A Case for Multimedia Streaming Over the Grid Infrastructure

Lambros Lambrinos¹ and Fotis Georgatos²

Dept. Of Communication and Internet Studies, Cyprus University of Technology Limassol, Cyprus lambros.lambrinos@cut.ac.cy 2National Technical University of Athens, Athens, Greece fotis@mail.cern.ch

Abstract. The grid infrastructure consists of nodes all over the world which are usually interconnected with high speed network links and provide storage and processing facilities. In this paper, we investigate how this massive infrastructure can be utilized to facilitate efficient and scaleable real-time multimedia streaming. The aim is to avoid the issues usually associated with one-to-many media streaming architectures through the use of a mechanism that initiates reflectors as and when they are needed thereby reducing bandwidth-related bottlenecks and ensuring that the delay between the last media distribution point and the receiving client is as low as possible.

Keywords: grid, multimedia, streaming, architecture

1 Introduction

Streaming multimedia data over the packet-switched internet is nowadays a highly popular application. End-user hosts have more than adequate processing power for the decoding of high quality audiovisual streams; mobile devices are also capable of displaying reasonable quality video. The material distributed is usually pre-recorded but an increasing number of live events (ranging from speeches to music festivals and sports events) are streamed in real-time.

The increased user demand for multimedia streaming results in increased network bandwidth requirements. As the data is predominantly carried on a best-effort basis, issues such as jitter and packet loss degrade the user's experience. These problems are not uncommon and are somewhat expected considering that the media data competes with other internet traffic. Researchers have always been examining ways to reduce the load on the infrastructure as this allows more users to be served.

An infrastructure that could be used for multimedia streaming is the grid [1]. So far, the grid is predominantly seen as a massive infrastructure that can be used to perform intensive processing operations and store vast amounts of data. To facilitate

these activities, the various nodes (clusters) are in many cases interconnected using high speed links. In this paper we propose an architecture that aims to utilize this exact characteristic and use the grid for scalable real-time multimedia streaming.

The paper starts with a brief introduction to the main characteristics of media streaming technologies and the grid. Our rationale for deploying multimedia streaming services over the grid is then presented followed by a description of our proposed architecture and its components and future work plans.

2 Background

The use of the internet for distribution of multimedia data to multiple clients has been a highly active research topic for many years. It is expected that it will remain so since commercial applications are already using IP networks for their data delivery [2] and the internet already "hosts" an enormous collection of audiovisual material.

2.1 Multimedia Streaming

In our work we concentrate on streaming data from a single source to multiple receivers. Streaming involves the transmission of media data over the packet-switched network. The source of such data may be a stored file or a live feed from audiovisual equipment. A receiving host expects a timely and loss-free delivery of the data packets in order to decode the media streams and present them to the user.

Solutions designed for this purpose are generally categorised as Video-on-Demand (VoD) systems. True Video-on-Demand systems [3] allow users to have total control: they can pause, rewind, fast forward the stream whenever they like. Such activities put additional strain on the provider's resources; to prevent that, near Video-on-Demand [4] systems allow such user actions but at discrete time intervals.

The one-to-many data distribution model prevents such systems from scaling as the media server has a finite bandwidth capacity available. Scaleable multimedia data distribution is a topic that has attracted a lot of research interest in the past years; various solutions [5,6] have been proposed to increase the scalability of systems. The ideal solution for the unnecessary bandwidth consumption (i.e. the multiple unicast streams) is the use of native multicast [7] data distribution; this is not yet supported in the majority of the internet.

2.2 The Grid

One definition of the grid is that it is a technology based on a system that collects users and resources in a common infrastructure even if they belong to multiple independent organizations, carriers, companies etc. In practice, this implies that we can view all these resources as a single entity independent of geographical location; the primary reason is that networks such as the internet, make it irrelevant. This method of organization can (and does) have tremendous impact in multiple scientific and/or commercial activities that have to manage systems, data and computations in

an intensive way. Furthermore, grids can impact end-users providing benefits both in content search applications as well as content storage and retrieval of any type currently in use: text-audio-video. One particular aspect of the grid is that it combines storage, networking and computational resources in a single infrastructure. This makes it possible to perform all functions necessary for:

- building and running Digital Libraries
- performing on-the-fly signal compression-decompression or de/multiplexing
- doing capacity provisioning/scheduling on both cpu and network resources
- allocating capacity dynamically either based on on-demand requests or on static traffic patterns

3 Video Distribution over the Grid

In this section we present the rationale behind the suitability of the grid as a platform for multimedia streaming and describe our proposed architecture. Briefly put, a fixed infrastructure investment is never the optimal solution so the grid is the right environment for applications that exhibit variable resource utilization.

3.1 The grid as a multimedia streaming platform

The concept of utility computing is highly applicable in multimedia content delivery; the consumers of the content are typical human beings following diurnal, weekly or monthly patterns with highs and lows in service utilization levels. This is highly evident in TV broadcasting during popular programs.

If we try to transfer this kind of media broadcasting service on the internet as we know it currently, the end result would be congestion, low latency and bad performance at the very moment of highest demand. The most critical aspect of serving an online user community is service stability. It is a good observation to point out that the Internet, as currently experienced by the vast majority of its users, is a network without guarantees and as such it is unsuitable to provide, for example, a better replacement for existing TV and Radio transmission technologies.

This happens because the dependability aspects of the latter are much more robust and resilient to usage by an excessively large audience population. In order for the internet to be able to reach similar levels of robustness and resilience, a number of items have to be addressed:

- a federation-capable solution is needed, if many content providers are going to coexist in the same infrastructure.
- reflectors have to be placed at or near network branching points and reflector capacity must be tunable on demand at run-time.
- bandwidth reservation must be a standard automatic network service, without a human in the loop, even across network boundaries.
- multiple paths should be provided for media data delivery, so that no single system failure or malfunction can disrupt service.

In fact, grid systems are supposed to be able to manage all these aspects at once:

- 1) grids are by definition a collection of resources spanning multiple administrative domains. As a result, technology provisioning for many Certification Authorities [8] and entities below them is considered as standard in grids. On top of that, under the concept of Virtual Organizations [9] one is able to collect entities (resources or users) at arbitrary sets together so an infinite number of policies can be globally applied.
- 2) grid clusters are typically located near backbone network edges or central junctions, for the very fact that this is the only way to make efficient use of them. Moreover, grid clusters contain multiple CPUs which makes it possible to run as many reflectors as needed on a given moment.
- 3) bandwidth reservation and network Quality-of-Service are not new topics and they play a major role in many distributed applications today. What is though a whole new topic, is the capability to be able to offer them in a fully automated manner. This will build guaranteed paths across multiple networks (autonomous systems in Internet parlance) querying in the process and following distinct network policies in a federated manner. Recent experience during summer 2007 in European networks attached to GEANT2 [10] showed that this is feasible and the current plan is to integrate this service as a part of the capabilities of grid systems.
- 4) the internet has provisioning for resilience through the automatic rerouting of traffic through multiple paths. Grids should be able to overcome transient or even permanent failures in a similar manner and in a way which is transparent for the user. The concept applies to all subsystems comprising the grid (CPUs, network links etc.)

It is important to note that although peer-to-peer solutions have been proposed as the underlying mechanisms for multimedia data streaming, the grid offers a far larger and more stable infrastructure that is already being used for data storage and computational purposes. To illustrate our point, a project driven by CERN that involves conducting the largest physics experiment ever, is exploiting grid technology exclusively for the distribution and management of its datasets.

3.2 Grid enhancements

One particular aspect that is related to points 2 and 3 in the previous section, which is worthy of further discussion, is that capacity usage in either the CPU or the network must be a schedulable resource. In the current grid implementation this is not always possible, for the very reason that existing grids are of the batch form, since this is the simplest way to build such a service. For example, in the EGEE grid currently it is not possible to request an allocation of 100 CPUs for the first day of the next month; the basic method of service is based on FIFO queues at Computing Elements (CEs).

In order to schedule CPU and network usage, a slightly revised CE architecture must be developed that includes a scheduler that supports job preemption. Although at first glance someone might comment that such a scheme is not directly supportable by the existing grid, the fact is that the implementation is feasible with only slight modifications. There exist local scheduling systems that are able to do so [11] and have already been tested in the real grid environment. Grid engineers are also currently building the components necessary to provide network performance guarantees for grid-based applications [12]. This is work of very high complexity since the underlying network infrastructure must be able to provide such a service.

3.3 System architecture

The architecture of our proposed solution is shown in Figure 1 below. Essentially our proposed solution falls under the category of application-level multicast data distribution systems.

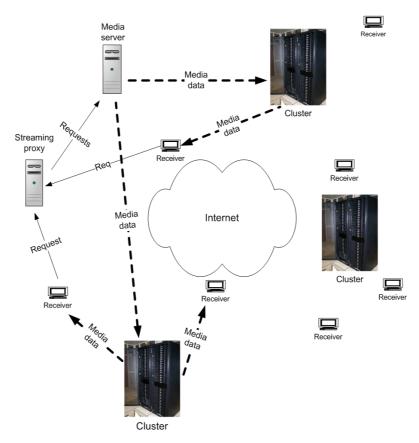


Fig. 1. A grid-based multimedia streaming architecture

The *streaming proxy* is the brain of our architecture. It is responsible for handling client requests and associating clients with the most appropriate reflector. The reflectors run on worker nodes in the grid *clusters*; they are responsible for sending the media data to the receivers that the streaming proxy has allocated to them. The media data is sent to the reflectors by the *media server*. For clarity purposes, some receivers in the diagram appear to be idle.

The different stages of the system's operation are as follows:

- the streaming proxy receives a client request
- proxy identifies the best reflector for the client
- media server requested to start sending data to the reflector (if its new)

reflectors stream data to their clients

The key research issue in the operation outlined here is the identification of the best reflector for the particular client. The definition of "best" here implies that the reflector can accommodate the client (i.e. has not reached its bandwidth capacity) and the round-trip time between client and reflector is reasonable and fairly static.

4 Conclusion and future work

In this paper we presented our ideas and potential solution for the provision of scalable multimedia streaming services by exploiting the grid infrastructure's wide area and stability characteristics. At the current stage, the proposed architecture does not designate any particular protocol or signaling approach since it is agnostic to them. The benefit of this approach is that it provides a model for diverse implementations.

Our ultimate goal is the development of an algorithm that dynamically optimizes the placement of reflectors as clients join and leave during the "broadcast" of a live event. In the initial parts of our work we will concentrate on analyzing the Grid to ascertain whether it fully meets our requirements as they were defined in this paper; if that is not the case then steps have to be taken towards that direction before a reflector placement algorithm can be defined. In our work we will utilize techniques and solutions that are already developed for VoD systems and modify the Videolan [13] streaming software to add extra messaging functionality as it will form our server and client software.

References

- 1. Foster, I., Kesselman, C.: The Grid: Blueprint for a New Computing Infrastructure, Morgan Kaufmann, San Francisco (2004)
- 2. IPTV Industry, http://www.ipty-industry.com
- Fonseca, N.L.S., Rubinsztejn, H.K.: Channel allocation in true video-on-demand systems. In: Globecom 2001.
- 4. Profeta, E, Shin, K.: Scheduling Video Programs in Near Video-on-Demand Systems, In: ACM Multimedia 1997
- 5. Nguyen, T and Zakhor, A.: Protocols for distributed video streaming, In: ICIP 2002
- 6. Gialama, E. et al.: Distributed Video Server for Streaming, In: CSCC 2001
- 7. Deering, S.: Multicast routing in internetworks and extended lans. In: SIGCOMM 1995
- 8. Astalos, J.: International Grid CA Interworking, In: EGC 2005
- 9. Foster, I., Kesselman, C., Tuecke, S.: The Anatomy of the Grid: Enabling Scalable Virtual Organizations, Journal of High Performance Computing Applications, 2001
- 10. GEANT, http://www.geant.net
- Etsion, Y., Tsafrir, D.: A Short Survey of Commercial Cluster Batch Schedulers, Technical Report 2005-13, Hebrew University, 2005
- Stewart G.: Grid data management. Reliable file transfer services performance. In: CHEP 2006
- 13. Videolan, http://www.videolan.org