

Special Issue of Journal of Computer Animation and Virtual World, (to appear), 2005.

Optimal pipeline decomposition and adaptive network mapping to support remote visualization

Mengxia Zhu[†] Qishi Wu[‡] Nageswara S.V. Rao[‡]

S.Sitharama Iyengar[†]

[†]Louisiana State University

Baton Rouge, LA 70802

Tel. (+1)225 578 1495 Fax. (+1)225 578 1465

email: {mzhu, iyengar}@bit.csc.lsu.edu

[‡]Oak Ridge National Lab

Oak Ridge, TN 37830

Tel. (+1)865 576 7210 Fax. (+1)865 574 0405

email: {wuqn, raons}@ornl.gov

Abstract

This paper discusses algorithmic and implementation aspects of an automatically configurable remote visualization system, which optimally decomposes and adaptively maps the visualization pipeline to a wide-area or dedicated network. The first node typically generates or stores raw datasets and a remote client resides on the last node equipped with a display device ranging from a personal desktop to a powerwall. Intermediate nodes can be located anywhere on the network and often include workstations, clusters, or rendering engines. We employ a regression model-based network daemon to estimate the effective bandwidth of a transport path using active traffic measurement. Based on the link measurements, node characteristics, and module properties, we strategically organize visualization pipeline modules into groups and dynamically assign them to various network nodes for minimal total delay or maximal frame rate. We propose polynomial-time algorithms using the dynamic programming method to compute the optimal solutions for the problems of pipeline decomposition and network mapping under different constraint conditions. A remote visualization system is implemented and deployed at several geographically distributed nodes for experimental testing. The proposed decomposition and mapping scheme is generic and can be applied to other distributed applications whose computing components form a linear arrangement.

Keywords: remote visualization, visualization pipeline, bandwidth measurement, network mapping

Introduction

A remote visualization system can potentially enable an end user equipped with a simple display device and network access to visualize large volumes of scientific data residing or generated at remote sites. A typical version of the remote visualization system consists of a remote data source acting as a server, a local rendering/display terminal acting as a client, zero or more intermediate hosts, and a network connecting them all together. The performance of such systems critically relies on how efficiently the visualization pipeline is *decomposed* or *partitioned* and mapped to the network nodes. In this paper, we address both analytical and implementation aspects of realizing the optimal visualization pipeline decomposition and adaptive network mapping under various circumstances.

Many existing remote visualization systems [1, 2, 3] employ a predetermined partition of the visualization pipeline and typically send fixed-type data streams such as raw data, geometric primitives, or framebuffer (FB) to remote client nodes. While such schemes are common, they are not always optimal for high performance visualizations that typically deal with large data sets of complex structures. Particularly over wide-area connections, this problem is further compounded by limited network bandwidths and time-varying dynamics of network conditions. A considerable amount of research efforts have been made to implement flexible remote visualization systems that distribute visualization modules across network nodes. Bowman *et al* [4] proposed a framework to predict the processing times of visualization modules using linear models and the network bandwidth using Network

Weather Service [5]; these predictions are used to obtain a suitable mapping of the visualization pipeline. Luke *et al* [6] constructed Semotus Visum, a flexible remote visualization framework capable of multiple partition scenarios, which is tested and evaluated on a local network. In these works, the mapping of visualization modules to network nodes is accomplished by empirical testing and manual configuration. Furthermore, their visualization systems are limited to a traditional client and server mode with no intermediate nodes considered.

In this paper, we analytically formulate a problem of optimizing the total delay or frame rate of a visualization pipeline by considering the times for transmission, computation and rendering. The resulted computational model enables us to analyze the algorithmic aspects such as complexity and optimality of mapping the visualization pipeline onto a given network. The first node serves as a data source with the capability of performing other tasks such as filtering, geometry generation, or rendering. The last node displays final images on the screen but also contains other visualization components; it could be equipped with a simple terminal or a more sophisticated device such as a powerwall or a tiled display. The intermediate nodes could be workstations, computational clusters, or custom rendering engines that are specially utilized to perform computation-intensive visualization subtasks. We present several polynomial-time solutions with increasingly stringent conditions based on dynamic programming to compute optimal mappings with minimum total delay or maximum frame rate. The time complexities of our methods range from $O(n \times k)$ to $O(n \times |E|)$ under different conditions, where $n + 1$ is the number of visualization modules in a pipeline,

$k + 1$ is the number of network nodes arranged linearly in a predetermined order, and $|E|$ is the number of edges in an arbitrary wide-area or dedicated network.

The rest of the paper is organized as follows. In Section , we describe a generic visualization pipeline and a framework for the design of remote visualization system, which form the basis for our analytical model. In Section , we first describe a bandwidth measurement method based on [7], and then present the optimal decomposition and mapping schemes using dynamic programming-based algorithms for total delay minimization or frame rate maximization under various constraint conditions. Implementation details and test results are provided in Section . We finally conclude our work in Section .

Visualization pipeline and analytical model

Visualization pipeline

Large volumes of data generated in scientific or medical applications need to be appropriately retrieved and mapped onto a 2D display device to be “visualized” by human operators. This visualization process involves several steps that form the so-called *visualization pipeline* or *visualization network* [8].

Fig. 1 shows a high-level general abstraction of a visualization pipeline along with the flow of data produced at each pipeline module. In many scientific applications, the raw data usually takes a multivariate format and is organized in structures such as CDF, NetCDF,

and HDF [9, 10, 11]. The *filtering* module extracts the information of interest from the raw data and performs necessary preprocessing to improve processing efficiency and save communication resources as well. The transformation module typically uses a surface fitting technique (such as isosurface extraction) to derive 3D geometries (such as polygons), or uses transfer functions to perform color and opacity classifications for each voxel based on its attribute. Rendering module converts the transformed geometric or composite volumetric data in 3D view coordinates to a pixel-based image in 2D screen coordinates. In most of the existing graphics systems, the visual properties of a raster image such as color and opacity is stored and carried in a frame buffer for final display on an end device.

It is worth pointing out that the visualization pipelines in real applications may significantly vary due to the disparate implementation procedures and the use of different visualization techniques; for instance, the rendering module itself might involve several stages of vertex transformations in the OpenGL environment.

Framework for design of remote visualization system

The block diagram in Fig. 2 illustrates a baseline framework for the design of a remote visualization system that employs pipeline *partitioning* or *decomposition* and automatic *mapping* to network nodes. Since many scientific applications generate terabyte or even petabyte datasets, it is currently very difficult to execute all visualization modules on a single desktop computer. Hence we first exploit the parallelism in the modules and the computational

capability provided by the nodes with High-Performance Computing (HPC) resources, such as clusters. The connection bandwidths over the underlying transport network are estimated using active traffic measurements. Based on the computing node capabilities, bandwidth estimates, sizes of datasets to be processed, and computational complexities of visualization modules, the visualization pipeline is decomposed into *groups*, which are then mapped one-to-one to the computing nodes distributed over a transport network.

In a practical implementation, we generally start with the computing node deployment and network topology construction. For example, the site of data source is *a priori* known, and the location of a remote client is determined whenever a “visualization” connection to the server is initiated. The information collected on networking and computing resources can be used to select optional intermediate nodes with specific visualization and/or computing capabilities. In this paper, we consider several problems of pipeline decomposition and network mapping with different optimization goals under various constraint conditions, as tabulated in Fig. 3.

Analytical model

We now describe an analytical model for the general problem of visualization pipeline decomposition and its network mapping. As shown in Fig. 4, the visualization pipeline consists of $n + 1$ sequential modules, $M_1, M_2, \dots, M_{u-1}, M_u, \dots, M_{v-1}, \dots, M_w, \dots, M_{x-1}, M_x, \dots, M_{n+1}$, where M_1 is a data source. Module $M_j, j = 2, \dots, n + 1$ performs a computational task of

complexity c_j on data of size m_{j-1} received from module M_{j-1} and generates data of size m_j , which is then sent over the network link to module M_{j+1} for further processing. An underlying transport network consists of $k + 1$ geographically distributed computing nodes denoted by $v_1, v_2, \dots, v_k, v_{k+1}$. Node $v_i, i = 1, 2, \dots, k, k + 1$ has a normalized computing power p_i^1 and is connected to its neighbor node $v_j, j = 1, 2, \dots, k, k + 1, j \neq i$ with a network link $L_{i,j}$ of bandwidth $b_{i,j}$ and minimum link delay $d_{i,j}$. The minimum link delay is mostly contributed by the link propagation and queuing delay, and is in general much smaller than the bandwidth-constrained delay of transmitting a large message of size m given by $m/b_{i,j}$. The transport network is represented by a graph $G = (V, E), |V| = k + 1$, where V denotes the set of nodes (vertices) and E denotes the set of links (edges). The transport network may or may not be a complete graph, depending on whether the node deployment environment is the Internet or a dedicated network.

We consider a path P of q nodes from a source node v_s to a destination node v_d in the transport network, where $q \in [2, \min(k + 1, n + 1)]$ and path P consists of nodes $v_{P[1]} = v_s, v_{P[2]}, \dots, v_{P[q-1]}, v_{P[q]} = v_d$. The visualization pipeline is decomposed into q visualization groups denoted by $g_1, g_2, \dots, g_{q-1}, g_q$, which are mapped one-to-one onto the q nodes on transport path P . The data flow between two adjacent groups is the one produced by the last module in the upstream group; for example in Fig. 4, we have

¹For simplicity, we use a normalized quantity to reflect a node's overall computing power without detailing its memory size, processor speed, and presence of co-processors, which may result in different performances for both numeric and graphics computations.

$m(g_1) = m_{u-1}, m(g_2) = m_{v-1}, \dots, m(g_{q-1}) = m_{x-1}$. The client residing on the last node v_d sends control messages such as simulation parameters, filter types, visualization modes, and view parameters to one or more preceding visualization groups to support interactive operations. However, since the size of control messages is typically in the order of bytes or kilobytes, which is often much smaller than the size of visualization data, its transport time is assumed to be negligible.

- **Minimal total delay**

The most important requirement in many collaborative visualization applications is the interactiveness of the system. For applications in which the remote visualization is employed, we consider the optimization problem minimizing the total end-to-end delay:

$$\begin{aligned}
T_{total}(\text{Path } P \text{ of } q \text{ nodes}) &= T_{computing} + T_{transport} \\
&= \sum_{i=1}^q T_{g_i} + \sum_{i=1}^{q-1} T_{L_{P[i],P[i+1]}} \\
&= \sum_{i=1}^q \left(\frac{1}{b_{P[i]}} \sum_{j \in g_i, j \geq 2} (c_j m_{j-1}) \right) + \sum_{i=1}^{q-1} \left(\frac{m(g_i)}{b_{P[i],P[i+1]}} \right).
\end{aligned} \tag{1}$$

Our goal is to minimize the total time incurred on the forward links from the source node to the destination node to achieve the fastest response. Note that in Eq. 1, we assume that the transport time between modules within one group on the same computing node is negligible.

- **Maximal frame rate**

Our goal is to maximize the frame rate by minimizing the time incurred on a bottleneck

link/node for applications with streaming data, given as in Eq. 2.

$$\begin{aligned}
& T_{bottleneck}(\text{Path } P \text{ of } q \text{ nodes}) = \\
& \max_{\substack{\text{Path } P \text{ of } q \text{ nodes} \\ i=1,2,\dots,q-1}} \left\{ \begin{array}{l} T_{computing}(g_i), \\ T_{transport}(L_{P[i],P[i+1]}), \\ T_{computing}(g_q) \end{array} \right\} \\
& = \max_{\substack{\text{Path } P \text{ of } q \text{ nodes} \\ i=1,2,\dots,q-1}} \left\{ \begin{array}{l} \frac{1}{p_{P[i]}} \sum_{j \in g_i \text{ and } j \geq 2} (c_j m_{j-1}), \\ \frac{m(g_i)}{b_{P[i],P[i+1]}}, \\ \frac{1}{p_{P[q]}} \sum_{j \in g_q \text{ and } j \geq 2} (c_j m_{j-1}) \end{array} \right\} \tag{2}
\end{aligned}$$

Design of Remote Visualization System

In this section, we present the technical solutions for the design of a remote visualization system. In particular, we introduce a linear regression model to estimate link bandwidth and minimum delay using active traffic measurement, and propose several approaches based on dynamic programming to solve the visualization pipeline partitioning and mapping problems under various conditions.

Bandwidth measurement

We present a linear regression model to estimate the bandwidth of a transport path using active traffic measurement based on [7]. Due to complex traffic distribution over wide-area networks and the non-linear nature of transport protocol dynamics (in particular TCP), the

throughput achieved in actual message transfers is typically different from both the link and available bandwidths, and typically contains a random component. We consider the *effective path bandwidth* (EPB) as the throughput achieved by a flow using a given transport module under certain cross traffic conditions. The notion of effective path bandwidth is specific to the transport protocol employed by the transport daemon. The active measurement technique we apply here is to estimate the effective path bandwidth and minimum delay for each virtual link. Note that a virtual link in the overlay network of transport daemons may correspond to a multi-hop data path in wide-area networks, which usually consists of multiple underlying physical links from different networks.

There are three main types of delays involved in the message transmission over computer networks, namely, link propagation delay d_p imposed at the physical layer level, equipment-associated delay d_q mostly incurred by processing and buffering at the hosts and routers, and bandwidth-constrained delay d_{BW} . The delay d_q often experiences a high level of randomness in the presence of time-varying cross traffic and host loads. Also, since the transport protocol reacts to the competing traffic on the links, the delay d_{BW} may also exhibit randomness particularly over congested wide-area connections. We use Eq. (3) to measure the end-to-end delay in transmitting a message of size r on a path P with l links:

$$d(P, r) = d_{BW}(P, r) + \sum_{i=1}^l (d_{p,i}(P) + d_{q,i}(P, r)). \quad (3)$$

Due to the large size of data transfer in high-performance visualization applications, only the first term of Eq. (3) is significant and therefore the delay $d(P, r)$ of transmitting a mes-

sage of size r along path P can be approximated by a linear model: $d(P,r) \approx r/EPB(P)$. The active measurement technique generates a set of test messages of various sizes, sends them to a destination node through a transport channel such as a TCP flow, and measures the end-to-end delays, on which we apply a linear regression to estimate the EPB. The delay measurements between Louisiana State University (LSU) and Oak Ridge National Laboratory (ORNL) as well as its corresponding linear regression are illustrated in Fig. 5, from which we estimate the EPB on this path to be about 1.0 Mbps. We emphasize that the same transport method used for collecting measurements will be used by the visualization modules. If measurements are collected by tools such as Iperf, or NWS [5], they must be appropriately translated into the effective bandwidth seen by the visualization pipeline. If the bandwidth measurements do not use the same transport module or the same transport parameters as the visualization modules, it is not sufficient to simply “plug-in” their estimates in Eqs. (1) and (2).

Pipeline partition and network mapping

We propose approaches based on dynamic programming to solve the visualization pipeline partition and mapping problems with the minimal total delay or maximal frame rate. Intuitively, a more general version of these problems is quite similar to the classical graph clustering problem, which is NP-complete and solved by approximate solutions using the branch-and-bound technique but in exponential time [12]. However, with the understand-

ing on the linear arrangement of the visualization modules, we are able to achieve optimal solutions to these problems in polynomial time using the dynamic programming method.

Minimal total delay

We consider Problems 2, 3, and 4 in Fig. 3 for total delay minimization under different stringent levels of constraint conditions. Problem 2 decomposes the modules of the visualization pipeline into a fixed number of groups that are then mapped *onto* a pre-selected linear arrangement of network nodes. Problem 3 maps all the modules in the visualization pipeline *one-to-one* to different nodes in an arbitrary network, i.e. a mapped node only runs one module. Problem 4 addresses a general case where the pipeline can be decomposed into any number of groups and the available computing nodes form an arbitrary network.

- **Surjective mapping with linearly arranged network nodes**

In this scenario, all $k+1$ nodes in the computer network are pre-selected and linearly arranged for mapping. We aim to decompose $n+1$ modules into $k+1$ groups and map them onto $k+1$ nodes to achieve the minimal total delay of the visualization pipeline. There are n messages flowing between $n+1$ visualization modules with sizes $m_j, j = 1, 2, \dots, n$, and k network links with bandwidths $b_i, i = 1, 2, \dots, k$ connecting $k+1$ computing nodes, each of which has processing speed $p_s, s = 1, 2, \dots, k+1$.

Let $T(i, j)$ denote the minimal total delay with the first j messages (namely the first $j+1$ visualization modules) mapped onto the first i network links. The dynamic programming

method uses the following recursion to compute $T(i, j)$:

$$T(i, j) = \min_{i=1 \text{ to } k, j=2 \text{ to } n, i \leq j} \left\{ \begin{array}{l} T(i-1, j-1) + m_j/b_i + c_{j+1}m_j/p_{i+1} \\ T(i, j-1) + c_{j+1}m_j/p_{i+1} \end{array} \right\}, \quad (4)$$

where the base conditions are computed as $T(0, t) = \sum_{i=1}^t c_{i+1}m_i/p_1, t = 1, 2, \dots, n$ in the first row and $T(t, t) = \sum_{i=1}^t m_i/b_i + \sum_{i=1}^t c_{i+1}m_i/p_{i+1}, t = 1, 2, \dots, k$ on the diagonal line. Note that the first module is assumed to be a data source that does not introduce computing time. The complexity of this algorithm is in the order of $O(n \times k)$.

We now establish the correctness of Eq. (4). At each step of the recursion, $T(i, j)$ takes the minimal of two subcases as illustrated in Fig. 6. In the first subcase, the last message m_j is mapped to the last link b_i so that the transport time on this mapped link together with the computing time of the last module on the last node is added to $T(i-1, j-1)$, which is the sub-problem of size $i-1$ and $j-1$. In the second subcase, we do not map the last message m_j on any link, which means that the last two modules are both executed on the last node, and then the problem directly reduces to the sub-problem of size i and $j-1$ with the addition of the computing time of the last module on the last node.

Fig. 7 shows the dynamics of 2D matrix construction for computing $T(i, j)$. The entries of the diagonal line and first row are entered beforehand as the base conditions. In order to fill up the whole upper triangle, the calculation sweeps across the matrix from left to right first and from top to bottom afterwards. For example, let us consider the back tracing path of $T(3, 4)$. We first visit $T(2, 3)$ and $T(3, 3)$, of which $T(3, 3)$ is already known from the

base case and $T(2,3)$ is calculated from $T(2,2)$ and $T(1,2)$, which is in turn calculated from $T(1,1)$ and $T(0,1)$. Once the base cases are reached, the tracing process will stop and bounce back to its starting point. On the diagonal line with an equal number of links and messages, we simply map all messages one-to-one onto links in a linear order. For each comparison step, the mapping scheme matrix either inherits the mapping scheme from that of $T(i, j - 1)$ by adding module M_{j+1} to the last group, or appends a separate group with module M_{j+1} to the mapping scheme of $T(i - 1, j - 1)$. Note that an additional matrix is needed to record the mapping scheme for the winner of each comparison along with the computation.

- **Injective mapping with arbitrary network**

In this scenario, all $n+1$ visualization modules are individually mapped to a set of $n+1$ nodes in an arbitrary network of $k+1$ nodes with the first module mapped to v_s , and the last module mapped to v_d . That is, we consider a mapping path P of $q = n + 1$ nodes and n hops. The mapping of intermediate modules to network nodes needs to be decided such that the total delay time is minimized. Since no grouping of visualization modules is allowed, the pipeline partition is actually predetermined.

With appropriate adjustments made on the evaluation of the objective function, this problem turns into an n -hops shortest path problem. Let $T^l(v_i)$ denote the minimal delay of a path with l hops from the source node v_s to the node v_i under consideration. We have the

following recursion leading to $T^n(v_d)$:

$$T^l(v_i) = \left\{ \min_{u \in \text{adj}(v_i)} \left(T^{l-1}(u) + c_{l+1}m_l/p_{v_i} + m_l/b_{u,v_i} \right) \right\}. \quad (5)$$

As Fig. 8 shows, this recursion follows from the observation that the minimal delay to v_i with l hops is the minimum of the delays to its neighbor with $l-1$ hops plus the cost incurred by

that link. The base conditions are computed as:
$$T^1(v_i) = \begin{cases} c_2m_1/p_{v_i} + m_1/b_{v_s,v_i}, & \forall e_{v_s,v_i} \in E \\ \infty & , \text{ otherwise} \end{cases}.$$

The complexity of this algorithm is $O(n \times |E|)$.

- **A general mapping scheme**

Now we consider a general mapping scheme that removes the previous constraints imposed on the visualization pipeline and computer network. The objective of pipeline partition and network mapping is to minimize the total end-to-end delay by decomposing the visualization pipeline into any number of groups and mapping them to an arbitrary network. This scenario imposes the least restrictions and is the most general case for our problem. Let $T^j(v_i)$ denote the minimal total delay with the first j messages (namely the first $j+1$ visualization modules) mapped to a path from the source node v_s to node v_i under consideration in the computer network. Then, we have the following recursion leading to $T^n(v_d)$ [20]:

$$T^j(v_i) = \min_{j=2 \text{ to } n, v_i \in V} \left\{ \begin{array}{l} T^{j-1}(v_i) + c_{j+1}m_j/p_{v_i}, \\ \min_{u \in \text{adj}(v_i)} \left(T^{j-1}(u) + c_{j+1}m_j/p_{v_i} + m_j/b_{u,v_i} \right) \end{array} \right\}, \quad (6)$$

with the base conditions computed as:
$$T^1(v_i) = \begin{cases} c_2 m_1 / p_{v_i} + m_1 / b_{v_s, v_i}, & \forall e_{v_s, v_i} \in E \\ \infty & , \text{ otherwise} \end{cases}$$

on the first column and $T^t(v_s) = \sum_{i=1}^t (c_{i+1} m_i / p_{v_s})$, $t = 1, 2, \dots, n$ on the first row in the 2D matrix of dynamic programming as shown in Fig. 9.

In Eq. (6), at each step of the recursion, $T^j(v_i)$ takes the minimal of two subcases. In the first subcase, we do not map the last message m_j to any network link; instead we directly place the last module M_{j+1} at node v_i itself. Therefore we only need to add the computing time of M_{j+1} on node v_i to $T^{j-1}(v_i)$, which is a sub-problem of node v_i of size $j-1$. This subcase is represented by the direct inheritance link from its left neighbor element in the 2D matrix. In the second subcase, the last message m_j is mapped to one of the incident network links from its neighbor nodes to node v_i . The set of neighbor nodes of node v_i is enclosed in a cloudy area in Fig. 9. We calculate the total delay for each mapping of an incident link of node v_i and choose the one with the minimum delay, which is further compared with the one calculated in the first subcase. For each comparison step, the mapping scheme of $T^j(v_i)$ is obtained as follows: we either directly inherit the mapping scheme of $T^{j-1}(v_i)$ by simply adding module M_{j+1} to the last group, or create a separate group for module M_{j+1} and append it to the mapping scheme $T^{j-1}(u)$ of the neighbor nodes $u \in adj(v_i)$ of node v_i . The complexity of this algorithm is $O(n \times |E|)$.

Maximal frame rate

In visualization applications producing streaming data such as animations with a number of time steps, visualization data of the same modality is continuously generated, manipulated, and rendered in a pipelined manner. The maximal frame rate that a pipelining system can achieve is limited by the slowest (bottleneck) transport link or computing node along the pipeline. Due to the similarity of the recursive processes in achieving optimization goals for surjective and injective mapping schemes, only the most general case is discussed here.

We adapt the above dynamic programming method to this problem by introducing necessary modifications. Let $1/T^j(v_i)$ denote the maximal frame rate with the first j messages (namely the first $j+1$ visualization modules) mapped to a path from source node v_s and node v_i in an arbitrary computer network. Let $S^j(v_i)$ represent the sum of the message sizes of all modules on node v_i with the first j messages mapped from node v_s to v_i . We have the following recursion leading to $T^n(v_d)$:

$$T^j(v_i) = \min_{j=2 \text{ to } n, v_i \in V} \left\{ \begin{array}{l} \max \left(T^{j-1}(v_i), (S^{j-1}(v_i) + c_{j+1}m_j)/p_{v_i} \right), \\ \min_{u \in \text{adj}(v_i)} \left(\max \left(T^{j-1}(u), c_{j+1}m_j/p_{v_i}, m_j/b_{u,v_i} \right) \right) \end{array} \right\} \quad (7)$$

with the base conditions computed as: $T^1(v_i) = \begin{cases} \max(c_2m_1/p_{v_i}, m_1/b_{v_s,v_i}), & \forall e_{v_s,v_i} \in E \\ \infty & , \text{ otherwise} \end{cases}$

and $T^t(v_s) = \sum_{i=1}^t (c_{i+1}m_i/p_{v_s})$, $t = 1, 2, \dots, n$.

Every transport link or computer node is a potential bottleneck and needs to be checked.

At each step of the recursion, the bottleneck times for all possible schemes are computed

and the scheme with the minimal bottleneck time is chosen as the optimal result for maximal frame rate. One inherent problem for the maximal frame rate algorithm is that since our algorithm doesn't exclude loops in network path selection, loops may compromise the optimality of frame rate for streaming data due to the non-exclusive computing node possession by different modules at the same time. If loops appear, optimality cannot be guaranteed for the maximal frame rate. However, such situations rarely occur.

It is worth pointing out that some additional constraints may arise in practical applications. For example, the maximal number of modules on a specific node may be limited, or some nodes are only capable of executing certain visualization modules. Such constraints can be conveniently handled by imposing feasibility checks at each step of the dynamic programming recursions described in Eqs. (4)-(7). The scenario with failed feasibility check is simply discarded. Note that a solution in this situation is not guaranteed, especially when all scenarios fail in the feasibility checks. This feasibility checking feature makes our system versatile in meeting diverse practical requirements.

Implementation and Case Study

Our remote visualization system is implemented in C++ on the Linux platform using GTK+ for client GUI design. We selected three Internet hosts, boba.sinrg.cs.utk.edu at University of Tennessee at Knoxville (UT), ozy4.csm.ornl.gov at ORNL, and robot.rrl.lsu.edu at LSU for system deployment, as shown in Fig. 10. The hosts ozy4 at ORNL and robot at LSU are

Linux workstations equipped with 3GHz CPU. The host boba deployed at UT is a cluster, on which we use four nodes for parallel visualization computations in our experiment.

Our system provides various visualization functionalities such as isosurface extraction, volume rendering, linear integral convolution (LIC), and streamline. Data transmission is currently carried out via TCP/IP sockets.

Experimental tests

Our distributed visualization system is a multi-threaded system and implemented in Java, Fortran, and C++ in Linux. The client GUI is based on GTK+ development. Our system consists of four virtual component nodes, i.e. client, central management (CM), data source (DS), and computing service (CS), which are connected together by network links to form a closed loop. A visualization loop always starts at a client that initiates a particular visualization task by sending a request (containing data file, variable name, visualization method, view parameters, etc.) to a designated CM node. CM determines the best system configuration to accomplish the visualization task. Based on the global knowledge of data source and system resource distributions within its service area, CM strategically partitions the pipeline of the selected visualization method into groups and selects an appropriate set of CS nodes to execute the visualization modules. The computation for pipeline partitioning and network mapping results in a visualization routing table, which is delivered sequentially over the rest of the loop to establish the communication channel for the requested visualization task.

Each CS node along the path receives intermediate results from upstream node, and executes the visualization modules designated by the routing table. Computation results are forwarded to the downstream node for further processing. This loop ends at client node for final image display.

Several experiments are conducted with ozy4 at ORNL as DS, boba cluster at UT as CS, and robot at LSU. We established three different routing tables for isosurface visualization of two raw datasets as shown in Fig. 11. Each column indicates corresponding computing node for that computing module.

In the first case, raw dataset was sent out by ozy4 node to client end. Robot at LSU received raw data performed the local isosurface extraction and rendering. Significant end-to-end delays were observed for this case due to large dataset transmission on a low bandwidth link. In the second case, ozy4 conducted isosurface extraction and off-screen rendering, and final image was shipped to robot for final display. End-to-end delays were reduced approximately by half of the first case. In the third case, ozy4 pushed raw data to boba for parallel isosurface extraction and off-screen rendering with 4 nodes. The resulted image was transmitted from boba to robot. Due to the speed up gained by parallel processing, end-to-end delays were only about one-third of the second case.

Our system is self recoverable in the sense that if any node experiences system error or is stuck in a certain computation task, a service failure signal will be triggered to clear up current network setup.

In dealing with large data sets, image is likely to be sent to the client especially when the

available bandwidth is limited. Such partition scheme can exploit the data locality, disk I/O bandwidth, memory, and parallel processing to accelerate the visualization pipeline. The size of image generated by the rendering process is mainly decided by the width and height of the display screen regardless of the dataset volume. In case of very low bandwidth links, the display window may need to be suitably scaled down to reduce the image size, and in addition, a lower level of image resolution may be chosen.

Conclusion and Future Plan

In this paper, we proposed a framework and a mathematical model for automatically mapping a visualization pipeline to computer networks. Our objectives are to minimize the total delay or maximize the frame rate of the visualization system. Dynamic programming-based approaches are proposed to compute optimal schemes for grouping and mapping the visualization modules onto a computer network, which could either be the Internet or a dedicated network.

It would be of our future interest to study various formulations of this class of problems from the viewpoint of computational criteria and practical implementation. The routing table is manually generated in our current implementation. In the future, we plan to strengthen our implementation to support automatic routing table construction after network condition monitoring daemon is embedded, and also deploy our system over dedicated networks, such as DOE UltraScience Net [13], for experimental testing.

As the dataset sizes reach terabytes, the transport delays within the pipeline could be prohibitively high even over dedicated high-throughput networks. One obvious solution is to speed up the transport process. But most existing remote visualization systems (including ours) use the default TCP for both data and control message transmission. Newer transport protocols based on stochastic approximation methods for throughput stabilization and maximization [14, 15] have been developed to overcome the limitations of default TCP or UDP in terms of throughput, stability and dynamics. We plan to incorporate these new transport methods in our remote visualization system at a later stage.

Acknowledgements

Authors thank Professor John Blondin of North Carolina State University for providing us the code for computing the hydrodynamics of supernova and also for providing us an access to their computational and networking facilities. This research is sponsored by the High Performance Networking Program of the Office of Science, U.S. Department of Energy, under Contract No. DE-AC05-00OR22725 with UT-Battelle, LLC, the Defense Advanced Projects Research Agency under MIPR No.~K153, and by National Science Foundation under Grants No. ANI-0229969 and No. ANI-335185.

References

- [1] Aspect. <http://www.aspect-sdm.org/>.
- [2] Computational engineering international. <http://www.ceintl.com/products/ensight.html>.
- [3] Paraview. <http://www.paraview.org/HTML/Index.html>.
- [4] I. Bowman, J. Shalf, K. Ma, and W. Bethel. Performance modeling for 3d visualization in a heterogeneous computing environment. In *Technical Report No. 2005-3, Department of Computer Science, University of California at Davis*, 2005.
- [5] Network weather service. <http://nws.cs.ucsb.edu>.
- [6] E.J. Luke and C.D. Hansen. Semotus visum: a flexible remote visualization framework. In *Proceedings of IEEE Visualization*, pages 61–68, 2002.
- [7] Q. Wu, N.S.V. Rao, and S.S. Iyengar. On transport daemons for small collaborative applications over wide-area networks. In *Proc. of the 24th IEEE International Performance Computing and Communications Conference*, Phoenix, Arizona, April 7-9, 2005.
- [8] W. Schroeder, K. Martin, and B. Lorensen. *The Visualization Toolkit (2nd Edition)*. Prentice Hall PTR, 1998.
- [9] Common data format. <http://nssdc.gsfc.nasa.gov/cdf>.

- [10] Network common data form. <http://my.unidata.ucar.edu/content/software/netcdf>.
- [11] Hierarchical data format. <http://hdf.ncsa.uiuc.edu>.
- [12] P. Fränti, Virmajoki O., and Kaukoranta T. Branch and bound technique for solving optimal clustering. In *Int. Conf. on Pattern Recognition*, pages 232–235, August 2002.
- [13] Doe ultrasciencenet. <http://www.csm.ornl.gov/ultranet>.
- [14] Q. Wu and N.S.V. Rao. A class of reliable udp-based transport protocols based on stochastic approximation. In *Proc. of the 24th IEEE INFOCOM*, Miami, Florida, March 13-17, 2005.
- [15] Q. Wu. *Control of transport dynamics in overlay networks*. PhD thesis, Dept of Computer Science, Louisiana State University, 2003.

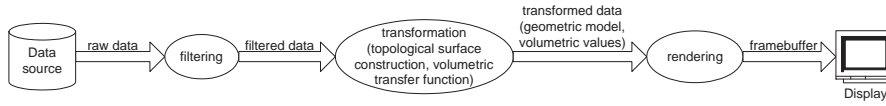


Figure 1: A general visualization pipeline.

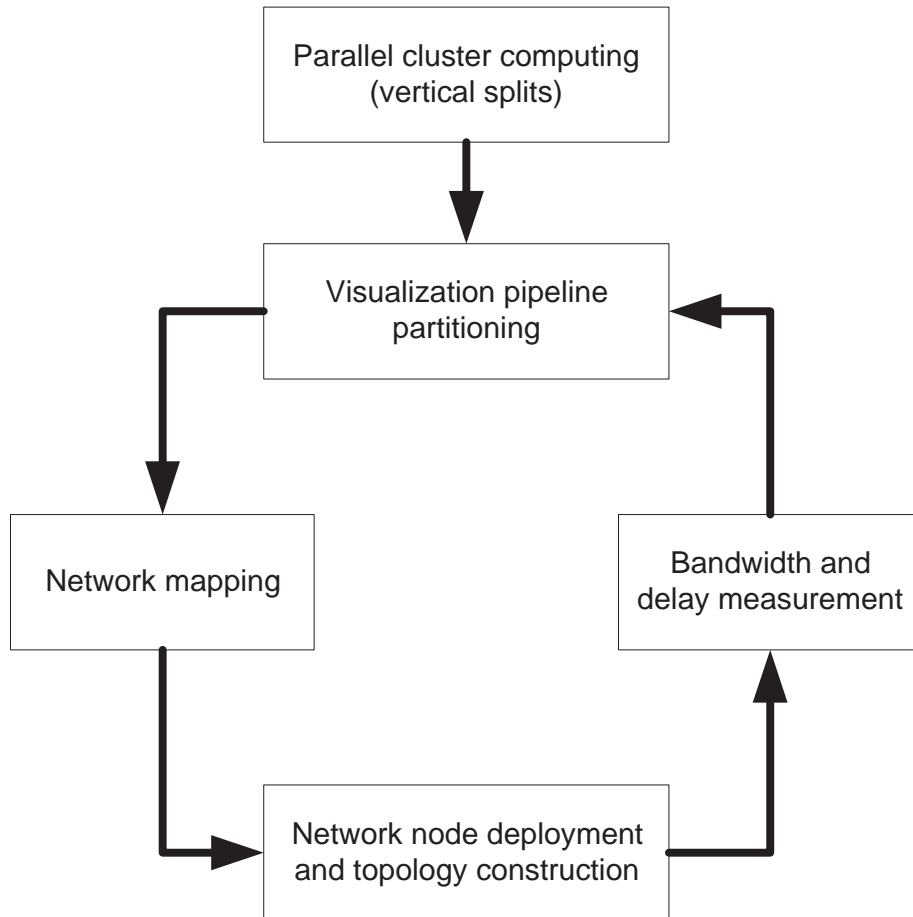


Figure 2: Framework for design of remote visualization system.

Problem No.	Mapping scheme	Module grouping in visualization pipeline	Node deployment in computer network	Optimization goals	Solutions
1	Bijjective mapping (one-to-one onto)	No	Linear arrangement of a predetermined number of nodes	N/A	Trivial
2	Surjective mapping (onto)	Yes			
3	Injective mapping (one-to-one)	No	Arbitrary network with undetermined number and order of nodes	Minimal total delay	Dynamic programming
4	General mapping for delay minimization	Yes			
5	General mapping for rate maximization	Yes			

Figure 3: Problem category of pipeline decomposition and network mapping.

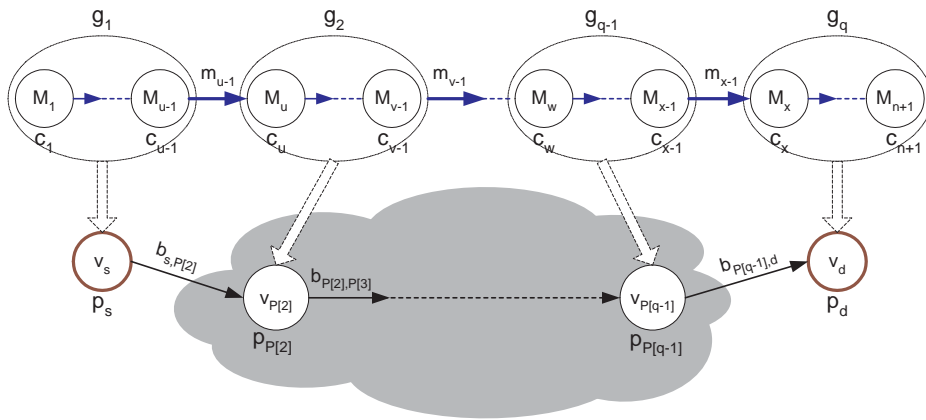


Figure 4: Mathematical model for pipeline partitioning and network mapping

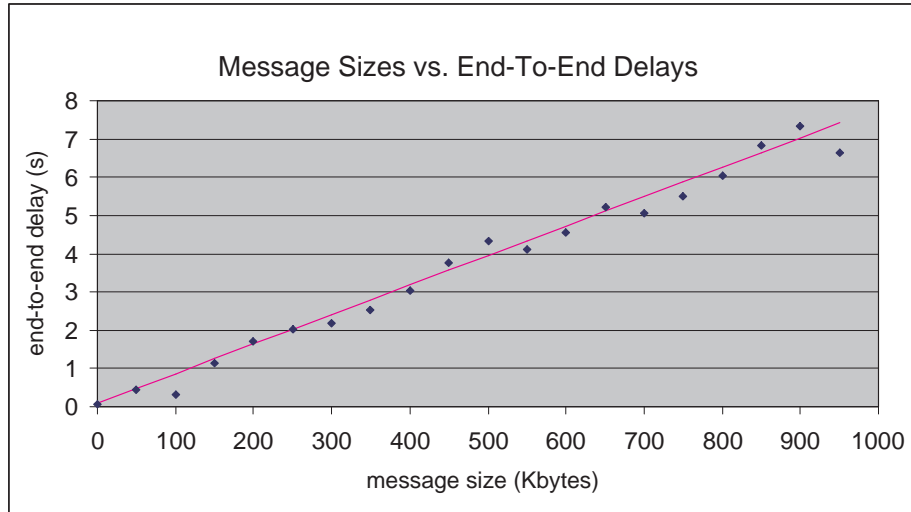


Figure 5: End-to-end message transmission delay measurements between LSU and ORNL.

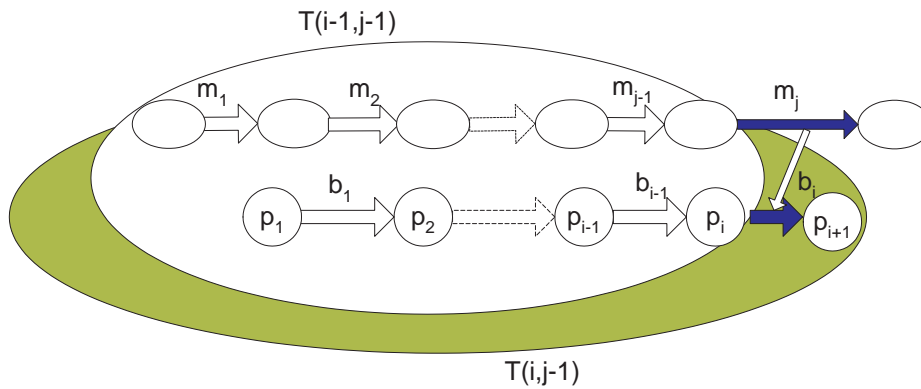


Figure 6: Surjective mapping for minimal delay time.

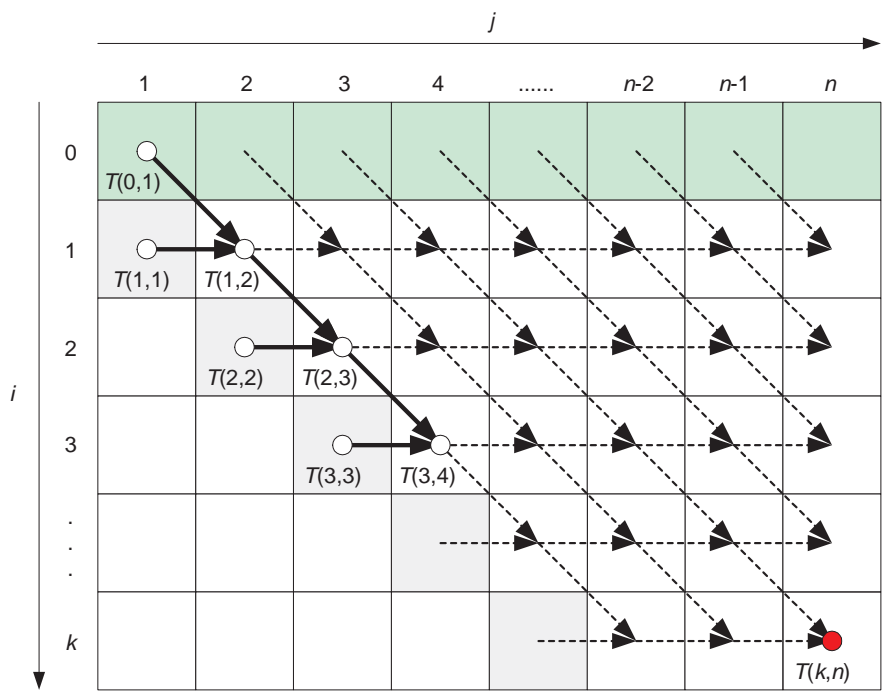


Figure 7: Construction of a 2D matrix of $T(i,j)$.

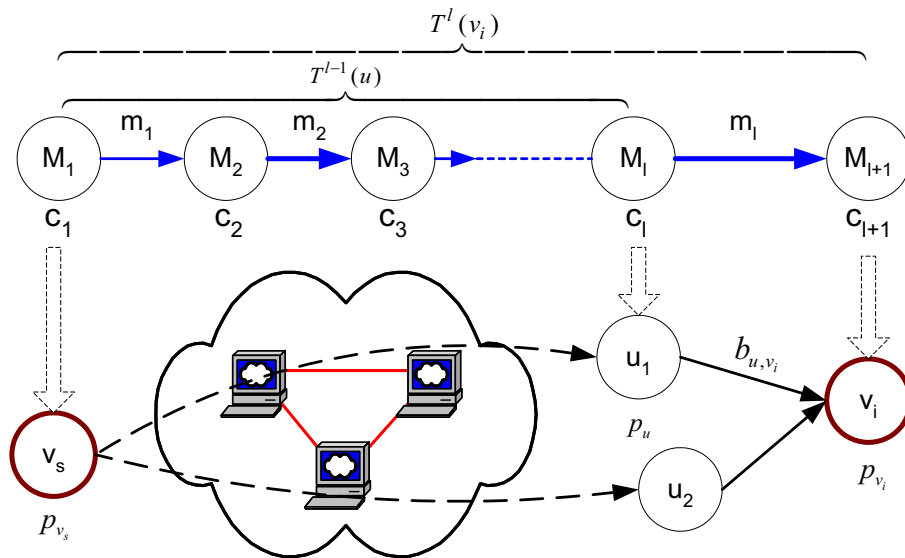


Figure 8: Illustration of n -hops shortest path algorithm.

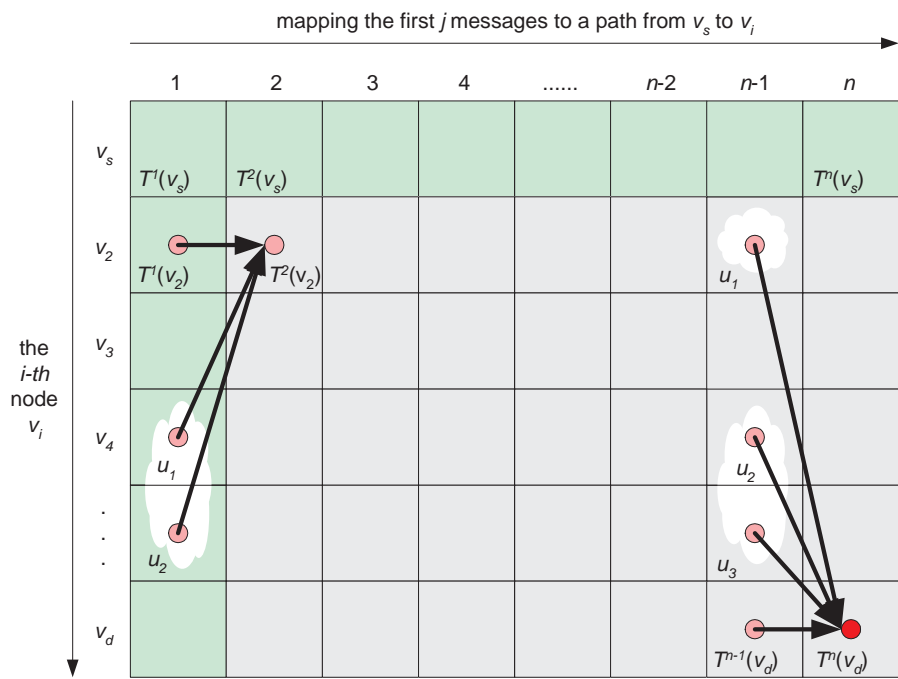


Figure 9: Construction of 2D matrix in dynamic programming for the most general case.

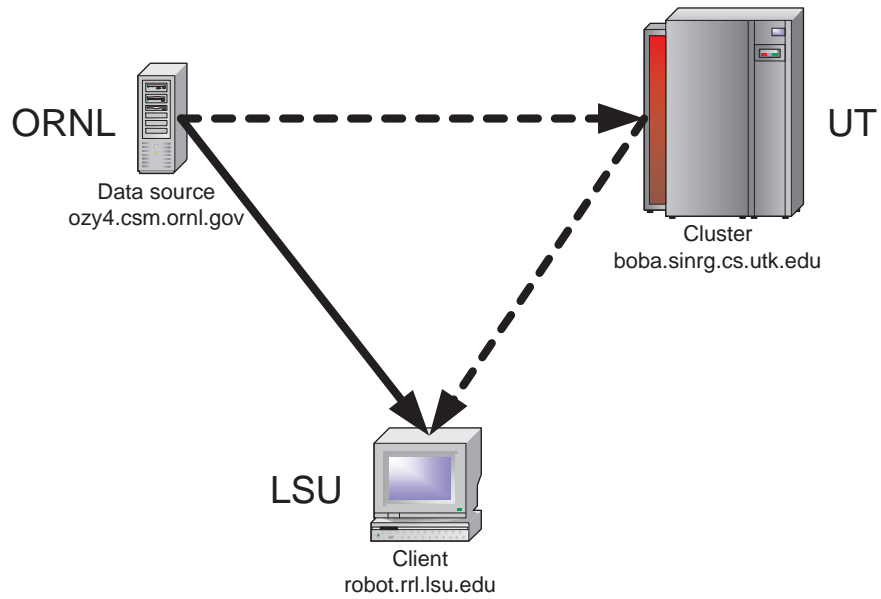


Figure 10: Remote visualization system deployment.

Raw data	Send data	Isosurface extraction	Rendering	Display	Interactive response (sec)
Hand (5.24MB)	ozy4 @ ORNL		robot @ LSU		27.07
Engine (7.21MB)					41.53
Hand (5.24MB)	ozy4 @ ORNL			robot @ LSU	15.21
Engine (7.21MB)					27.37
Hand (5.24MB)	ozy4 @ ORNL	boba @ UT	robot @ LSU		5.84
Engine (7.21MB)					8.15

Figure 11: Experimental results.