

Proposta de Serviço Piloto

Grupo de Trabalho – Segunda Fase

GT-CNC – Computação em Nuvem para Ciência: Armazenamento de Dados

Roberto Samarone dos Santos Araújo - UFPA

05/09/2012

1. Concepção do Serviço

1.1 Resumo

Serviços de armazenamento de dados em nuvem estão cada vez mais comuns. Estes serviços possibilitam a seus usuários o armazenamento de dados em um ambiente remoto e altamente disponível. Em sua primeira fase, o GT-CNC avaliou ferramentas para construção de nuvens de

armazenamento e estabeleceu um protótipo baseado em uma das ferramentas avaliadas. O protótipo é semelhante a serviços de nuvens conhecidos como o Dropbox e assim realiza funções como a criação de pastas, o envio de arquivos para nuvem, etc. Para sua segunda fase, o GT propõe-se a refinar este protótipo para a realização de um piloto de serviço de armazenamento em nuvem. Este refinamento inclui a adição de autenticação através do serviço de autenticação CAFe e o acesso a arquivos via Web. Além disso, o GT propõe-se a avaliar a adoção e a integração da nuvem com o projeto Cloud Drive do TERENA.

1.2 Abstract

Cloud storage services are very common nowadays. These services allow users to store their data in a remote and high available environment. In its first phase, GT-CNC evaluated tools for making a cloud storage service. Based on one of these tools, GT-CNC established a cloud storage prototype. The prototype is similar to known cloud storage services as Dropbox and thus it can perform the creation of folders, the users can send files, etc. For the second phase, the GT aims at improving the prototype in order to establish a cloud storage service pilot. These improvements include authentication through the identification service CAFe and will allow users to access their files via Web browser. In addition, the GT will evaluate the adoption and the integration of the cloud storage solution with the TERENA Cloud Drive Project pilot.

1.3 Descrição do Serviço Proposto

A computação em nuvem vem se tornando cada vez mais atrativa devido a seus inúmeros benefícios. Isso tem motivado o surgimento de uma variedade de serviços que facilitam tarefas do cotidiano. Um dos serviços muito utilizados atualmente é o armazenamento de dados em nuvem. Ele possibilita o armazenamento de dados em um ambiente remoto e altamente disponível. Exemplos destes serviços são inúmeros como o Google Drive [4], o Dropbox [2] e o Microsoft SkyDrive [7].

Em sua primeira fase, o grupo de trabalho computação em nuvem para ciência (GT-CNC) realizou um amplo estudo sobre sistemas de armazenamento em nuvem e seus serviços. Neste estudo foram identificadas e avaliadas ferramentas de código aberto. Além disso, foi realizado um levantamento sobre ferramentas e opções privadas para armazenamento em nuvem. O estudo também objetivou o estabelecimento de um protótipo de nuvem de armazenamento de dados, seguindo o modelo de serviço IaaS (infraestrutura como serviço).

Baseado na avaliação destas ferramentas, o GT selecionou uma delas e estabeleceu um protótipo de serviço de armazenamento de dados em nuvem. O protótipo foi implantado através

de servidores dispostos na UFPA e na UFSC. Ele possibilitou a equipe do GT avaliar melhor as características da ferramenta bem como testá-la. O acesso a nuvem foi possível a partir de computadores *desktops* e de dispositivos móveis (e.g. celulares).

O protótipo possui características fundamentais para serviços de armazenamento em nuvem. Ele possibilita o armazenamento de arquivos grandes (e.g. 3GB), a criação de pastas, a sincronização de arquivos entre computadores e a nuvem, o agrupamento de usuários para compartilhamento de arquivos na nuvem, a edição de arquivos diretamente na nuvem, etc. Além disso, a ferramenta utilizada no protótipo permite, dentre outros benefícios, a extensão da capacidade da nuvem sobre demanda e a possibilidade de utilização em conjunto com nuvens de processamento de dados.

A partir das características apresentadas, o protótipo é semelhante a serviços de armazenamento em nuvem comerciais (e.g. o Dropbox). Esse tipo de serviço tem atraído um número cada vez maior de usuários. Isso é consequência principalmente da disponibilidade dos dados, que podem ser acessados de qualquer lugar. O protótipo, no entanto, vai além da disponibilização do serviço de armazenamento a usuários.

Todavia, serviços de armazenamento de dados em nuvem comerciais têm recebido críticas recentemente. Embora grande parte destes serviços atendam as necessidades dos usuários, muitos são questionados quanto a sua licença de uso. A licença do Google Drive [3], por exemplo, diz o seguinte: “*Quando você faz upload ou de algum modo envia conteúdo a nossos Serviços, você concede ao Google (e àqueles com quem trabalhamos) uma licença mundial para usar, hospedar, armazenar, reproduzir, modificar, criar obras derivadas.*”. Dessa forma, há preocupações sobre a forma de como as licenças para estes serviços são definidas. Além dos problemas com as licenças de uso, alguns serviços não tem se mostrado confiáveis. Há relatos de falhas de segurança [6] [1] e como os dados dos usuários não são criptografados, não há garantias de que eles dados não sejam acessados por terceiros.

Apesar das críticas, este tipo de serviço vem sendo muito empregado por professores e pesquisadores no Brasil. Sua utilização varia entre o armazenamento de dados comuns ao armazenamento de informações importantes como dados de pesquisas.

Um serviço de armazenamento de dados em nuvem criado especialmente para o meio acadêmico (i.e. instituições de ensino e pesquisa) traria vários benefícios: ele poderia eliminar o armazenamento de dados sensíveis em nuvens externas e reduziria os riscos inerentes a serviços de nuvens comerciais. Além disso, ele poderia ser integrado a outros serviços da RNP, como por exemplo, a utilização de provedores de identidade da CAFe para a autenticação dos clientes do serviço.

Adicionalmente, um serviço próprio que possibilitasse a integração a serviços pagos seria menos custoso para as universidades. Esse serviço permitiria, por exemplo, compartilhar dados em projetos acadêmicos (com interfaces combinadas pelos pesquisadores) ou até mesmo a inclusão do serviço como módulo em outras ferramentas (ou aplicativos) em uso.

A fim atender a esses requisitos, para sua segunda fase, o GT-CNC propõem-se a implantar um serviço piloto de nuvem de armazenamento de dados. O piloto objetiva disponibilizar aos usuários da RNP, principalmente professores e pesquisadores, um serviço de nuvem de armazenamento de dados.

O piloto será baseado no protótipo implantado na primeira fase do GT. Ele permitirá aos usuários da RNP o armazenamento de seus dados em uma nuvem. A partir do piloto será possível avaliar o seu nível de disponibilidade em larga escala bem como a sua confiabilidade e segurança. Através deste serviço, os usuários poderão realizar operações comuns a serviços de armazenamento de dados como a criação e a sincronização de pastas. O piloto também objetiva a integração da nuvem de armazenamento de dados a ser criada com o sistema de identificação federada CAFe.

Para a realização do piloto serão dispostos servidores na UFPA, UFSC e UFMG. Ele herdará características da ferramenta utilizada no protótipo da primeira fase do GT e assim possibilitará:

- 1- A extensão da capacidade da nuvem de acordo com a demanda;
- 2- O armazenamento de grande quantidade de dados;
- 3- A possibilidade de integração com nuvens de processamento de dados, ou seja, nuvens que possibilitem a execução de máquinas virtuais;
- 4- Um serviço de baixo custo com o uso de soluções de baixo custo, mas de alta qualidade;
- 5- O suporte ao protocolo S3 da Amazon para comunicação entre os clientes e a nuvem;

6 O armazenamento e a recuperação de arquivos via clientes (incluindo dispositivos móveis);

7 A possibilidade de desenvolvimento de novos clientes através de APIs.

O piloto proposto poderá vir a ser um dos serviços de armazenamento em nuvem a ser utilizado pelo projeto TERENA Cloud Drive [8]. Neste projeto está sendo construído um middleware (broker) de acesso a nuvens de armazenamento. O middleware faz uma interface entre o usuário e o serviço de nuvem de armazenamento. Ele possibilita a autenticação federada e o sigilo de dados. Dentro do piloto de serviço aqui apresentado pretende-se também avaliar sua adoção e integração com o projeto Cloud Drive do TERENA.

A seguir é apresentado um resumo dos objetivos da segunda fase deste GT:

1 O aprimoramento do protótipo (melhorando algumas de suas características) para utilizá-lo como piloto;

2 A integração da nuvem de armazenamento de dados com o sistema de identificação federada CAFé;

3 A avaliação e a implantação de ferramentas de monitoramento no piloto;

4 A implantação do piloto de armazenamento de dados em nuvem;

5 Avaliar a adoção e a integração do piloto com o projeto Cloud Drive do TERENA;

6 A realização de testes no piloto;

7 Avaliar a disponibilidade, segurança e estabilidade do piloto para uso em larga escala;

1.4 Identificação do Público Alvo

O serviço a ser implantado através do piloto visa o armazenamento de dados em nuvem. Seu público alvo será primariamente professores e pesquisadores de instituições de ensino e pesquisa. No entanto, dado as características do serviço, ele pode ser empregado também por qualquer usuário da RNP (incluindo seus funcionários) que necessitem de armazenamento.

2. Definição do Serviço Piloto

2.1 Arquitetura do Serviço Piloto

O piloto será composto por uma nuvem de armazenamento de dados baseado no protótipo da primeira fase deste GT. Ele terá servidores dispostos na UFPA, UFSC e UFMG. Estes servidores serão divididos em servidores *proxy* e servidores de armazenamento, conforme a Figura 1.

Os servidores *proxy* serão encarregados pela autenticação e pelo atendimento das requisições dos usuários. No piloto serão dispostos três servidores *proxy*, um em cada instituição. Estes servidores funcionarão com balanceamento de carga. Dessa forma, caso um deles não seja acessível em um dado momento, os outros poderão atender as requisições.

Diferentemente, os servidores de armazenamento serão responsáveis pelos dados dos usuários. No piloto serão dispostos 5 servidores de armazenamento: dois na UFSC, dois na UFPA e um na UFMG. Devido sua estrutura em *cluster*, caso um destes servidores venha a

falhar, os outros servidores poderão atender as solicitações. Ressalta-se que a nuvem pode ser expandida com a adição de mais servidores.

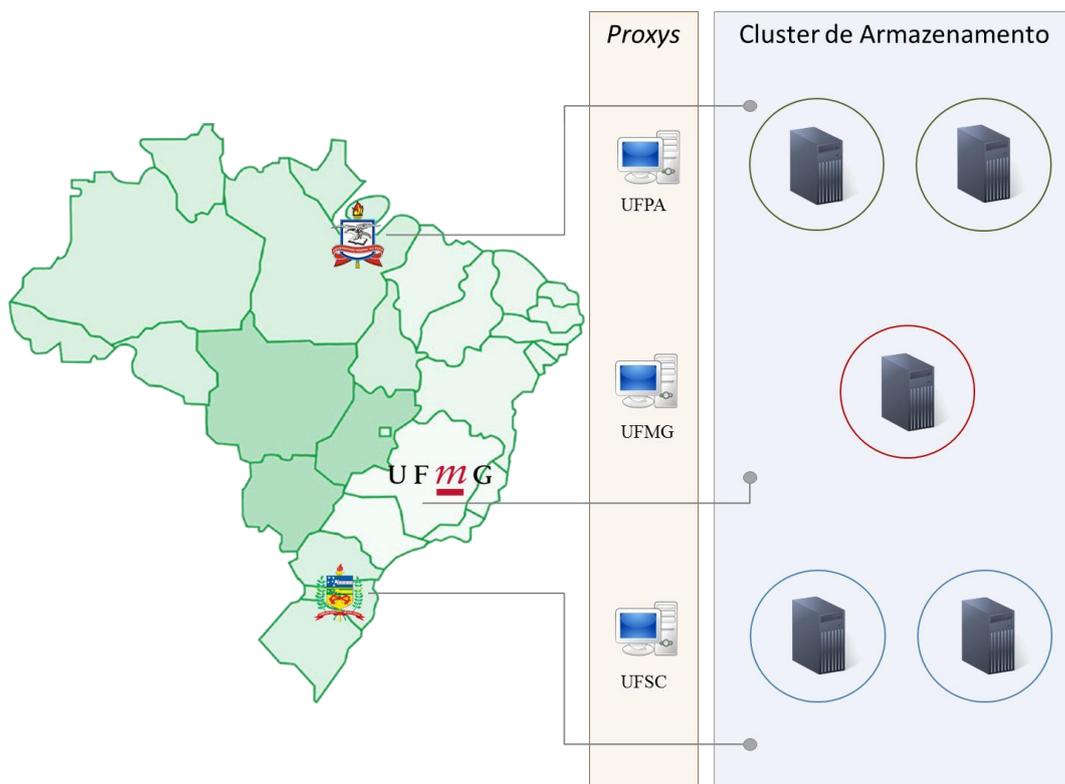


Figura 1. A arquitetura do piloto. Os usuários se comunicam com os servidores *proxy* através de programas clientes. Estes servidores encaminham as requisições recebidas aos servidores de armazenamento.

Na primeira fase do GT foram dispostos 1 servidor na UFPA (POP-PA RNP) e 1 servidor na UFSC (LabSEC). Na segunda fase, serão utilizados os servidores da fase 1. Além disso, serão adicionados um novo servidor na UFSC e outro na UFPA. Na UFMG serão instalados dois servidores, um para o *proxy* e outro para armazenamento de dados.

Visando uma nuvem de baixo custo, tanto os servidores *proxy* como os servidores de armazenamento utilizarão sistema operacional Linux. Além disso, estes servidores utilizarão o *software OpenStack Swift* [11] para formar a nuvem de armazenamento. De forma a gerenciar o ambiente de nuvem, o GT pretende utilizar ferramentas de monitoramento como descrito adiante.

De forma a possibilitar o acesso à nuvem pelos usuários, o piloto também envolverá programas clientes. Estes programas são instalados nos computadores dos usuários. O piloto utilizará o cliente *Cyberduck* [9] (Mac OS e Windows) e o cliente *Rackspace* [10] para dispositivos móveis baseados em IOS. O GT pretende desenvolver um cliente para dispositivos móveis baseados em Android. Além disso, serão realizados estudos para a inclusão de acesso aos arquivos da nuvem via navegadores *Web*. A Figura 2 ilustra de uma pasta no piloto através de um cliente *Desktop*. A Figura 3 ilustra a listagem de pastas e arquivos no piloto por meio de um cliente IOS.

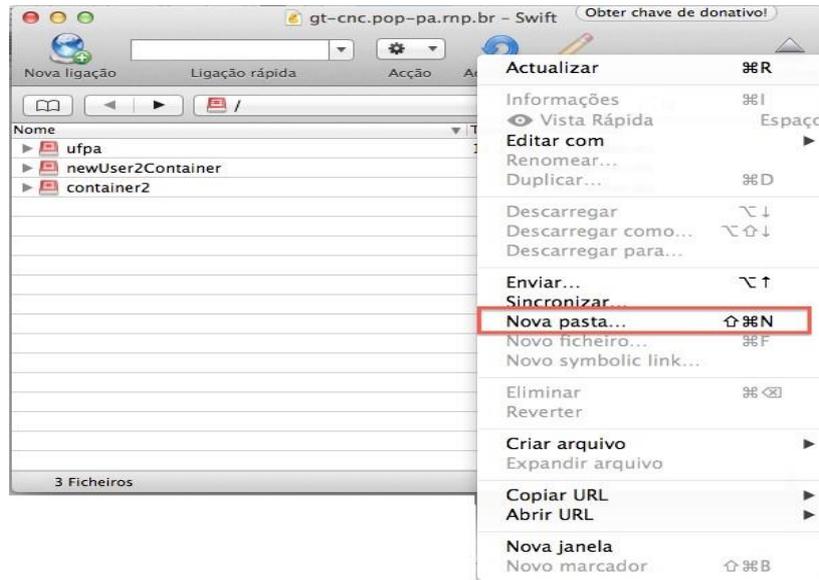


Figura 2. A criação de uma pasta no piloto.

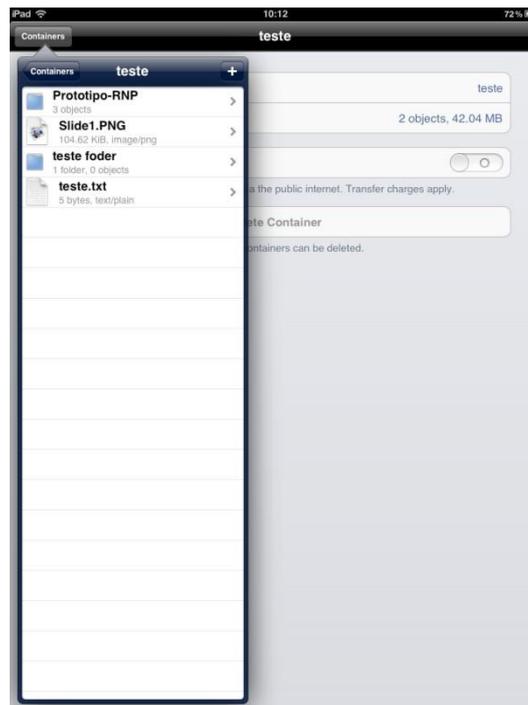


Figura 3. A listagem de pastas e arquivos na nuvem de armazenamento via cliente *mobile IOS*.

2.2 Instituições Participantes

Em sua segunda fase, o GT contará com as seguintes instituições participantes:

- **Universidade Federal do Pará**

Contato: Roberto Araújo

- **Universidade Federal de Santa Catarina**

Contato: Ricardo Custódio

- **Universidade Federal de Minas Gerais**

Contato: Jeroen van de Graaf

2.3 Refinamento do Protótipo

O protótipo definido na primeira fase do GT possui funções semelhantes à de serviços de nuvem disponíveis na Internet. Ele possibilita o armazenamento e a recuperação de arquivos na nuvem a partir de clientes. Além disso, o protótipo permite estender sua capacidade de armazenamento para atender a um grande volume de dados e usuários. No entanto, de forma a ser utilizado no piloto, alguns refinamentos são necessários como listado a seguir:

- ✦ **Autenticação:** O piloto utilizará um novo serviço de autenticação de usuários à nuvem, o *Keystone*. Este serviço possui uma melhor integração com sistema operacional e fornece suporte a diversos tipos de autenticação (LDAP, PAM, etc). Adicionalmente, pretende-se avaliar e integrar o *Keystone* ao serviço identificação por federação CAFé. Para isso, serão realizados estudos relativos à estrutura interna do *Keystone* e do CAFé para a verificar a melhor forma de integração destes serviços.

- ✦ **Acesso a Arquivos via Web:** No protótipo, o acesso aos dados poderia ser realizado através *desktops* e de dispositivos móveis. Isso foi possível graças a programas clientes instalados nesses equipamentos. No entanto, como o acesso aos dados através de programas clientes requer a instalação desses, a utilização de um navegador Web facilitaria o acesso a estes dados. Assim, no piloto pretende-se avaliar e incluir a possibilidade de acesso aos dados via *Web*. Para isso será utilizada os módulos *Tempurl* e *Formpost* do *Swift* como base de desenvolvimento.

- ✦ **Monitoramento do Serviço de Nuvem:** De forma a se tornar um serviço, a nuvem deve dispor de uma ferramenta de monitoramento. Esta ferramenta possibilitaria, por exemplo, verificar quais servidores estão sendo mais utilizados. De forma a prover monitoramento no piloto, o GT pretende avaliar e incluir no piloto uma ferramenta para este fim;

- ✦ **Cliente para Dispositivo Móvel:** No protótipo foram utilizados um cliente para IOS. Esse cliente possibilitou a conexão com a nuvem, a visualização de arquivos, a criação

de pastas, etc. A fim de possibilitar a conexão à nuvem através de dispositivos baseados em *Android*, o GT pretende avaliar e desenvolver um cliente Android baseado em uma API para a nuvem *OpenStack Swift*, por exemplo o *Jclouds*.

- ✎ **Cloud Drive do TERENA:** O GT realizará uma avaliação sobre sua adoção e a integração com o projeto piloto do *Cloud Drive TERENA*. Embora este projeto objetive a integração com nuvens de armazenamento como a que está sendo proposta, uma API de integração ainda encontra-se em desenvolvimento.

2.4 Ferramentas de Suporte à Operação

Na primeira fase do GT foram desenvolvidos *scripts* de instalação e configuração. Estes *scripts* facilitaram a realização destas operações no *software* de nuvem do protótipo. Além disso, o GT também desenvolveu um programa para facilitar a configuração do cliente de acesso à nuvem, *Cyberduck*.

Para a segunda fase, o GT propõe-se a empregar ferramentas para obtenção de métricas (e.g. utilização da CPU, consumo de memória, uso de disco, etc) que facilitem a administração da nuvem como as listadas a seguir:

1. **StatsD:** Para facilitar a operação do piloto na segunda fase, o GT empregará o *StatsD*. Esta ferramenta possibilita a geração de gráficos a partir de métricas do serviço;
2. **Stadslog:** O funciona em conjunto com o StatsD. Ele provém dados sobre o estado das máquinas *proxy* e servidores de armazenamento;
3. **Swift Recon:** Este *middleware Swift* provém métricas gerais dos servidores e específicas do swift como: o hash MD5 de cada arquivo de anel, o tempo de replicação de objetos mais recente, etc.
4. **Swift Informant:** Este *middleware* possibilita a visualização das requisições dos pedidos de objetos pelos clientes em tempo real e pode enviar métricas para o *StatsD*;
5. **Swift Zenpack para o Zenoss:** Extensão para o Zenoss que permite o monitoramento dos serviços do *Swift*.

Além das ferramentas acima, o GT se propõem a fazer uma avaliação sobre a necessidade de programas para gerenciamento de usuários e, se necessário, pretende desenvolver *scripts* para este fim. No entanto, estes programas podem não ser necessários devido à integração do piloto com o serviço de identidades federadas.

O GT também se propõe a desenvolver *scripts* para facilitar a administração da nuvem. Isso inclui a instalação de um novo computador de armazenamento à nuvem e o rebalanceamento dos anéis que formam a nuvem.

3. Cronograma

A seguir é apresentado um cronograma para o piloto.

	2012		2013									
	Nov	Dez	Jan	Fev	Mar	Abr	Mai	Jun	Jul	Ago	Set	Out
Refinamendo Protótipo												
<i>1. (Autenticação) Integração do serviço de autenticação KeyStone</i>	x	x										
<i>2. (Autenticação) Estudos relativos a estrutura interna do KeyStone e do CAFe</i>		x	x									
<i>3. (Autenticação) Integração do</i>			x	x	x	x	x	x				

<i>CAFe com o Keystone</i>												
<i>4. Desenvolvimento de acesso a arquivos via Web</i>		x	x	x	x							
<i>4. Estudo e integração de monitoramento ao serviço de nuvem</i>			x	x	x	x						
<i>5. Desenvolvimento de cliente mobile Android</i>			x	x	x	x	x					
<i>6. Avaliação do piloto</i>			x	x	x	x	x	x	x	x	x	x
<i>Implantação do piloto de armazenamento de dados em nuvem</i>			x	x	x	x	x					

<i>Avaliação da integração do protótipo com o projeto TERENA Trusted Cloud Drive</i>						X	X	X	X	X	X	X
<i>Testes no piloto</i>					X	X	X	X	X	X	X	X

Referências

1. **BOOT, Ed. Sorry, Dropbox, I still don't trust you. ZDNet.** Disponível em: <<http://www.zdnet.com/blog/bott/sorry-dropbox-i-still-dont-trust-you/4173/>>. Acesso em: 03 set. 2012.
2. **DROPBOX.** Site de Serviço de Armazenamento em Nuvem. Disponível em: <<https://www.dropbox.com/>>. Acesso em: 25 ago. 2012.
3. **GOOGLE. Termos de Serviço do Google.** Disponível em: <<http://www.google.co.uk/intl/pt-BR/policies/terms/regional.html>>. Acesso em: 04 set. 2012.
4. **GOOGLE DRIVE.** Site de Serviço de Armazenamento em Nuvem. Disponível em: <<https://drive.google.com/start>>. Acesso em: 25 ago. 2012.
5. **MALPASS, Ian. Code as Craft - Measure Anything, Measure Everything.** Disponível em: <<http://codeascraft.etsy.com/2011/02/15/measure-anything-measure-everything/>>. Acesso em: 04 set. 2012.
6. **MARSHALL, Matt. Dropbox has become “problem child” of cloud security.** Venturebeat. Disponível em: <<http://venturebeat.com/2012/08/01/dropbox-has-become-problemchild-of-cloud-security/>>. Acesso em: 03 set. 2012.
7. **SKYDRIVE. Site de Serviço de Armazenamento em Nuvem.** Disponível em: <<https://skydrive.live.com/>>. Acesso em: 25 ago. 2012.
8. **TERENA. Terena Trusted Cloud Drive.** Disponível em: <<https://confluence.terena.org/display/CloudStorage/TERENA+Trusted+Cloud+Drive>>. Acesso em: 25 ago. 2012.
9. **CYBERDUCK.** Disponível em: <<http://cyberduck.ch/>>. Acesso em: 25 de agosto de 2012.
10. **RACKSPACE HOSTING (IOS).** Disponível em: <<http://itunes.apple.com/br/app/rackspace-cloud/id327870903?mt=8>>. Acesso em: 31 de ago. de 2012.
11. **OPENSTACK (SWIFT).** Disponível em: <<http://swift.openstack.org/>>. Acesso em: 31 de ago. de 2012.

