

Globus Platform-as-a-Service for Collaborative Science Applications

Rachana Ananthkrishnan, Kyle Chard, Ian Foster, Steven Tuecke

Computation Institute
Argonne National Laboratory & University of Chicago
Chicago, IL 60637, USA

Abstract— Globus, developed as Software-as-a-Service (SaaS) for research data management, also provides APIs that constitute a flexible and powerful Platform-as-a-Service (PaaS) to which developers can outsource data management activities such as transfer and sharing, as well as identity, profile and group management. By providing these frequently important but always challenging capabilities as a service, accessible over the network, Globus PaaS streamlines web application development and makes it easy for individuals, teams, and institutions to create collaborative applications such as science gateways for science communities. We introduce the capabilities of this platform and review representative applications.

Keywords—*platform as a service, identity, group, authentication, authorization, profile, transfer, sharing, science gateways, collaboration, cloud*

I. INTRODUCTION

Developers of science gateways [1] and other collaborative science applications frequently need to assign identities to their users, manage user profiles, and organize users into groups—for example, to control access to resources and to track usage. Indeed, this requirement is so fundamental to effective scientific collaboration that it is hard to imagine building a useful collaborative science application without such capabilities. Building on such fundamental capabilities, the same developers frequently then want to provide additional more specialized functionality to their users: for example, the ability to upload, download, and transfer files; to share files with colleagues; to run simulations and so forth.

However, providing high-quality implementations of such capabilities can be extremely challenging, due to the complexity of dealing with varied and sometimes rapidly evolving security protocols, ensuring best practices implementation approaches are applied, and developing efficient and reliable implementations. Custom implementations of such important features often require significant expertise and infrastructure for development, deployment and continued support.

The advent of cloud computing has seen a movement towards outsourcing important functionality to dedicated and professionally hosted providers. Many commonly used applications are now available over the Internet using Software-as-a-Service (SaaS) approaches. Examples include

email, photo and music storage, and collaborative document editing. The advantage of service-based approaches is that software is managed by expert operators who control the entire infrastructure and environment to provide reliable systems; leverage economies of scale to keep prices low; and use elastic compute infrastructure to scale service delivery to peak demands. Globus, formerly Globus Online [2], has successfully applied such SaaS approaches to research data management for the past three years. Thousands of researchers now use Globus's Web and command line interfaces to transfer, synchronize, and share scientific data. In aggregate, more than 30 petabytes and 1 billion files have been moved by Globus in its first three years.

We explain in this paper how the benefits enjoyed by users of Globus can also be leveraged by external applications such as science gateways. Through REpresentational State Transfer (REST) Application Programming Interfaces (APIs) developers of such applications can leverage Globus as a platform. These APIs implement what is commonly referred to in the cloud industry as a Platform-as-a-Service (PaaS): a set of building block operations that are operated reliably by a third party, and to which application developers can outsource important, yet difficult, tasks. PaaS has proved popular because it enables developers to focus on developing and operating their applications, reducing time spent developing support for mundane tasks and increasing reliability and availability through reliance on dedicated services designed to perform specific tasks efficiently. Many commercial vendors now offer PaaS capabilities for commercial users, such as IBM Smart Cloud, Amazon Web Services and Google AppEngine. However, unlike Globus these commercial platforms are not developed for the research community.

The PaaS operations that Globus provides are concerned with assigning and managing identities, managing user profiles, and organizing groups. Other related operations provide for moving, synchronizing, and sharing data. In each case, Globus provides application developers with powerful and flexible management interfaces and REST APIs. Users of applications that make use of these capabilities encounter intuitive web interfaces with a common look and feel across different services. Both developers and users benefit from high-quality implementations, support for a wide range of security protocols, implementation of best practices security approaches, and a highly reliable platform based on

replicated state and services distributed over multiple commercial cloud data centers.

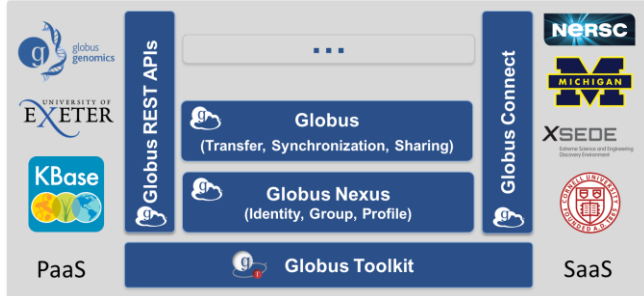


Figure 1: Globus services provide users with research data management functions (right) but also implement REST APIs that permit their use as a platform (left).

The Globus family of services is depicted in Figure 1. Globus platform capabilities are grouped into 2 service interfaces: Globus identity, profile and group management (Globus Nexus) and Globus data transfer, synchronization and sharing (Globus transfer).

Globus Nexus [3] provides identity management capabilities that allow users to create a unique global Globus identity which can be used for single sign-on across Globus services. Users can then link external identities to their Globus identity, enabling them to manage their many identities in a single location and to also use these linked identities to authenticate with other Globus services. Its group management capabilities permit users to create and manage their own groups, assign roles to group members, and use groups for authorization decisions, for example to control access to data using Globus transfer. We provide more details on these capabilities in Section III below.

Globus transfer [4] provides high performance, secure, third party data movement and synchronization between Globus “endpoints”. An endpoint is the name given to a resource on which users may transfer and share data. Endpoint creation requires the installation of a simple agent called Globus Connect to make the resource accessible to Globus transfer. Globus Connect can be used to create either a multi-user server endpoint for Linux or a lightweight personal endpoint for Linux, Windows or MacOS. Globus transfer handles all the difficult aspects of data transfer between endpoints allowing a user to “fire and forget” data transfers, while tuning parameters to maximize bandwidth usage, managing security configurations, providing automatic fault recovery, and notifying users of completion and problems. It also provides in place sharing of data, allowing researchers to share large datasets from their existing storage repositories on which the data resides. In effect, this sharing capability provides Dropbox-like functionality without the need to replicate data in the cloud to be shared.

This paper is structured as follows. Section II presents common usage scenarios of Globus as a platform. Section III and IV describe the capabilities of Globus Nexus and Globus

transfer, respectively. Section V describes how Globus services can be leveraged as a platform by science gateways. Globus deployment infrastructure is summarized in Section VI. Finally, related work is presented in Section VII and the paper is summarized in Section VIII.

II. PLATFORM USAGE SCENARIOS

Globus is used as a platform by a variety of different user communities. The following examples demonstrate the range of scenarios in which Globus functionality is used at different scales, thus providing examples of how Globus services can be leveraged by science gateways.

Earth System Grid Federation (ESGF) [5] is a global federation that provides infrastructure supporting management and analysis of climate data. ESGF infrastructure spans sites across the globe and supports single sign-on across its resources through a web portal. ESGF also provides a metadata catalog along with tools to aid data discovery and integration. Underpinning its primary role as a data distribution portal, ESGF uses Globus transfer as a platform to achieve high performance data distribution. ESGF deploys and manages public Globus endpoints on their resources to access their data. The ESGF portal includes an interface for using Globus transfer via its REST API and endpoint selection web pages. The advantage of this approach is that users of ESGF obtain access to high performance and reliable data transfer of large climate datasets in the context of their data access and through the familiar ESGF portal interface.

The University of Exeter launched, in 2011, the OpenExeter project to investigate how researchers manage their data. The project found that while there was huge demand for the institutional 1PB DSpace [6] repository there were significant technical hurdles when ingesting large datasets. DSpace operated well with small files, in small numbers, using real time transfer over HTTP, but was not equipped to handle large datasets that may be many terabytes in size. To address this shortcoming the OpenExeter project developed a Globus platform-based approach to provide high performance upload and download of data between users and their DSpace repository. In this model they use Globus Nexus to provide single sign-on to their data management services, and OAuth for MyProxy to bridge to their Identity Provider. They also provide high performance asynchronous data upload using Globus transfer through a branded site (<https://go.exeter.ac.uk>).

Globus Genomics [7] is itself a PaaS offering that provides scalable genome analysis pipeline creation and execution on the cloud. Globus Genomics addresses the challenges of running Next Generation Sequencing (NGS) analyses on a large scale by providing state of the art algorithms, sophisticated data management tools, a graphical workflow environment and an elastic cloud-based infrastructure. Globus Genomics leverages both Globus transfer and Globus Nexus as a platform to provide high

performance data movement and identity management. Given that large NGS datasets are often terabytes in size, Globus Genomics uses Globus transfer within its data management framework to enable users to move data from acquisition through analysis and archival. Users can move data asynchronously and reliably through the integrated transfer interfaces implemented in Globus Genomics' Galaxy workflow environment. Through integration with Globus Nexus, users of Globus Genomics can log in with any of their supported external identities (e.g., campus identities, Google account, or username/password). This same identity is also used to authenticate with Globus transfer when moving data to/from Globus Genomics to facilitate single sign-on across the system.

Systems Biology Knowledgebase (KBase) is a large software and data environment designed to enable collaborative generation, testing and sharing of hypotheses about gene and protein functions. KBase provides large scale computing infrastructure to enable analysis and modeling of interactions between microbes, plants and their communities. KBase currently supports more than 400 users spread over several different research groups across the US. To address the challenges associated with managing hundreds of user identities and providing best practices authentication and authorization support, KBase outsources all identity and group management to Globus Nexus. Through a branded Globus site (<https://gologin.kbase.us/>) users are able to create and manage identities and groups. KBase uses the Nexus REST API to access user identities and groups to provision policies within local compute and storage resources. This model allow KBase to define group-specific admission policies and to have group membership be validated in real time using the Nexus APIs.

Biomedical Research Informatics Network (BIRN) [8] is a national initiative to advance biomedical research through an infrastructure designed to support data sharing and collaboration. BIRN relies on Globus platform services to enable these capabilities. Sharing and movement of large datasets is provided by Globus transfer while Globus Nexus manages all BIRN identities and enables creation of user-defined groups used to control sharing and access to collaborative resources. BIRN uses the Nexus REST APIs and LDAP interface to enable user identities and groups to be provisioned across BIRN's management and collaboration services; for example, Globus identities and groups are used in their consortium-wide Confluence Wiki. All identity and group management capabilities are offered through an integrated branded site (<https://access.birncommunity.org>).

Branded sites provide the same Globus capabilities through a customized web presence. Several of the examples described above use branded Globus sites to expose access to Globus services for their communities. These branded sites employ the same Globus infrastructure and provide the same capabilities available on the Globus website, but differ in that they have customized logos, color schemes, and styles to match the branding of the community. This approach saves

significant investment for communities who would otherwise have to recreate web infrastructure that mimics existing Globus interfaces. Other groups that offer branded Globus sites to their communities include the Blue Waters supercomputer at NCSA, the National Energy Research Scientific Computing Center (NERSC), the University of Chicago Research Computing Center, and Indiana University.

Custom integration with the Globus platform is also possible through the published APIs and CLI interfaces. For example, researchers at Pacific Northwest National Laboratory (PNNL) used Globus transfer to support near real-time remote analysis of data collected using synchrotron microtomography at the Advanced Photon Source at Argonne National Laboratory (ANL). The team streamed terabytes of data for analysis and visualization to provide feedback to the experiment in progress. They developed a streaming transfer utility using the transfer API to move data as it is curated from Windows-based acquisition machines at ANL through to PNNL's high performance computing cluster. The Cardio Vascular Research Grid (CVRG) uses the Nexus API to provision users and groups in their Liferay portal and Galaxy deployment. They aim to leverage Nexus's single sign-on capabilities to provide seamless use of a single identity for their distributed user base across their diverse infrastructure.

III. GLOBUS NEXUS

Globus Nexus allows developers of science gateways to outsource identity, profile, and group management functionality. As a platform, Globus Nexus addresses four major obstacles to the creation and operation of high-quality collaborative applications:

1. **Identity provisioning:** Create and manage identities for gateway users.
2. **Identity hub:** Link different user identities, so that, for example a user can authenticate to a gateway with a campus (InCommon) credential.
3. **Group hub:** Create and manage user defined groups which can then be used in authorization decisions.
4. **Profile management:** Manage profile attributes and visibility for those attributes. Profile attributes can be used in authorization decisions, for example to determine who is allowed to join a group.

Importantly all of these features are provided using best practices implementations. For example, user passwords are salted and hashed, all service instances are hosted behind firewalls and use active intrusion detection monitoring, and all identities are backed up daily and stored encrypted.

A. Identity Provisioning

Globus Nexus can act as an identity provider for a project or collaborative science application, providing convenient Web interfaces for identity creation and supporting common

workflows such as email validation and password retrieval. KBase is an example of a project that uses Globus Nexus for identity provisioning: see Figure 2. Nexus currently manages approximately 400 identities for kBase.

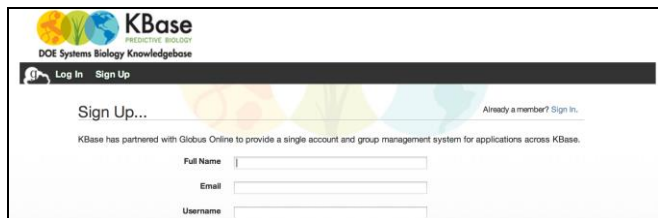


Figure 2: kBase sign up page, showing part of sign-up dialog.

B. Identity Hub

Having created a Globus identity, it is straightforward to link identities from other federated identity providers: see Figure 3. For example, InCommon (via the SAML protocol) [9], Google (via OpenID), XSEDE (via OAuth MyProxy [10]), an IGTF-certified X.509 certificate authority, or SSH key pair.

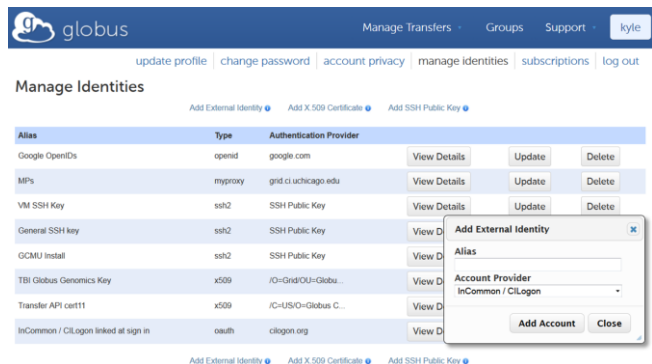


Figure 3: Globus Nexus identity hub, showing user engaged in linking an InCommon identity.

Having linked an identity, the user can then use that identity to authenticate to Nexus and can get mapped as any of the linked identities, including the Globus identity: see Figure 4. Thus, for example, it is straightforward for a user to authenticate to Globus Nexus with their InCommon campus identity: a feature that would provide single sign-on for science gateways using any of the external identities supported by Globus Nexus.

Globus Nexus can act as a federated identity provider to other services, via the OAuth protocol. Various groups have leveraged this capability to enable authentication to Confluence, Zendesk, and other Globus services, for example.

Nexus can also cache, on the user's behalf, delegated credentials (X.509 Certificate) obtained from a third-party service. Thus, for example, a gateway that must access an XSEDE service repeatedly on a user's behalf need not interact with the user repeatedly after a first authentication.

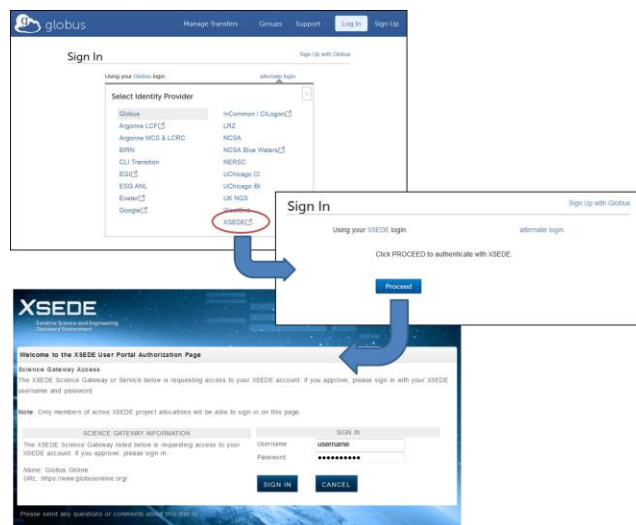


Figure 4: Using a linked identity for authentication to Globus Nexus. Here, an XSEDE identity is used. The OAuth protocol is used to delegate authentication to XSEDE.

We use an example from the BIRN project to illustrate the power of the Globus Nexus identity hub: see Figure 5. BIRN uses Globus Nexus for identity provisioning. In this example, Dr. Smith has created a BIRN identity (a Globus identity created via a BIRN-tailored interface) to which she has linked her campus identity and XSEDE identity. Dr. Smith can then:

- Authenticate to BIRN with her campus identity
- Query a BIRN catalog (BIRN identity)
- Request data transfer from BIRN to campus (BIRN and campus identities)
- Request transfer from BIRN to XSEDE (BIRN and XSEDE identities)
- Repeat these tasks without repeated authentication, thanks to the use of cached credentials.

C. Group Hub

Having created a set of identities, it is natural to want to group them for authorization and related purposes. Globus Nexus provides powerful group management functions, with a particular emphasis on putting users in control of group creation, membership, and properties. Any authorized user can use intuitive Web and REST interfaces to create a group, define its properties (e.g., admission policies, visibility), and

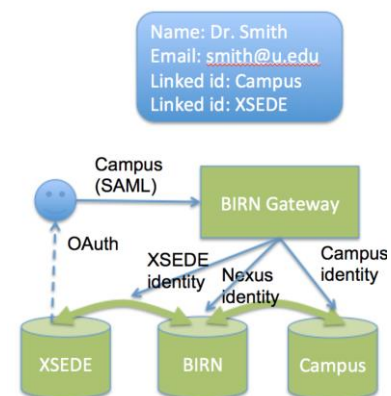


Figure 5: Globus Nexus identity hub example, showing use of multiple identities for different purposes.

Any authorized user can use intuitive Web and REST interfaces to create a group, define its properties (e.g., admission policies, visibility), and

invite other users to join. Groups can be used in authorization decisions, for example to control data sharing via Globus transfer. As with other Globus Nexus services, interfaces can be skinned to meet the needs of specific communities.

Figure 6 shows the users' view of the Globus Nexus group management interface, here skinned for KBase. The KBase project automatically enrolls every user who signs up for a KBase identity in the "kbase_users" group. Subgroups defined within that group are used to organize KBase users who participate in specific KBase project functions.

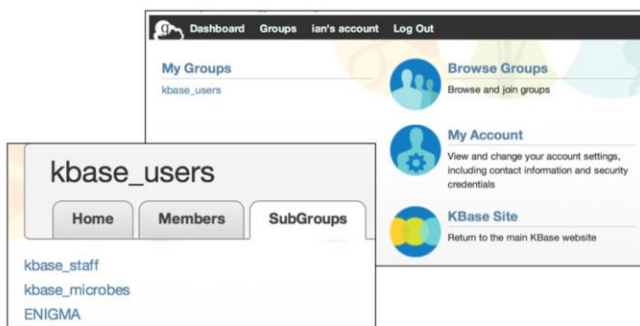


Figure 6: KBase group management interfaces.

Figure 7 shows a different Globus Nexus group management interface, here indicating all groups that include the user "kyle." The subscreen on the right shows the policies that apply to one particular group. Note that the group owner can specify policies that govern visibility (the group is only visible to members) and membership (users must be invited to join and must be approved by administrators and managers).

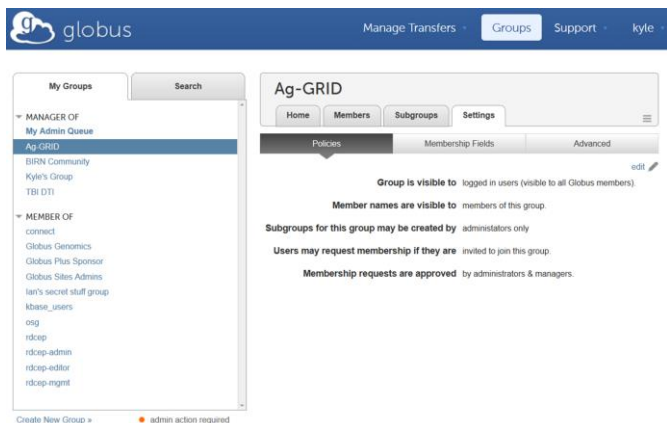


Figure 7: A user view of groups associated with user "kyle."

IV. GLOBUS TRANSFER

Collaborative science applications often require sophisticated data management capabilities across distributed resources. For example, science gateway users frequently need to move large amounts of data between acquisition, storage, analysis and archival. While this task may appear trivial, doing so efficiently and reliably is often a difficult undertaking. Consider, for example, that most users would like to maximize available bandwidth, restart

transfers that fail (due to network unreliability), validate that files have been transferred correctly, and securely transfer data between resources with different security configurations. Globus transfer addresses these challenges by providing an easy to use "fire and forget" data transfer mechanism. In particular, the platform addresses several major challenges for data transfer and sharing in collaborative applications:

1. High performance and reliable: All transfers are automatically managed to maximize bandwidth usage, check data integrity, and restart failed transfers
2. Provide third-party transfers across security domains: Support authentication at either end of transfer and facilitate "third party" transfers between two remote endpoints.
3. Share data with collaborators: Enable data sharing with gateway users and groups directly from existing endpoints without moving data to the cloud.
4. Easily expose local resources as Globus endpoints: Enable users to expose their local resources as Globus endpoints using lightweight client software deployed on their resources.

A. High performance and reliable data transfer

GridFTP [11], the underlying protocol used by Globus transfer, provides dozens of different options that a user can configure when transferring data. For example, developers can configure the TCP buffer size, the number of concurrent control channel connections, and the number of TCP channels used. While this flexibility can enable high performance transfer under a range of conditions, most users lack the expertise required to tune parameters. Globus transfer tunes these parameters before and during transfer to ensure that good settings are used. In particular, it uses a set of heuristics to configure parameters based on the number and size of files in a transfer. For many small files Globus transfer will use pipelining to ensure that multiple files are transferred at a time. For larger files Globus transfer uses concurrent streams to parallelize transfers. Parameters can also be tuned between batches within the same transfer.

Globus transfer includes a number of user-configurable reliability and integrity options. By default, all file transfers initiated via Globus transfer will recover from intermittent network failures, automatically restart when endpoints are disconnected or reconnected, and alert users of other issues (credential expiry, quota exceeded, permission denied), allowing restart when problems are resolved. In addition, a range of other options can be configured based on the requirements of a particular transfer. Users can specify that transfers must validate the integrity of files after completion using individual checksums, that transfers be encrypted, that modification times be maintained, and that complex synchronization semantics be applied. Figure 8 shows the

different user-configurable options supported by Globus transfer.

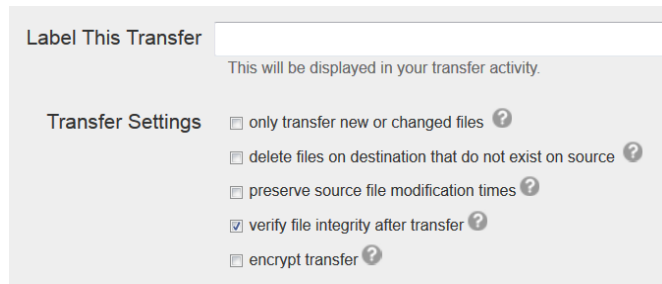


Figure 8: Globus transfer configuration options.

B. Third party transfers across domains

When moving large datasets it is common to require movement between two remote computers (a third party transfer) rather than between a client and a remote server. An example of such a third party transfer between the University of Chicago’s Research Computing Center and Texas Advanced Computing Center’s (TACC) Ranch storage system is show in figure 9.

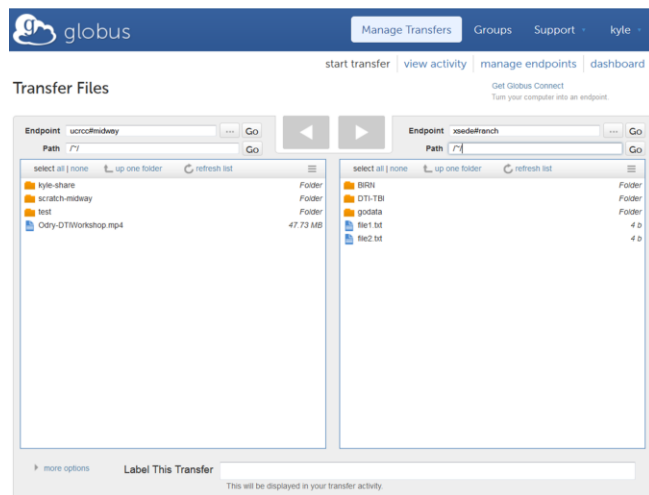


Figure 9: Globus transfer interface showing a transfer between a research computing center and an XSEDE resource.

The difficulty with third party transfers is that a third party system must orchestrate the transfer and users must authenticate with both ends of the transfer, potentially using different security protocols and credentials. In the previous example, UChicago RCC uses CILogon campus credentials while TACC Ranch supports XSEDE proxy credentials. Globus transfer tracks the security credentials required by different endpoints, and prompts the user to authenticate using the appropriate mechanism. In the case of MyProxy, authentication short-term proxy credentials are granted that can be used to access the endpoint for the duration of the credential (often 12 hours). Where possible, Globus transfer (using Globus Nexus) caches these “activated” proxy credentials so that transfers can be executed from the same endpoints or endpoints that rely on the same identity

provider (i.e., MyProxy server) throughout the lifetime of the proxy certificate. Users may choose to deactivate an endpoint at any time which will remove the active proxy credential from Nexus.

C. In place data sharing

Collaborative research frequently requires the sharing of potentially sensitive data among collaborators. Dropbox and Box.net are often used for this purpose, but both have limited storage capacity, charge users for data storage used (even if researchers have significant storage resources available themselves), and require that users upload data to the cloud. Thus, it is infeasible for large research-scale datasets. Alternatively, researchers may create accounts for their collaborators on their internal data repository; however, this approach is often not permitted by administrators, presents an additional burden of account lifecycle management for ephemeral collaborations, and represents an ad hoc process that complicates the tasks of auditing data access and ensuring security protocols are associated with data residing on the repository.

Globus transfer simplifies this common sharing scenario by allowing research-scale datasets to be shared in place on existing storage servers without creating user accounts on those servers. Users can create virtual “Shared endpoints” rooted at any file system location on their existing endpoint using the Globus transfer web, CLI, and REST interfaces. Users may manage access to different levels of the shared file system defining both read and write access to files and directories for Globus users and groups. They may also view, revoke or change these permissions at any time. Figure 10 shows a web interface used to modify sharing permissions through Globus. In, this example data is shared as read only with user ‘Ian’.

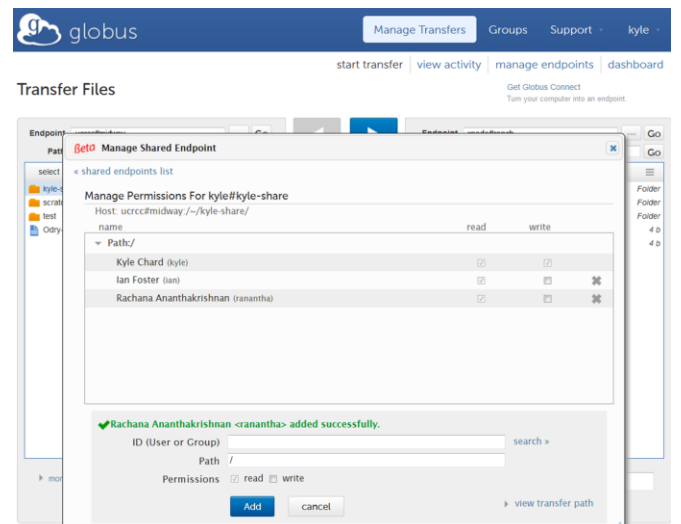


Figure 10: Globus transfer sharing permissions.

D. Expose local resources using Globus Connect

Globus has a well-established network of existing GridFTP servers distributed across the world; all of which can be exposed as Globus endpoints to Globus transfer. Examples of existing GridFTP servers include national cyberinfrastructure providers such as XSEDE and OSG, large compute providers such as NERSC and TACC, and university and group resources such as the University of Chicago's RCC. Globus transfer users therefore do not need to install any software to transfer data between these existing Globus endpoints. However, when users need to move or share data from machines that do not have GridFTP servers installed, such as a user's personal computer, they must install a GridFTP server. To minimize the difficulty of installing a GridFTP server, we provide a packaged client install called Globus Connect. Globus Connect is available in two varieties: Globus Connect Personal and Globus Connect Server.

Globus Connect Personal is designed to enable users to add their personal computers to the Globus ecosystem. It features a one click install and can be configured to run on startup, much like the clients that must be installed to access Google Drive and Dropbox. Globus Connect Personal is available for Linux, Windows, and MacOS. It includes a packaged GridFTP server that runs with user permissions and a GSI-OpenSSH client that is preconfigured to establish an authenticated connection to Globus transfer. Globus Connect Personal is also able to operate behind firewalls and Network Address Translation (NAT) through its use of outbound connections and relay server.

Globus Connect Server is designed for multi user server deployments. It is available as native packages for Debian and RPM systems. The installation process enables administrators to create a fully configurable endpoint including a full GridFTP server and the optional deployment and configuration of a MyProxy server for authentication. The MyProxy server enables end users to authenticate using their local accounts on the server, it can also be configured to use other forms of authentication such as LDAP or campus credentials via CILogon. For added security the MyProxy server can be configured to use OAuth, in this mode users are redirected from Globus Nexus to the MyProxy server to authenticate via a web form before starting a transfer. This mode has the advantage that no passwords are ever seen by any Globus services; instead a short-term proxy credential provides access to the endpoint.

V. GLOBUS PLATFORM FOR SCIENCE GATEWAYS

There are several mechanisms by which science gateways and other collaborative applications can use Globus as a platform. They can access Globus capabilities directly through its various interfaces, they can offer customized branded sites exposing access to Globus functionality, and—for many of the Globus web interfaces—they can use web interfaces directly via redirection from a gateway.

A. Interfaces

Globus provides multiple interfaces through which its capabilities can be integrated in science gateways including a web interface, REST API, programming clients, and LDAP. The primary interface to both Globus Nexus and Globus transfer is through the Globus website (www.globus.org). This web interface provides an easy to use and intuitive means for users to manage identities, groups, and transfers.

Both Globus Nexus and Globus transfer provide REST APIs for programmatic invocation, enabling developers to interact with all features of the respective service. Two sample programming language clients (Python and Java) are provided for each service to simplify access directly from external applications. Full REST API Documentation and sample programming clients are available online (www.globus.org/platform).

Science Gateways can integrate complete identity, profile, groups, transfer, synchronization and sharing capabilities by using the Globus Nexus and Globus transfer REST APIs in their applications. Gateways can create, retrieve and manage user identities and groups, create and manage shared endpoints, and start, monitor and manage transfers. Authentication to the APIs is based on the OAuth 2 protocol allowing users to grant permission to the requesting gateway to access a specific Globus service on their behalf. Importantly users do not enter their Globus (or external identity) passwords in the gateway application directly. Rather the user is redirected to Globus Nexus to authenticate (using any of their linked identities). Globus Nexus in turn grants the requesting gateway a short term access token through which the gateway can act on behalf of the user.

Globus Nexus implements a dynamic LDAP interface to satisfy the requirement for simple read-only access to user profiles and groups. LDAP is a common protocol for implementing user and group directories and is supported by many third party applications. Using the LDAP interface gateway developers can quickly and easily leverage Globus identities and groups in their own applications. The LDAP interface allows users to bind using their Globus identity username and password. After successful authentication the LDAP directory can be explored using common "ldapsearch" commands to view profiles and groups visible to the authenticated user. Unlike the REST API, the LDAP interface provides read-only access to identities and groups; it also does not enable users to authenticate using external identities.

Globus transfer implements a command line interface (CLI) for advanced users. This interface enables client-side script-based invocation of the Globus transfer API which can be easily included in user scripts. Rather than rely on the typical CLI approach requiring installation of client-side libraries, Globus transfer instead provides a novel SaaS-based CLI [4] through which users can connect (using SSH) to a restricted shell that enables execution of transfer commands. For example, a user can transfer files between

two endpoints (go#ep1 and go#ep2) using the CLI with the following command:

```
ssh kyle@cli.globus.org scp \  
go#ep1:/share/godata/file1.txt \  
go#ep2~/myfile.txt
```

B. Branded sites

Science gateways may want to use Globus interfaces directly in their own website. To address this need Globus supports two models: 1) creation of branded sites, or 2) integration of Globus web interfaces via redirection.

A branded site is a full Globus deployment that has been tailored to a particular project's branding with customized logos, color schemes, and layouts. While the underlying Globus Nexus service is shared among many branded sites, the approach enables different communities to have their own interface to the Globus ecosystem. Examples of different branded sites are presented in Section II. Each group can use their own URL and branding on their instance of Globus; however, the entire software stack is provided via Globus, leveraging the community support provided by Globus Nexus. Branded site developers can preselect their community identity provider, manage policies around access to their branded site, and also provide integrated access to other Globus services. Each branded site has its own root community group that is preselected under that branded view. This feature allows users of the branded site to have a community specific view of their groups and members.

Several Globus web pages have been designed to be accessible directly by external applications to facilitate integration into custom workflows. That is, external applications can compose URLs to the Globus website that include preselected options and redirection commands to return the user to the calling application. In this model application users are redirected to Globus to perform their action (e.g., file or endpoint selection) and then they are redirected back to the external application. One example is the "Select Destination" page which allows application developers to create a download page by preselecting a source endpoint and path. End users of the application are then presented only with a destination selection interface to download datasets from a fixed endpoint and path. This approach provides a simple way for gateways to use complex Globus transfer functionality without having to develop their own user interfaces from scratch.

VI. PRODUCTION DEPLOYMENT

Providing reliable platform and software services requires highly scalable and reliable infrastructure. To meet this need, we deploy all Globus services on a commercial cloud and implement best practices approaches to deployment (e.g., automated deployment processes, replicated instances, spread across availability domains) and operation (e.g., monitoring and intrusion detection).

The Globus infrastructure, like any cloud service, must be able not only to grow to support an increasing user base but also to overcome failures to provide high reliability and availability. Based on user experience requirements Globus Nexus and transfer have different availability goals. The availability goal for Globus Nexus is 99.99% as it underpins the entire service and without it users would be unable to log in. Globus transfer, on the other hand, aims for 99.9% availability as short outages can be more easily tolerated by users. In 2013, both Globus Nexus and transfer met these goals: Nexus achieved 99.99% availability while transfer achieved 99.96% availability. The combined availability of both services was 99.95% with only six unplanned outages.

To maintain these levels of availability we leverage experiences learned by commercial platform and service providers by using Amazon Web Services (AWS), a commercial Infrastructure-as-a-Service provider. Globus Nexus and Globus transfer use a range of different services provided by AWS including Elastic Compute Cloud (EC2), Simple Storage Service (S3), Elastic Load Balancing (ELB), Simple Notification Service (SNS), and Scalable Domain Name System (Route 53). These services are paid for through a combination of federal grants and a subscription-based sustainability model [12].

The different availability goals of Globus Nexus and transfer are reflected in their respective deployment architectures: Globus Nexus leverages multiple load balanced services and databases spread across Availability Zones (AZs), whereas Globus transfer has only single service and database instances in a single AZ. Globus Nexus and its constituent databases (Cassandra and ElasticSearch) are hosted on a group of EC2 instances in the US East region. The Nexus application consists of three replicated services managed behind a load balancer that assigns requests to these three instances. There is no state stored directly on these instances, rather a pool of data storage nodes are operated using replicated backups between the nodes. Nexus relies on three Cassandra instances and two ElasticSearch instances to provide storage and search capabilities respectively. Globus transfer includes several instances deployed in the US East region for the transaction database, transfer agents, history database, transfer REST API server, CLI server, and backups. Collectively, 13 amazon instances are used to deliver Globus.

A further 8 instances are used for services that are directly responsible for supporting the production infrastructure, such as deployment, monitoring, logging, and intrusion detection across Globus infrastructure. Nagios monitors each deployed service and, upon failure, sends alerts to the Globus operations team. An external logging service ensures that all logs generated from all nodes are stored reliably; if an instance fails, developers can still investigate the logs generated immediately before failure. An Open Source SECurity (OSSEC) service is also hosted to perform intrusion detection across the deployment infrastructure using a number of approaches.

Security is a primary concern for the Globus architecture, given that user identities and active credentials are stored in the Nexus database. Thus, best practices approaches are used when storing data and securing instances. For instance, Amazon's security groups are used to restrict communication between instances and the general Internet. These security groups also restrict access to predefined ports. For example, the Cassandra and ElasticSearch nodes are only accessible from the Nexus nodes, and only on the published API ports. All sensitive information is encrypted before being stored, and interfaces that can access this information are restricted using access tokens. All communication with the REST API is encrypted using SSL.

VII. RELATED WORK

Several service providers offer capabilities similar to individual Globus components such as identity management or data transfer. The Agave API [13], which originated in the iPlant Collaboration [14], provides authentication and authorization, job submission, and data management through a REST interface. Like Globus, its major goal is to simplify common yet difficult IT operations associated with research. Rather than expose a web interface, Agave includes a sample GUI that demonstrates how its REST interfaces can be used.

Many commercial service providers offer identity management and authentication services. Social network identities in particular, are commonly used for authentication by other services, examples include Facebook Connect and Google Accounts. In research domains, the CILogin authentication framework allows users to authenticate using a campus identity. Amazon Identity and Access Management (IAM) [15] enables user management across Amazon services and resources. It also enables federation of public identities such as Facebook and Google. Like Nexus, these capabilities allow service developers to integrate their different identities into their services. Atlassian Crowd [16] provides identity management capabilities for web applications. It enables user identities to be sourced from several directories (e.g., LDAP) and exposes different authentication interfaces that can be embedded in external applications (e.g., OpenID). Both IAM and Crowd are commercial applications that require paid subscriptions; they are also designed to support commercial identity providers.

Group management and authorization services are also available commercially and academically. For example, Google Groups provides user defined groups that can be used for authorization to Google services. Previous work from Grid computing such as the Virtual Organization Management Service (VOMS) [17] provides group based authorization capabilities using short-lived proxy credentials to Grid resources. Atlassian Crowd provides user defined groups that can be incorporated in external applications. Grouper [18] is perhaps most similar to Globus Nexus in its group management capabilities. Nexus is distinguished by its focus on user-driven group management.

File transfer is a common activity for many researchers. Most rely on tools that are available on their computer such as rsync, SCP, and FTP. These tools enable movement between a client and a remote location, but typically they do not support third party transfers or provide the optimization and reliability enhancements available when using Globus transfer. GridFTP, the underlying protocol used by Globus, offers tools (e.g., globus-url-copy) that can move large amounts of data; however they requires installation locally and configuration of complex parameters to optimize transfers. Research projects such as bbftp [19], Kangaroo [20], CATCH [21] and Stork [22] provide data movement over wide-area networks. While the focus of these works is performance oriented (like Globus), they do not use SaaS and therefore represent complex integration challenges when incorporating in a gateway, moreover none have the huge network of existing endpoints available in Globus.

Commercial applications such as Dropbox (www.dropbox.com) and Box (www.box.com) provide mechanisms to upload and share data from a Cloud based storage repository. However, neither is well equipped to handle large datasets as the cost and overhead of moving data may become prohibitive. Like Globus transfer, Aspera (<http://asperasoft.com/>) provides high performance data transfer between Aspera servers, it also provides a sharing mechanism that enables sharing directly from its servers; however Aspera is a paid service and primarily targets commercial users.

VIII. SUMMARY

The Globus family of services and APIs provide a powerful platform to which developers of collaborative science applications can outsource important identity, group, and data management operations. The availability of this platform allows developers to focus on their core application logic while creating rich user environments that can manage different identities, provide user-managed groups, transfer data among the gateway and other endpoints, and share data in place with other community members. Importantly, developers can rely on high availability of these services due to the Globus professional hosting and support model and deployment on Amazon cloud. These capabilities, while certainly not challenging from a scientific perspective, represent a significant undertaking for application developers, particularly if they wish to ensure that best practices are met and to provide robust and reliable implementations. Globus services have been successfully leveraged by a wide range of small to large collaborations, representing supercomputers, research computing centers, national collaborations, and even universities, with more than 13,000 registered users, several hundred groups, and over 30 petabytes moved to date.

ACKNOWLEDGMENTS

We thank the Globus team for their work implementing and operating Globus. We also thank the XSEDE

architecture team for their contributions to our understanding of requirements and Von Welch for his Globus security review. This work was supported in part by the NIH through NIGMS grant 5U24RR025736, the NSF through grants OCI-1053575 and OCI-0534113, and the DOE through grant DE-AC02-06CH11357.

REFERENCES

1. Wilkins-Diehr, N., *Science Gateways – Common Community Interfaces to Grid Resources*. Concurrency and Computation: Practice and Experience, 2007. **19**(6): p. 743–749.
2. Foster, I., *Globus Online: Accelerating and democratizing science through cloud-based services*. IEEE Internet Computing, 2011(May/June): p. 70-73.
3. Ananthakrishnan, R., et al. *Globus Nexus: An identity, profile, and group management platform for science gateways and other collaborative science applications*. in *Science Gateway Institute Workshop, co-located with IEEE Cluster*. 2013.
4. Allen, B., et al., *Software as a service for data scientists*. Commun. ACM, 2012. **55**(2): p. 81-88.
5. Cinquini, L., et al. *The Earth System Grid Federation: An open infrastructure for access to distributed geospatial data*. in *E-Science (e-Science), 2012 IEEE 8th International Conference on*. 2012.
6. Smith, M., et al., *DSpace An Open Source Dynamic Digital Repository*. D-Lib Magazine, 2003. **9**(1).
7. Bo, L., et al. *Deploying Bioinformatics Workflows on Clouds with Galaxy and Globus Provision*. in *High Performance Computing, Networking, Storage and Analysis (SCC), 2012 SC Companion*:. 2012.
8. Helmer, K.G., et al., *Enabling collaborative research using the Biomedical Informatics Research Network (BIRN)*. Journal of the American Medical Informatics Association, 2011.
9. Barnett, W., et al., *A Roadmap for Using NSF Cyberinfrastructure with InCommon*. 2011.
10. Anderson, K.M. *Integrating Open Hypermedia Systems with the World Wide Web*. in *Eighth ACM Conference on Hypertext*. 1997. Southampton, UK.
11. Allcock, W., et al., *The Globus Striped GridFTP Framework and Server*, in *Proceedings of the 2005 ACM/IEEE conference on Supercomputing*. 2005, IEEE Computer Society. p. 54.
12. Foster, I., V. Vasiladiadis, and S. Tuecke, *Software as a Service as a path to software sustainability*. 2013, figshare. <http://dx.doi.org/10.6084/m9.figshare.791604>.
13. Dooley, R., et al. *Software-as-a-Service: The iPlant Foundation API*. in *5th IEEE Workshop on Many-Task Computing on Grids and Supercomputers (MTAGS)*. 2012.
14. Goff, S.A., et al., *The iPlant Collaborative: Cyberinfrastructure for Plant Biology*. Frontiers in Plant Science, 2011. **2**.
15. Bajaj, C. and S. Cutchin. *Web based Collaborative Visualization of Distributed and Parallel Simulation*. in *IEEE Parallel Symposium on Visualization*. 1999.
16. Anderson, K.M., R.N. Taylor, and E.J.W. Jr., *Chimera: Hypermedia for Heterogeneous Software Environments*, in *ACM Transactions on Information Systems*. 2000. p. 211-245.
17. Alfieri, R., et al., *From gridmap-file to voms: managing authorization in a grid environment*. Future Generation Computer Systems, 2005. **21**(4): p. 549-558.
18. Anderson, K.M. *Supporting Industrial Hyperwebs : Lessons in Scalability*. in *21st International Conference on Software Engineering*. 1999. Los Angeles, CA.
19. Hanushevsky, A., A. Trunov, and L. Cottrell. *Peer-to-Peer Computing for Secure High Performance Data Copying*. in *2001 International Conference on Computing in High Energy and Nuclear Physics 2001*. Beijing.
20. Thain, D., et al. *The Kangaroo approach to data movement on the Grid*. in *High Performance Distributed Computing, 2001. Proceedings. 10th IEEE International Symposium on*. 2001.
21. Monti, H.M., A.R. Butt, and S.S. Vazhkudai, *CATCH: A Cloud-Based Adaptive Data Transfer Service for HPC*, in *Proceedings of the 2011 IEEE International Parallel & Distributed Processing Symposium*. 2011, IEEE Computer Society. p. 1242-1253.
22. Kosar, T. and M. Livny, *A framework for reliable and efficient data placement in distributed computing systems*. J. Parallel Distrib. Comput., 2005. **65**(10): p. 1146-1157.