# Exact approaches for static data segment allocation problem in an information network

Goutam Sen [a], Mohan Krishnamoorthy [a,b,*], Narayan Rangaraj [c], Vishnu Narayanan [c]

[a] *IITB-Monash Research Academy, IIT Bombay, Powai, Mumbai 400076, India*
[b] *Department of Mechanical Engineering, Monash University, Clayton, VIC 3800, Australia*
[c] *Industrial Engineering and Operations Research, IIT Bombay, Powai, Mumbai 400076, India*

## ARTICLE INFO

## ABSTRACT

In a large distributed database, data are geographically distributed across several separate servers (or data centers). This helps in distributing load in the access network. It also helps to serve data locally where it is required. There are various approaches based on the granularity of data for efficient data distribution in a communication network. The *file allocation problem* (FAP) locates files to servers, the *segment allocation problem* (SAP) locates database segments, and the *mirror location problem* (MLP) locates replicas of the entire database. The placement of such data to multiple servers can be modeled as an optimization problem. The major decisions influencing optimization involves the location of servers, allocation of content and assignment of users. In this paper, we study the segment allocation problem (SAP), which is also known as the partial mirroring problem. This approach is more tractable than the file allocation problem in realistic cases and also eliminates the overhead of (constant) update costs that is incurred in the mirror location problem. Our contribution is two-fold: Firstly, earlier works on SAP assume pre-defined segments. We build a data partitioning method using well-known facility location models. We quantify the performance of the partitioning method. We show that the method partitions the database within a reasonable limit of error. Secondly, we introduce a new model for the segment allocation problem in which the segments are completely connected to each other by high-bandwidth links and contains a cost benefit for inter-segment traffic flows. We formulate this problem as an MILP and build exact solution approaches to solve large scale problems. We demonstrate some structural properties of the problem that make it solvable, using a Benders decomposition algorithm. Computational results validate the superiority of the decomposition approach.

© 2014 Published by Elsevier Ltd.

## 1. Introduction

Internet usage has grown dramatically in recent times and the emergence of mobile technology is indicative of the fact that this phenomenon is not going to decrease. The result is a huge increase in data traffic in communication networks. Cisco Visual Networking Index (VNI), an IP traffic tracking and forecasting initiative, predicts an exponential increase in traffic at a compound annual growth rate (CAGR) of 23 percent from 2012 to 2017 [17]. Moreover, storage costs have been falling consistently over the last decade, due to advances in storage technologies, and are projected to fall further [31]. However, an expansion of network bandwidth to satisfy an ever-increasing demand is limited by budget availability. The International Data Center (IDC) projects a huge shortfall between the global Internet traffic growth (at 32% for the period 2010–2015) and growth in budgets for network equipment (at 8%) [35]. As a result, users could face high congestion in the network leading to high latency or download times, and service providers could incur high transmission costs for routing data. It is, therefore, necessary to study the problem of locating and accessing data in a cost-effective manner.

One popular approach for this problem is to establish multiple servers across the geography. A query is likely to be served locally, thereby reducing the overall congestion in the network. However, the locations of these servers and their contents could affect the transmission costs significantly. A small change in the placement scheme could lead to significant cost differences. Therefore, instead of using ad-hoc strategies, it is essential to model this problem and study it thoroughly.

We introduce the problem and two modeling approaches – an integrated approach and a 2-phased approach – in Section 2. We review the literature in Section 3 and place it appropriately in the

* Corresponding author at: IITB-Monash Research Academy, IIT Bombay, Powai, Mumbai 400076, India.

*E-mail addresses:* goutam.sen@iitb.ac.in (G. Sen),
mohank@iitbmonash.org (M. Krishnamoorthy),
narayan.rangaraj@iitb.ac.in (N. Rangaraj), vishnu@iitb.ac.in (V. Narayanan).

context of similar problems that have been studied in the past. Section 4 introduces a formulation for our integrated approach. The 2-phased approach is studied in detail in Section 5. A Benders decomposition algorithm for the 2-phased approach is developed in Section 6. Finally, we present detailed computational results in Section 7, and conclude the paper by summarizing our contributions and pointing out several directions of research in Section 8.

## 2. Problem description and modeling approach

Consider a very large database of files that need to be accessed regularly by users who are distributed in a network. Our objective is to locate these files in a cost-efficient manner in the network so that access to these files is easier. For example, consider a large online movie database that has a geographically distributed set of users. Netflix is one such provider. MovieLens is a movie recommender website (for data visualization, see [52]) which contains users and the movies that they recommend, review or rate. These user ratings can be easily interpreted as a proxy for access statistics to provide useful insights into user-file interaction in an information network. We could use this data to then decide where and how a provider of online movies must locate their movie files, given the likely locations of the users of this service.

Approaches to this problem in the literature are based on different levels of granularity. Fine-grained approaches consider the smallest units of allocations (the files) while coarse-grained approaches place replicas of the entire database. The *segment allocation problem* (SAP) is one way to arrive at an intermediate solution. The SAP distributes parts (or segments) of the database across the network.

We present the SAP and then define and study a problem, which we have termed the *static data segment location problem in information networks* (SDSLPIN). We consider segment allocation in a *static* environment, where access patterns do not change over time. Although the segment allocation problem has been modeled in the literature, this has been carried out under the assumption that segments (or groups of files) themselves are pre-defined. However, in realistic problems, we do not usually have pre-specified segments. So, SDSLPIN consists of two problems: the *data partitioning problem* (DPP) and the *segment allocation problem* (SAP). The DPP is a decision problem that decides how the files are clustered to form segments. The SAP, as we define it, is a decision problem that identifies potential sites, locates the segments, and routes the requests/queries in such a way that cost is minimized. DPP and SAP can be solved in an integrated manner as a joint decision problem; they can be solved sequentially too.

In an integrated approach, we consider files as units of allocation in a cost-optimization model. The number of servers are pre-specified. The optimization problem solves three problems: locations of servers, file placement and user assignment. The model output is an allocation schema of files to servers. In this process, the model generates optimal file clusters and their placement in the network. Files and users are usually large in number. Hence, such an approach may not be tractable by exact models for realistic data instances. Nevertheless, building an integrated model is important from a theoretical point of view. In addition, it might be important to study the integrated model, especially for small instances, to ascertain the performance of any sequential (or phased) approach.

In a phased approach, we solve two problems, DPP and SAP, sequentially. In the first phase, the DPP partitions the database using a partitioning method that clusters files into a pre-specified number of segments. In the second phase, the segments are given as inputs to the SAP. This results in significant reduction in the dimension of the problem compared to the integrated model. This
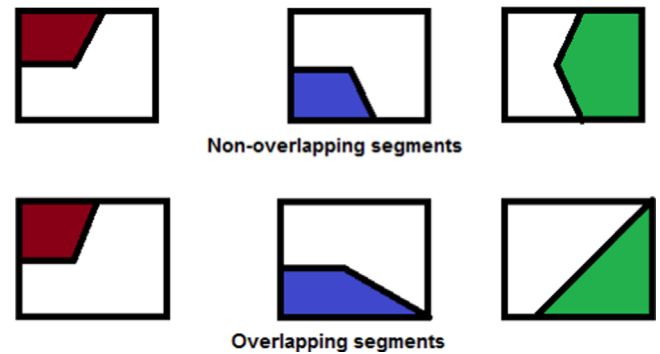


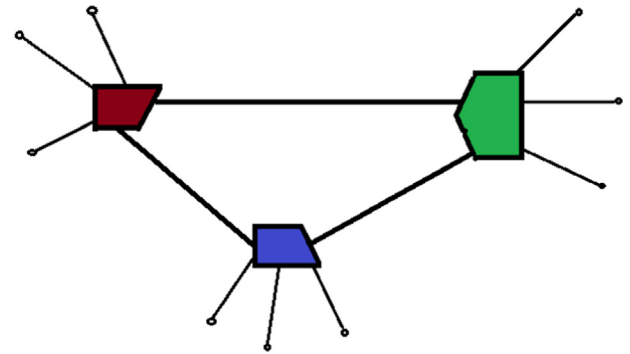**Fig. 1.** Database partitioning.



**Fig. 2.** Schematic diagram for SAP.

enables us to solve much larger instances, since the SAP model optimally places segments, and not the files (note that a file in a segment might be more profitably located in another segment). We are essentially treating two problems (partitioning and location) independently. Our contribution is in suggesting a reasonably strong partitioning approach to form the segments. Other contributions to the literature are in the form of developing a new model for SAP, formulating the problem, and building exact solution strategies to enhance the scalability of the approach.

Many models in data location literature are built on the *p*-median problem and its extensions. Our model is inspired by the *hub location problem* (see [23,10]) in physical logistics (examples are postal and airline networks). In this model, the user nodes, which send and receive information, communicate with each other via one or more hubs. These hubs are established to consolidate flows from the user nodes and to transfer them in bulk to another hub, where flow is de-consolidated and distributed to the destination nodes. The backbone network connecting the hubs is built with high-bandwidth links, and cost benefit accrues by considering economies of scale for the high volumes of inter-hub traffic flows.

We model our problem by treating the segments as data hubs connected via a high-bandwidth, fully connected backbone network. A user has 'direct access' to all the files in the segment that they are *uniquely* assigned to. We call this direct access a *hit*. If the requested content is absent in the segment that a user is allocated to, we term this a *miss*. The 'miss' content is accessed from the segment that is located at a different node that hosts it (see Figs. 1 and 2). This indirect access (or a miss) translates into an inter-hub flow of traffic via a high-bandwidth link adding a cost advantage to our model.

We make the following assumptions about our problem:

1. We consider that each file has a single copy and appears in only one segment. This results in a set of segments in which the contents do not overlap. The variant of this, with overlapping