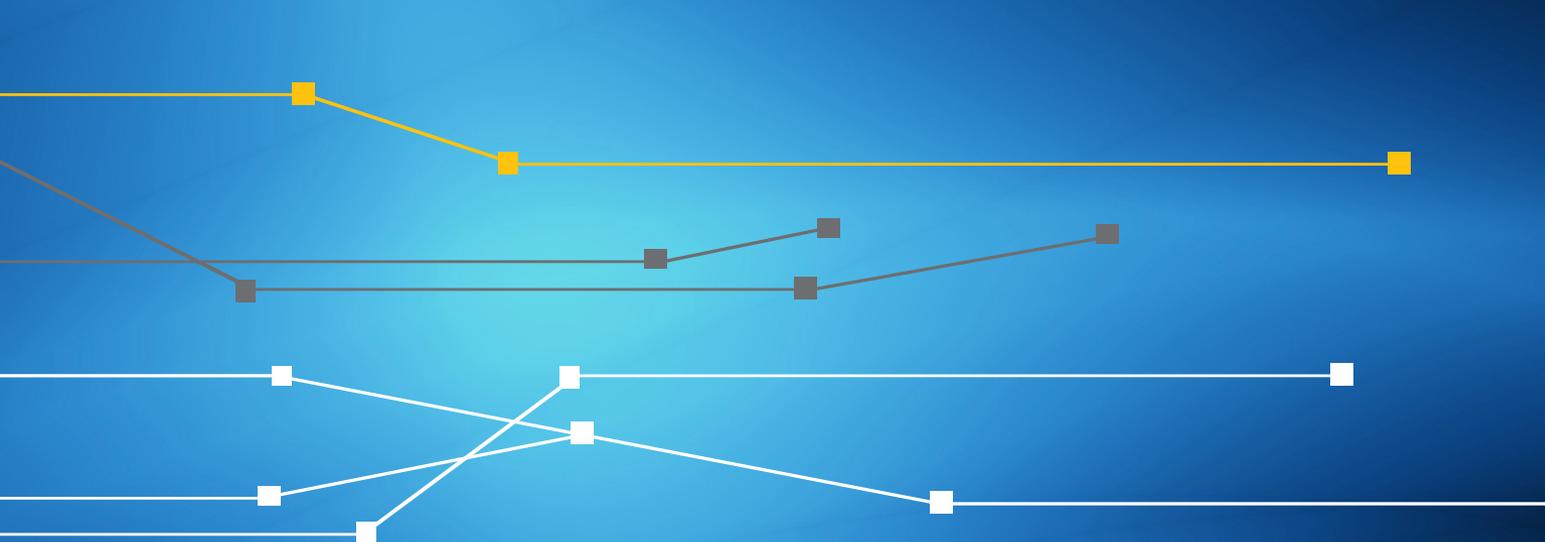


UFSC



UFBA





Obrigado!

Alex Soares de Moura
alex.moura@rnp.br

Rede Nacional de Ensino e Pesquisa – RNP
Ponto de Presença da RNP em Santa Catarina - PoP-SC
Universidade Federal de Santa Catarina – UFSC
Superintendência de Governança Eletrônica e Tecnologia da
Informação e Comunicação - SeTIC/UFSC



Ministério da
Cultura

Ministério da
Saúde

Ministério da
Educação

Ministério da
**Ciência, Tecnologia
e Inovação**