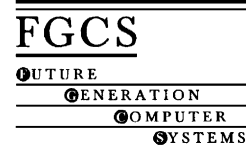




ELSEVIER

Available at
www.ComputerScienceWeb.com
 POWERED BY SCIENCE @ DIRECT®

Future Generation Computer Systems 1005 (2003) 1–15


www.elsevier.com/locate/future

Quanta: a toolkit for high performance data delivery over photonic networks

Eric He, Javid Alimohideen, Josh Eliason, Naveen K. Krishnaprasad,
 Jason Leigh*, Oliver Yu, Thomas A. DeFanti

Electronic Visualization Laboratory, University of Illinois at Chicago, MC152, 1120 SEO, 851 S. Morgan St., Chicago, IL 60620, USA

Abstract

Quanta is a cross-platform adaptive networking toolkit for supporting the data delivery requirements of interactive and bandwidth intensive applications, such as Amplified Collaboration Environments. One of the unique goals of Quanta is to provide applications with the ability to provision optical pathways (commonly referred to as Lambdas) in dedicated photonic networks. This paper will introduce Quanta's architecture and capabilities, with particular attention given to its aggressive and predictable high performance data transport scheme called Reliable Blast UDP (RBUDP). We provide an analytical model to predict RBUDP's performance and compare the results of our model against experimental results performed over a high speed wide-area network.

© 2003 Published by Elsevier Science B.V.

Keywords: Quanta; Photonic network; Reliable Blast UDP

1. Introduction

Amplified Collaboration Environments (ACE) are physical meeting spaces that enable distantly located groups to work in intensive collaboration campaigns that are augmented by advanced collaboration, computation, and visualization systems. One example of an ACE is the *Continuum* (Fig. 1) at the Electronic Visualization Laboratory [4], at the University of Illinois at Chicago, and at the Technology Research, Education and Commercialization Center in DuPage County, Illinois [25]. ACEs are based on the concept of the "War Room" or "Project Room" which have been shown to increase the productivity of collocated work-

ing teams by a factor of 2 [29]. The goal of the Continuum is to provide the same, if not greater, benefits for distributed teams. To this end, the Continuum integrates a broad range of technologies that include: multi-party video conferencing (via the AccessGrid [1]), electronic touch screens (for intuitive shared white-boarding), passive stereoscopic displays (such as the GeoWall, for displaying data sets in true 3D [13]), high resolution tiled displays (for displaying large visualizations or mosaics of visualizations), and PDAs and laptops for wireless control of these systems. Taken as a whole, each of these systems requires one or more computers to support. Hence a full Continuum will require a compute cluster per site. Furthermore this compute cluster must also be connected to other possibly distributed computing clusters, which might house massive data sets that are being shared in the collaborative environment. At EVL, we have developed a computing paradigm called the *Optiputer*

* Corresponding author.

E-mail addresses: spiff@evl.uic.edu, cavern@evl.uic.edu (J. Leigh).

URL: <http://www.evl.uic.edu/cavern>.



Fig. 1. The continuum—an amplified collaboration environment.

50 as the primary means for supporting future genera-
 51 tion networked applications such as the Continuum.
 52 Quanta is the networking middleware for supporting
 53 applications modeled after the Optiputer.

54 The Optiputer [18] is a National Science Founda-
 55 tion funded project to interconnect distributed storage,
 56 computing and visualization resources using photonic
 57 networks. The main goal of the project is to exploit
 58 the trend that network capacity is increasing at a rate
 59 far exceeding processor speed, while at the same time
 60 plummeting in cost. This allows one to experiment
 61 with a new paradigm in distributed computing—where
 62 the photonic networks serve as the computer's system
 63 bus and compute clusters taken as a whole, serve as
 64 the peripherals in a potentially, planetary-scale com-
 65 puter. We differentiate photonic networks from optical
 66 networks as networks comprised of optical fibers and
 67 MEMS optical switching devices. There is no trans-
 68 lation of the photons to electrons and hence no rout-
 69 ing within photonic switches. Applications that con-
 70 trol these networks will direct photons directly from
 71 the start point to the end point of a series of photonic
 72 switches and hence will have full control of the avail-
 73 able bandwidth in these allocated light paths.

74 In order to optimize data delivery in Optiputer ap-
 75 plications such as ACEs, advances need to be made
 76 at several of the OSI network layers. At the phys-
 77 ical layer, shared packet-switched Internet should be
 78 replaced by exclusive photonic switched networks

to guarantee high bandwidth. We are currently de-
 80 veloping the Photonic Interdomain Negotiator (PIN)
 81 to support this capability. At the data link layer,
 82 Multiple Protocol Label Switching (MPLS) or Vir-
 83 tual LAN (VLAN) replaces the slow and inefficient
 84 layer 3 switching, while at the same time provid-
 85 ing Quality-of-Service. The Internet Protocol (IP) is
 86 still used at layer 3 in order to maintain compatibil-
 87 ity with the Internet. At the transport layer, there is
 88 already consensus among network researchers that
 89 the current TCP implementations are not suitable for
 90 long distance high performance data transfer. Either
 91 TCP needs to be modified radically or new transport
 92 protocols should be introduced.

93 We intend to address the data transport problem with
 94 Quanta, a cross-platform adaptive networking toolkit
 95 for supporting the diverse networking requirements of
 96 interactive and data intensive Optiputer applications.
 97 The goal is to provide an easy to use system that will
 98 allow programmers to specify the data transfer char-
 99 acteristics of their application at a high level, and let
 100 Quanta transparently translate these requirements into
 101 appropriate networking decisions. The decisions will
 102 include making necessary QoS reservations, and adap-
 103 tively utilizing the transport protocols to fulfill the
 104 user's data transfer requirements.

105 Quanta currently uses an optical network (Starlight)
 106 and a photonic network (Omnet) as experimen-
 107 tal testbeds. Starlight [24] is a project managed by

108 the University of Illinois at Chicago, to provide an
109 IP-over-Dense Wave Division Multiplexing (DWDM)
110 peering point for national and international optical
111 networks. Plans are underway to convert Starlight to
112 a photonic network. OMNInet is a project supported
113 by Nortel Networks, SBC Communications Inc. and
114 Ameritech to assess and validate next-generation pho-
115 tonic technologies, architectures and applications in
116 metropolitan area networks [17].

117 This paper begins with a description of Quanta's
118 architecture and capabilities, with particular attention
119 given to the development of high performance data
120 transport schemes. We describe an algorithm for an
121 aggressive bulk data transfer scheme called Reliable
122 Blast UDP (RBUDP), provide an analytical model to
123 predict its performance, and compare the results of
124 our model against our implementation of RBUDP. Fi-
125 nally we extend the analytical model to support high
126 throughput reliable data streaming, and compare it
127 with the graphics streaming experiments performed
128 during the IGrid 2002 conference in Amsterdam, The
129 Netherlands.

130 2. Related work

131 In 1995, Gilder [7] predicted that network band-
132 width would triple every year for the next 25 years.
133 So far his prediction seems to be approximately cor-
134 rect. Each fiber optical wavelength channel can run at
135 10 Gbps. Wavelength division multiplexing gives 128
136 or more channels per fiber, resulting in a combined
137 bandwidth of 1 terabits per second (almost 20,000,000
138 times faster than 56 Kbps modem connection). Con-
139 sequently, this has led to a situation where straight-
140 forward use of the BSD socket library cannot take
141 advantage of the high bandwidth available, making
142 commonly used networking protocols unsuitable for
143 high-end applications. Even if networked applications
144 could make Gigabit "lambda reservations", it does not
145 however guarantee that they will be able to make full
146 use of that bandwidth. This problem is particularly evi-
147 dent when one attempts to perform large bulk data
148 transfers over long distance, high-speed networks (of-
149 ten referred to as "long fat networks" or LFNs) [26].

150 LFNs such as those between the US and Europe or
151 Asia have extremely high round-trip latencies (at best
152 120 ms). This latency results in gross bandwidth under

153 utilization when TCP is used for data delivery. This is
154 because TCP's windowing mechanism imposes a limit
155 on the amount of data it will send before it waits for an
156 acknowledgement. International networks have long
157 delays causing TCP to spend an inordinate amount of
158 time waiting for acknowledgments. Consequently, the
159 client's data transmission will never reach the peak
160 available capacity of the network. Traditionally this
161 is "remedied" by adjusting TCP's window and buffer
162 sizes to match the $bandwidth \times delay$ product (or ca-
163 pacity) of the network. For example, for a 1 Gbps con-
164 nection between Chicago and Amsterdam, with an av-
165 erage Round Trip Time (RTT) of 110 ms, the capacity
166 is $1024 \times 0.11/8 = 14.1$ MB. Adjusting TCP win-
167 dow size is problematic for several reasons: firstly, on
168 some operating systems (such as IRIX for the SGI)
169 the window size can only be modified by building a
170 new version of the kernel—hence this is not an op-
171 eration a user-level application can invoke. Secondly,
172 one needs to know the current capacity of the network
173 in order to set the window size correctly. The current
174 capacity varies with the amount of background traffic
175 already on the network and the path to the destination.

176 Several alternative solutions are possible. One so-
177 lution is to provide TCP with better estimates of the
178 current capacity of a link. The WEB100 Consortium
179 [28], which takes this approach, is developing tech-
180 niques to modify router operating systems to report
181 available bandwidth over a network link. They are also
182 modifying operating systems kernels to allow better
183 monitoring of TCP performance. Another solution is
184 to use parallel TCP. In parallel TCP, the payload being
185 delivered is divided into N partitions, which are de-
186 livered over N TCP connections. Leigh et al. [12,19]
187 and Allcock et al. [2] have shown that parallel TCP
188 can provide throughput as high as 80% of a network's
189 available bandwidth, but its performance is unstable
190 when excessive numbers of sockets are used. More-
191 over, it is difficult to predict the correct number of
192 sockets to use. However, there is a growing commu-
193 nity of high bandwidth network users that are real-
194 izing that there is no need for a congestion control
195 mechanism if the application is able to reserve a de-
196 dicated network path such as in the case of photonic
197 networks. As a consequence, there is now great in-
198 terest in developing UDP-based protocols to improve
199 bandwidth use. Simple Available Bandwidth Utiliza-
200 tion Library (SABUL) [8] and Tsunami [27] are two

201 recent examples. Our Reliable Blast UDP (RBUDP)
 202 protocol which we developed in 2000 is another [12].
 203 The unique contribution of RBUDP is that we are able
 204 to provide an analytical model to predict its perfor-
 205 mance. This kind of predictability is important for data
 206 intensive, interactive applications.

207 3. Overall design of quanta

208 Quanta emerged from almost a decade's experi-
 209 ence in connecting immersive CAVE systems [3]
 210 to each other and to supercomputers—this concept
 211 is called Tele-Immersion [14]. Quanta's predeces-
 212 sor is CAVERNsoft [19], which has been widely
 213 used by the CAVE community to develop advanced
 214 tele-immersive applications. Consequently, Quanta
 215 inherits all of the data sharing abstractions that have
 216 been found to be useful for developing these ap-
 217 plications; and networked applications in general.
 218 Quanta aims to provide an Adaptive Network Con-
 219 troller (ANC) (Fig. 2) and three supporting services:
 220 a Resource Monitor, a Quality-of-Service provisioner
 221 and a collection of data transport mechanisms and
 222 data sharing abstractions. The ANC's first role is to
 223 take application-specified data delivery requirements
 224 (e.g. bandwidth, latency, jitter, reliability, etc) and
 225 translate them into networking and computational
 226 resource allocations needed to meet the applications'
 227 demands. The ANC will monitor the current state

of the network or QoS capability, select an optimal
 transmission protocol, and make QoS requests (if
 available.) If QoS is available, the ANC contacts the
 Admission Control system to determine if the desired
 bandwidth is available, and then makes a reservation
 using the Signaling Controller. Once the strategy has
 been activated, the ANC will monitor the progress
 of the data transmission and adjust networking and
 computational parameters to sustain the desired per-
 formance. To accommodate multiple simultaneous
 and heterogeneous network flows, the ANC may alter
 some of the low-level transport protocol parameters
 (such as buffer size) or may adjust QoS reserva-
 tions dynamically. The Signaling Controller main-
 tains device-independence via a plug-in architecture
 that dynamically loads-in service-provider-specific li-
 braries to signal for QoS. In the case of our photonic
 testbeds, the signaling controller interacts with PIN
 to make light path reservations.

3.1. Photonic interdomain negotiator: interdomain light path provisioning for quanta

Work is currently underway to develop a software
 infrastructure for light path provisioning on photonic
 networks (Fig. 3). While Quanta can ensure that
 data is optimally delivered over these light paths, it
 presently does not have the ability to allocate these
 dedicated light paths. This is the role of PIN. An ap-
 plication wishing to allocate a light path between two

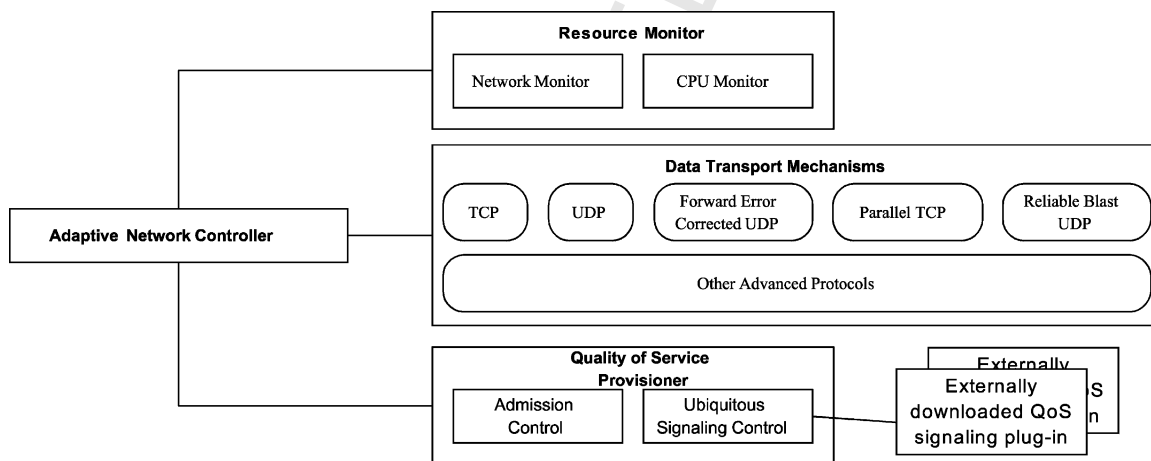


Fig. 2. Quanta's adaptive networking system.

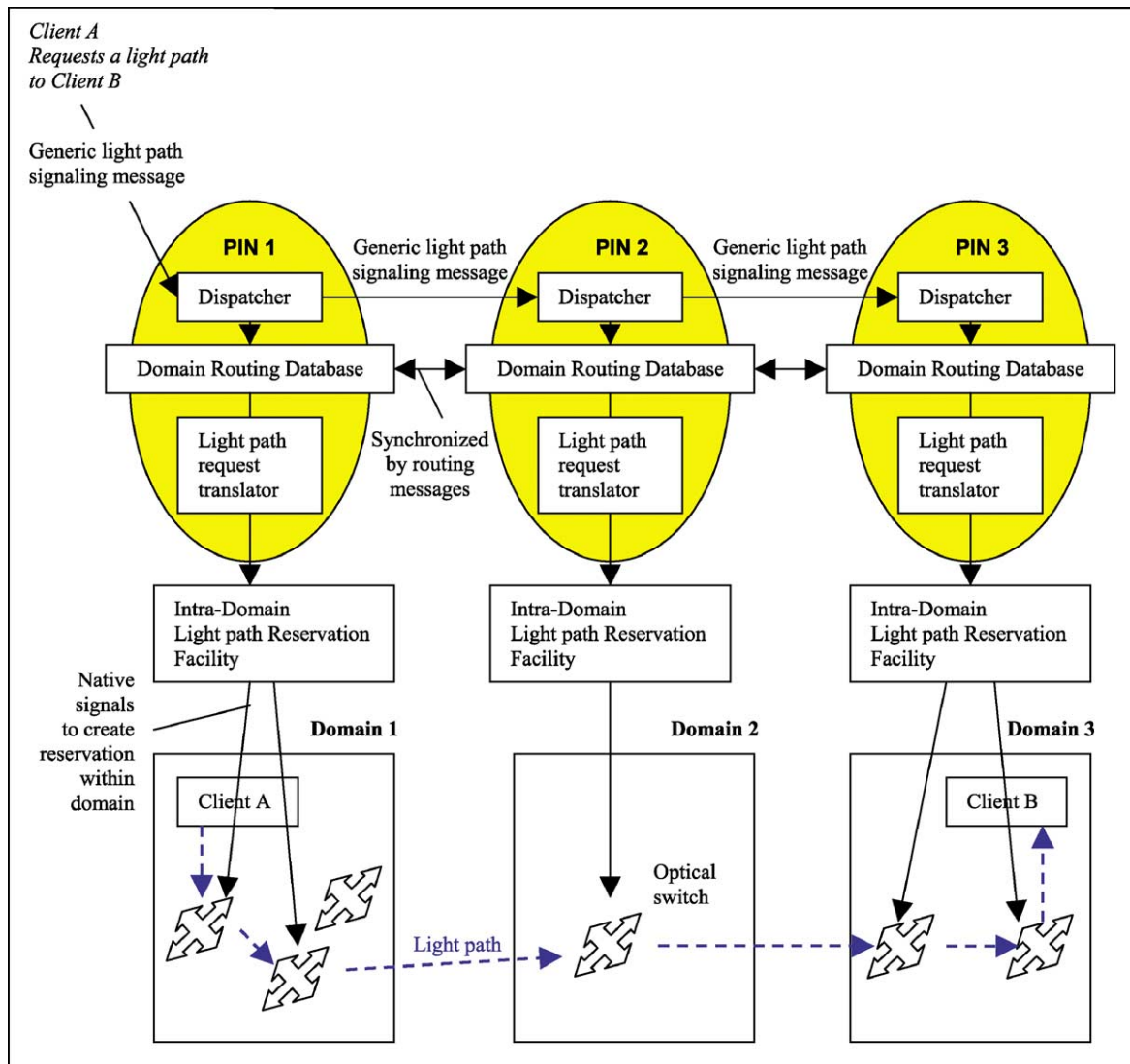


Fig. 3. Architecture of the interdomain light path provisioner.

255 end points, contacts its local PIN which will dispatch
 256 generic light path signaling messages to neighboring
 257 PINs until the final destination is reached. Each PIN
 258 will translate the generic light path signaling message
 259 into the native photonic signaling message that is un-
 260 derstood by the local intradomain light path signaling
 261 facility. This facility then signals the photonic switch
 262 to make adjustments to its internal MEMS switches
 263 to establish the connection. At the present time a pro-
 264 totype of PIN is being developed, and TL1 command

sets and APIs from multiple vendors such as Nortel, 265
 Glimmerglass, Calient and IMMI, are being examined 266
 to identify common commands that PIN will need to 267
 support. 268

3.2. Quanta's data transport and data sharing capabilities 269

Quanta's data transport capabilities include 271
 C++ classes that simplify socket-level programming 272

of TCP, UDP and multicast communications (these are encapsulated in the C++ classes: `QUANTAnet_tcp.c`, `QUANTAnet_udp.c`, `QUANTAnet_mcast.c`, respectively). The reader is encouraged to examine the Quanta API manual [20] for a detailed explanation of how one goes about using the individual C++ classes. The names of the C++ classes are provided here as a reference. All the data transport classes have performance monitoring built into them so that an application can easily determine how much bandwidth it is using and how much latency it is experiencing. As Quanta is a cross-platform toolkit, it provides a data packing API that allows applications to ensure that their transmissions are correctly translated into the format of the target computer system (`QUANTAnet_datapack.c`). Quanta also provides a set of threading and mutual exclusion classes (`QUANTAts_mutex.c`, `QUANTAts_thread.c`, `QUANTAts_condition.c`).

Quanta provides a number of data sharing abstractions. These are described below:

QUANTAnet_tcpReflector.c and QUANTAnet_udpReflector.c: Data reflection is a unicast method for emulating multicast. Clients send information to a central server rather than a single multicast address and the reflector repeats/reflects that same information to all other subscribing clients. From our experience we have found that this is one of the most heavily used capabilities for supporting data sharing in collaborative applications. The UDP reflector provides both unicast reflection and multicast bridging. This enables groups of clients to operate multicast within separated domains and share information across them using a bridge rather than having to set up a multicast tunnel, which often requires system administrator privileges. The TCP reflector is similar to the UDP reflector in that it places boundaries on TCP messages (making them discrete) instead of broadcasting them as a continuous stream.

QUANTAdb.c, QUANTAmisc_observer.c and QUANTAmisc_subject.c: Quanta provides persistent distributed shared memory emulation via the `QUANTAdb` (or database) class. This is essentially a client/server database with automatic data reflection. Hence any updates to the database are propagated to all subscribers of the database. Clients

are notified either via a traditional callback function or via a subject/observer mechanism [6]. This is essentially an object-oriented replacement for callbacks. The subject maintains a list of its observers for specific events and each observer will be triggered whenever the specific event occurs. The database assumes a Unix-like directory hierarchy with the leaf nodes containing the individual data values. These data values are intended to be small to expedite state information sharing rather than bulk data sharing.

QUANTAnet_rpc.c: To complement Quanta's distributed shared memory and message passing capabilities, remote procedure calling is also provided. This allows clients and servers to invoke each other's functions and procedures. This is a widely used technique for distributed computing.

QUANTAnet_http.c: This is a C++ class to access WEB servers.

QUANTAnet_parallelTcp.c: This class works like Quanta's regular TCP socket class except a data buffer is partitioned and transmitted over several sockets rather than just one. Parallel TCP has been shown to be able to overcome the LFN problem, however the performance becomes unstable when too many parallel sockets are used [19].

QUANTAnet_remoteFileIO32.c, QUANTAnet_remoteFileIO64.c, QUANTAnet_remoteParallelFileIO32.c and QUANTAnet_remoteParallelFileIO64.c classes: The Remote File I/O classes provide the capability for uploading and downloading files from a remote server. The provision of both 32 and 64 bit versions as well as parallel socket versions of the class allows for the efficient delivery of all file sizes, including those larger than 2 GB. The 64 bit version effectively allows one to deliver Terabyte files.

QUANTAnet_fecClient.c and QUANTAnet_fecServer.c: For long distance networks such as international networks, latencies are high (on the order of hundreds of milliseconds). In advanced collaborative applications, we would ideally like state updates in the shared environment to occur with a minimum amount of latency and with a high degree of reliability. Data should be transmitted reliably over long distances without the acknowledgement typically used in protocols such as TCP. We have applied Forward Error Correction (FEC) to achieve this [5].

368 FEC collects between 1 and N (typically 2 or 3) data
 369 packets and performs a bit-wise operation on the
 370 packets (such as XOR), to produce a “redundant”
 371 packet. This packet is delivered along with the
 372 regular UDP traffic as a separate UDP stream. If
 373 any data packets are lost, FEC packets can be used
 374 to reconstruct the missing packet. By using such
 375 a scheme the latency and jitter can be reduced for
 376 reliable transmission over long distance networks.
 377 *QUANTAnet_rbudpSender.c* and *QUANTAnet_rbudpReceiver.c*:
 378 When operating over dedicated networks the probability of packet loss is low. To take
 379 advantage of this opportunity one can use UDP
 380 augmented with acknowledgements. The Reliable
 381 Blast UDP (RBUDP) scheme works by “blasting”
 382 the contents of a data file at just below the avail-
 383 able bandwidth without asking the remote site to
 384 acknowledge any of the packets [12]. Hence, all the
 385 available bandwidth is used for pure data transmis-
 386 sion. At the remote site, a tally is kept for all the
 387 packets that have arrived and, after some timeout
 388 period, a list of missing packets is sent back to the
 389 sending client. The sender reacts by resending all
 390 the missing packets and again waiting for another
 391 negative acknowledgement, and so on. The next
 392 section focuses deeply into RBUDP, as it has re-
 393 cently gained significant importance as a technique
 394 in overcoming TCP’s inability to fill high band-
 395 width networks.
 396

4. Reliable Blast UDP

397

398 Reliable Blast UDP [9] has two main goals. The
 399 first is to keep the network pipe as full as possible
 400 during bulk data transfer. The second goal is to avoid
 401 TCP’s per-packet interaction so that acknowledgments
 402 are not sent per window of transmitted data, but aggre-
 403 gated and delivered at the end of a transmission phase.
 404 Fig. 4 illustrates the RBUDP data delivery scheme. In
 405 the first data transmission phase (A to B in the figure),
 406 RBUDP sends the entire payload at a user-specified
 407 sending rate using UDP datagrams. Since UDP is an
 408 unreliable protocol, some datagrams may become lost
 409 due to congestion or an inability of the receiving host
 410 from reading the packets rapidly enough. The receiver
 411 therefore must keep a tally of the packets that are re-
 412 ceived in order to determine which packets must be
 413 retransmitted. At the end of the bulk data transmis-
 414 sion phase, the sender sends a DONE signal via TCP
 415 (C in the figure) so that the receiver knows that no
 416 more UDP packets will arrive. The receiver responds
 417 by sending an acknowledgment consisting of a bitmap
 418 tally of the received packets (D in the figure). The
 419 sender responds by resending the missing packets, and
 420 the process repeats itself until no more packets need
 421 to be retransmitted.

422 In RBUDP, the most important input parameter is
 423 the sending rate of the UDP blasts. To minimize loss,
 the sending rate should not be larger than the band-

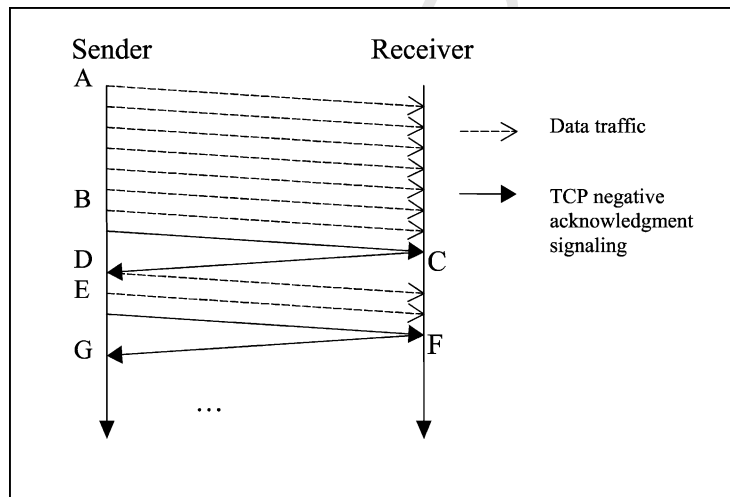


Fig. 4. Time sequence diagram of RBUDP.

width of the bottleneck link. Tools such as iperf [10] and netperf [16] are typically used to measure the bottleneck bandwidth. In theory if one could send data just below this rate, data loss should be near zero. In practice, however, other factors need to be considered. In our first implementation of RBUDP, we chose a send rate of 5% less than the available network bandwidth predicted by iperf. Surprisingly this resulted in approximately 33% loss! After further investigation we found that the problem was in the end host rather than the network. Specifically, the receiver was not fast enough to keep up with the network while moving data from the kernel buffer to application buffers. When we used a faster computer as the receiver, the loss rate decreased to less than 2%. The details of this experiment are further discussed in Section 4.2.

The chief problem with using iperf as a measure of possible throughput over a link is that it does not take into account the fact that, in a real application, data is not simply streamed to a receiver and discarded. It has to be moved into main memory for the application to use. This has motivated us to produce app_perf (a modified version of iperf) to take into account an extra memory copy that most applications must perform. We can therefore use app_perf as a more realistic bound for how well a transmission scheme should be able to reasonably obtain. In the experiments detailed in Section 4.2, we will include both iperf and app_perf's prediction of available bandwidth.

Three versions of RBUDP were developed:

- *RBUDP without scatter/gather optimization*: This is a naïve implementation of RBUDP where each incoming packet is examined (to determine where it should go in the application's memory buffer) and then moved there.
- *RBUDP with scatter/gather optimization*: This implementation takes advantage of the fact that most incoming packets are likely to arrive in order, and if transmission rates are below the maximum throughput of the network, packets are unlikely to be lost. The algorithm works by using readv() to directly move the data from kernel memory to its predicted location in the application's memory. After performing this readv() the packet header is examined to determine if it was placed in the correct location. If it was not (either because it was an out-of-order packet, or an intermediate packet was lost), then

the packet is moved to the correct location in the user's memory buffer. This optimization can improve the throughput by 10% when the receiving host is slower than the network. [9]

- *"Fake" RBUDP*: This implementation is the same as the scheme without the scatter/gather optimization except the incoming data is never moved to application memory. This was used to examine the overhead of the RBUDP protocol compared to raw transmission of UDP packets via iperf.

Experiments that compare these versions of the protocol, and an analytical model of RBUDP, will be presented next.

4.1. Analytical model for RBUDP

The purpose of developing an analytical model for RBUDP is twofold. Firstly we wanted to develop an equation similar to the "bandwidth \times delay product" equation for TCP, to allow us to predict RBUDP performance over a given network. Secondly we wanted to systematically identify the factors that influenced the overall performance of RBUDP so that we can predict how much benefit any potential enhancement in the RBUDP algorithm might provide.

Firstly, all variables are defined as follows: $B_{\text{achievable}}$ is the achievable bandwidth, B_{send} the chosen send rate, S_{total} the total data size to send (i.e., payload), T_{total} the total predicted send time, T_{prop} the propagation delay, T_{udpSend_i} the time to send UDP blast on i th iteration, N_{resend} the number of times to resend (depends on loss%), T_{ack} the time to acknowledge a blast (at least 1 ACK is always needed), L_i the % packet loss on i th iteration.

In our model we are attempting to predict the achievable bandwidth ($B_{\text{achievable}}$) of RBUDP:

$$B_{\text{available}} = \frac{S_{\text{total}}}{T_{\text{total}}} \quad (1)$$

Following the RBUDP algorithm, we estimate T_{total} as

$$T_{\text{total}} = (T_{\text{prop}} + T_{\text{udpSend}_0}) + \left(\sum_{i=1}^{N_{\text{resend}}} (T_{\text{prop}} + T_{\text{udpSend}_i}) \right) + (N_{\text{resend}} + 1)(T_{\text{ack}} + T_{\text{prop}}) \quad (2)$$

511 In (2), the first term is the time to send the main pay-
 512 load, the second term is the time to transmit missing
 513 packets, called T_{resend} , the last term is the time to send
 514 each acknowledgement.

515 Specifically:

$$517 \quad T_{\text{udpSend}_0} = \frac{S_{\text{total}}}{B_{\text{send}}}, \quad T_{\text{udpSend}_i} = \frac{L_{i-1} S_{\text{udpSend}_{i-1}}}{B_{\text{send}}},$$

$$518 \quad T_{\text{ack}} = \frac{S_{\text{ack}}}{B_{\text{send}}}, \quad S_{\text{ack}} = \frac{S_{\text{total}}/S_{\text{packet}}}{8},$$

$$519 \quad T_{\text{ack}} = \frac{S_{\text{total}}/8S_{\text{packet}}}{B_{\text{send}}}, \quad S_{\text{packet}} = 1.5 \text{ KB}$$

520 Consequently:

$$522 \quad T_{\text{total}} = \left(T_{\text{prop}} + \frac{S_{\text{total}}}{B_{\text{send}}} \right)$$

$$523 \quad + \left(N_{\text{resend}} T_{\text{prop}} + \sum_{i=1}^{N_{\text{resend}}} \frac{L_{i-1} S_{\text{udpSend}_{i-1}}}{B_{\text{send}}} \right)$$

$$524 \quad + \left((N_{\text{resend}} + 1) \left(\frac{S_{\text{total}}}{8S_{\text{packet}} B_{\text{send}}} + T_{\text{prop}} \right) \right)$$

$$525 \quad (3)$$

526 Given this equation, let us consider two possible
 527 situations—one where no loss occurs, and one where
 528 loss does occur. If no loss occurs, we can eliminate the
 529 middle term so that the best achievable performance
 530 can be computed using:

$$532 \quad T_{\text{best}} = \left(T_{\text{prop}} + \frac{S_{\text{total}}}{B_{\text{send}}} \right) + \left(\frac{S_{\text{total}}}{8S_{\text{packet}} B_{\text{send}}} + T_{\text{prop}} \right),$$

$$533 \quad B_{\text{best}} = \frac{S_{\text{total}}}{S_{\text{total}}/B_{\text{send}} + S_{\text{total}}/8S_{\text{packet}} B_{\text{send}} + 2T_{\text{prop}}}$$

$$534 \quad (4)$$

535 In the denominator, $S_{\text{total}}/8S_{\text{packet}} B_{\text{send}}$ is very small
 536 compared to other factors and can be omitted.

537 We can then derive the ratio of B_{best} and B_{send} as

$$538 \quad \frac{B_{\text{best}}}{B_{\text{send}}} = \frac{1}{1 + (\text{RTT } B_{\text{send}}/S_{\text{total}})} \quad (5)$$

539 where $2T_{\text{prop}}$ is RTT.

540 This ratio shows that in order to maximize through-
 541 put, we should strive to minimize $\text{RTT } B_{\text{send}}/S_{\text{total}}$ by
 542 maximizing the size of the data we wish to deliver.
 543 For example, given T_{prop} for Chicago to Amsterdam
 544 is 50 ms, and B_{send} is 600 Mbps, and if we wish to

545 achieve a throughput of 90% of the sending rate, then
 546 the payload, S_{total} needs to be at least 67.5 MB.

547 In Section 4.2 (Fig. 5), we will use Eq. (3) to com-
 548 pare the theoretical best rate B_{best} against experimen-
 549 tal results, over a variety of send rates (B_{send}). Fur-
 550 thermore we will compare B_{best} against experimental
 551 results with varying payload sizes (S_{total}) (Fig. 7).

552 Now let us turn to consider the situation where loss
 553 does occur. We will take a simplifying assumption that
 554 a constant loss rate of L occurs at every pass of the al-
 555 gorithm. We realize that in a real network subsequent
 556 losses in the retransmit phases is likely to be smaller,
 557 rather than constant, because we will be retransmit-
 558 ting a significantly smaller payload at each iteration.
 559 However to estimate that accurately would require us
 560 to develop a model for the buffer in the intervening
 561 routers too. Hence we can take our simplifying as-
 562 sumption as a worst-case estimate.

563 So, given loss rate L , retransmits will occur until the
 564 amount of data left is less than one packet. Therefore
 565 the number of retransmits required can be estimated as

$$566 \quad N_{\text{resend}} = \left\lceil \log L \left(\frac{S_{\text{packet}}}{S_{\text{total}}} \right) \right\rceil \quad (6)$$

567 The data size of all retransmits is therefore:

$$568 \quad S_{\text{resend}} = S_{\text{total}} \frac{L(1 - L^{\lfloor \log L(S_{\text{packet}}/S_{\text{total}}) \rfloor})}{1 - L} \quad (7)$$

569 We can now plug (6) and (7) back into Eq. (3) to
 570 produce our new estimate of $B_{\text{achievable}}$ given constant
 571 loss rate L . In Fig. 7, we will put this prediction to
 572 use comparing an experimental situation where packet
 573 loss was observed.

574 4.2. Experimental results

575 The testbed network consisted of an OC-48 link
 576 (2.5 Gbps) brought by SURFnet from Amsterdam to
 577 the StarLight facility in Chicago. There was little-to-no
 578 traffic on the link when the experiments were per-
 579 formed. Linux PCs were placed at each end of the link.
 580 The specifications of each PC is shown in Table 1.
 581 Wgsara (in Amsterdam) was the slower PC, Charyb-
 582 dis (in Chicago) was the faster one. The network bot-
 583 tleneck resides in the Gigabit Ethernet cards of host
 584 computers.

585 In the first set of experiments, data was sent via
 586 RBUDP from the faster PC to the slower PC (from

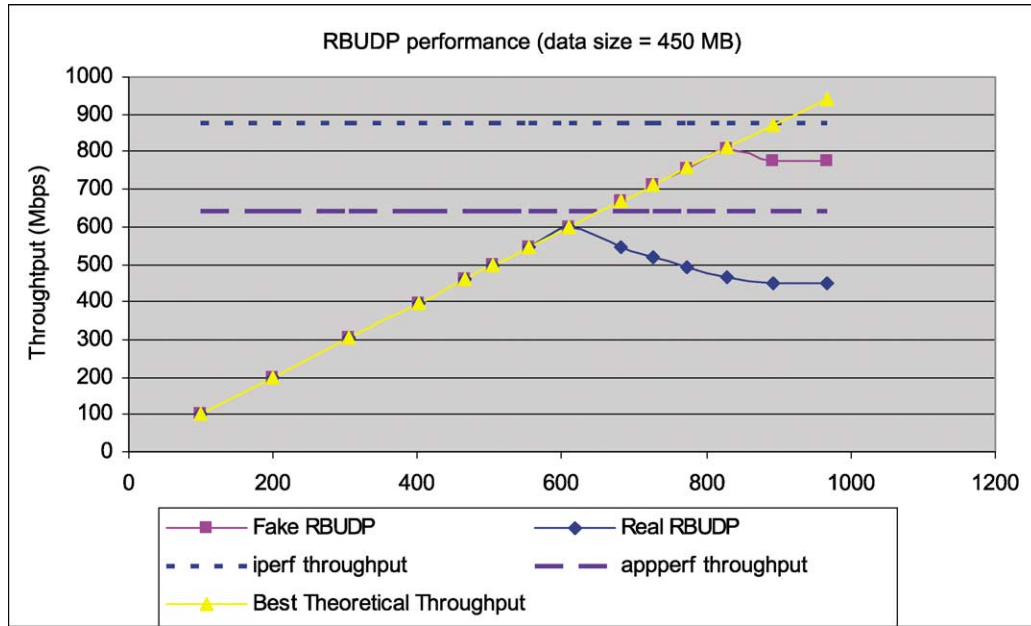


Fig. 5. RBUDP throughput from Chicago to Amsterdam. Payload is 450 MB. Bottleneck is in the receiving host. The lines indicating iperf and app_perf throughput show the maximum performance when the tools are sending at the network's full data rate. app_perf is a more realistic indication of the rate at which an application can absorb incoming data packets as it takes into account the additional overhead involved in most applications that need to take the data off the network and use it.

587 Chicago to Amsterdam). In the second set of experi-
 588 ments data was sent in the opposite direction. This al-
 589 lowed us to examine the performance of RBUDP when
 590 the bottleneck was either at the processor or in the
 591 network. The result was compared against predicted
 592 results from our analytical model. A third set of ex-
 593 periments examined RBUDP throughput for different
 594 payload sizes.

595 4.2.1. From the fast PC to the slow PC (Chicago to 596 Amsterdam)—when the bottleneck is in the receiving 597 host computer

598 In this experiment, iperf measured maximum avail-
 599 able bandwidth at 878 Mbps, and app_perf measured
 600 maximum possible throughput at 643 Mbps. In Fig. 5,

601 we plot these thresholds as lines across the top of
 602 the graph. Plotting the achieved throughput at various
 603 sending rates for the fake and real RBUDP algorithms,
 604 we notice that at sending rates below the network ca-
 605 pacity, RBUDP performs well, i.e., RBUDP gives the
 606 application exactly what the application asks for. We
 607 also notice that as the sending rates approach the ca-
 608 pacity of the network, Fake RBUDP achieves almost
 609 the same throughput as iperf, and the real RBUDP be-
 610 gins to hurt in performance because the underpowered
 611 CPU is unable to keep up with handling the incoming
 612 packets. However, as real RBUDP is able to match the
 613 maximum performance of app_perf, this means that
 614 RBUDP is making as much use of the network for use-
 615 ful data transfer as the CPU will allow. Finally, notice

Table 1
 Specification of host PCs in the experimental testbed

Host name	CPU	Memory size	System bandwidth
Wgsara2.phys.uu.nl (Amsterdam)	Pentium III 800 MHz	512 MB	238 MB/s
Charybdis.sl.startap.net (Chicago)	XEON 1.8 GHz	512 MB	844 MB/s

615 that there is a close match between our experimental
616 results and our prediction from Eq. (4) (which esti-
617 mated RBUDP performance when loss rate is zero).

618 4.2.2. From the slow PC to the fast PC (Amsterdam 619 to Chicago)—when the bottleneck is in the sending 620 host computer

621 We repeated the experiment in the opposite direc-
622 tion. This time the bottleneck was in the sending PC
623 rather than in the receiving PC. Fig. 6 shows that
624 when the host computer is fast enough, iperf and
625 app_perf performances match, as do the different im-
626 plementations of RBUDP. Fake RBUDP is able to
627 reach the maximum performance obtained by iperf;
628 and Real RBUDP is able to reach the maximum
629 performance obtained by app_perf—again confirming
630 RBUDP’s ability to maximize bandwidth utilization
631 for useful data delivery.

632 4.2.3. Effect of payload size on throughput

633 From the analysis in Section 4.1, we know that
634 the propagation time is the primary factor affecting
635 RBUDP overhead. For smaller payloads, the time
636 spent in the acknowledgement phase is almost con-

stant while the time spent blasting UDP packets
637 decreases. In Fig. 7, we compare an experimental situ-
638 ation where we send data at 611 Mbps (experiencing
639 no loss) against our theoretical prediction, which as-
640 sumes no loss (Eq. (3)). Furthermore we compare an
641 experimental situation sending data at 682 Mbps ex-
642 periencing 12% loss, against our theoretical prediction
643 where we assume a constant 12% loss per iteration.
644

645 Firstly, the results show that RBUDP performs best
646 for large payloads. Secondly, the results show that a
647 12% packet loss does not impact throughput greatly for
648 large payloads. Finally, our analytical models provide
649 good boundaries for our experimental results for 0 and
650 12% loss.

651 4.3. Adapting RBUDP for high speed data 652 streaming

653 Even though the initial motivation of RBUDP is
654 for bulk data transfer over long distance, some appli-
655 cations require high performance reliable streaming
656 transport. In Section 4.2.3, we showed that in order
657 to achieve fairly high throughput, the payload needs
to be large. In streaming applications, if the size of

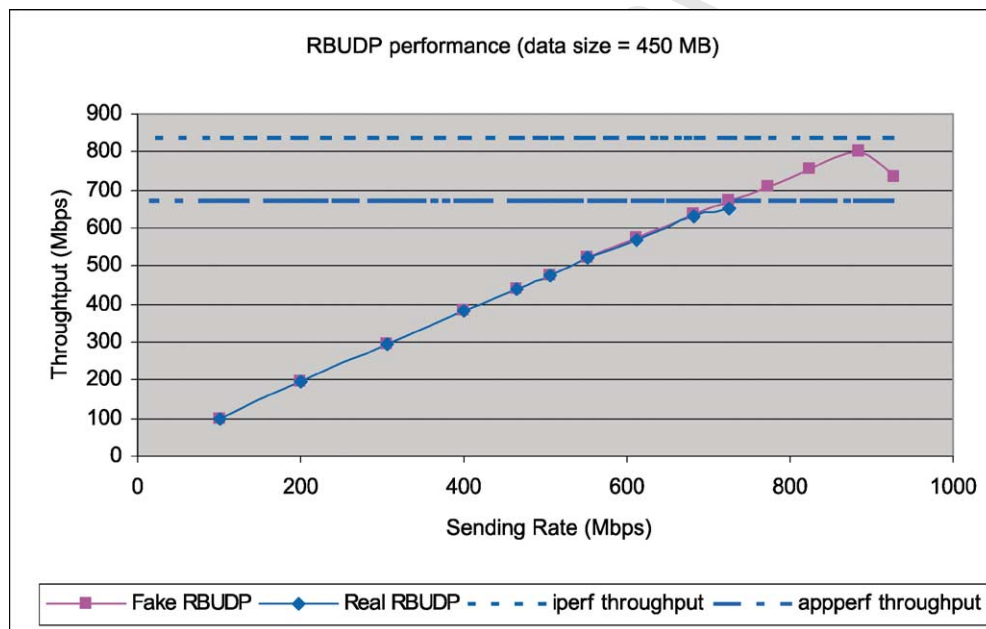


Fig. 6. RBUDP throughput from Amsterdam to Chicago. Payload is 450MB. Bottleneck is in the sending host. The maximum of the sending rate is 725 Mbps. See Fig. 5 for an explanation of the iperf and app_perf lines in the graph.

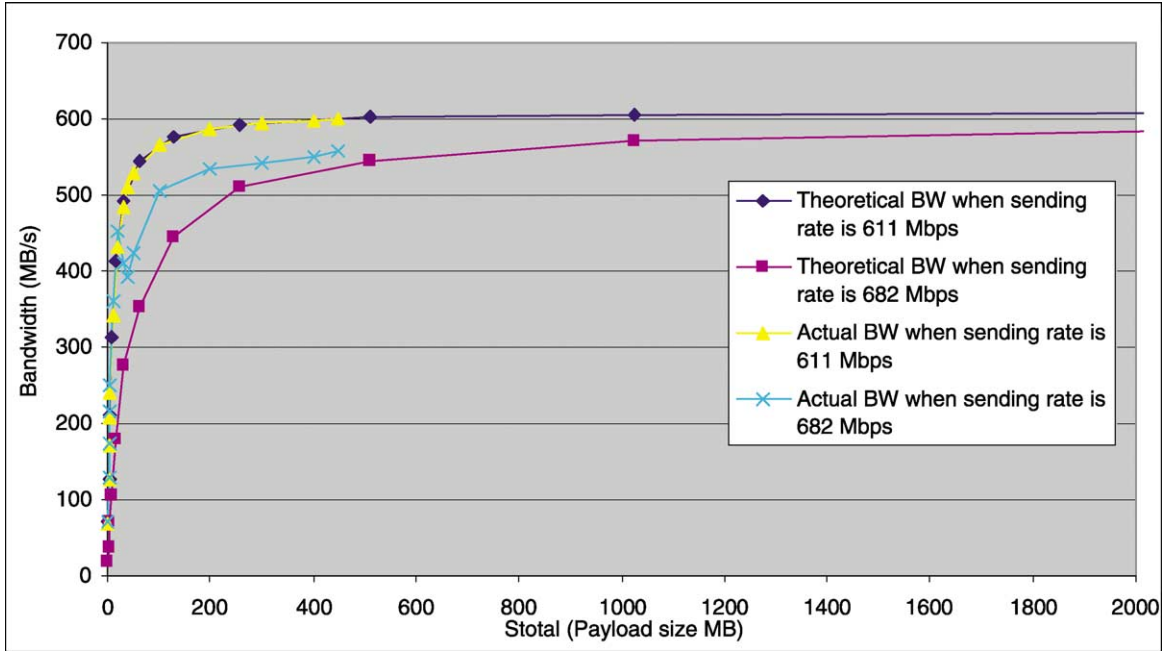


Fig. 7. Throughput vs. payload size. Larger payloads produce better network utilization.

658 objects to be streamed is small, we combine mul-
 659 tiple objects to form a large payload. However this
 660 will cause end-to-end latency to increase because
 661 of the buffering needed to form the large payloads.
 662 Based on our analytical model, we can determine the
 663 minimum sending rate needed to ensure a desired
 664 object throughput rate, given the maximum delay the
 665 application is able to tolerate.

666 Let S_{obj} is the size of streamed objects, N_{obj} the
 667 number of objects per payload, B_{obj} the required
 668 throughput of objects (number of objects per second).
 669 For example, in the case of graphics streaming, object
 670 throughput rate is measured in frames per second, D is
 671 the maximum extra delay the application can tolerate.

672 Then the size of a payload is

$$673 S_{total} = S_{obj}N_{obj} \quad (8)$$

674 where

$$675 N_{obj} = B_{obj}D \quad (9)$$

676 The required raw bandwidth is

$$677 B_{best} = B_{obj}S_{obj} \quad (10)$$

678 Assuming we are operating over an over-provisioned
 679 network, we plug (8), (9) and (10) back in Eq. (5) to
 680 compute the rate at which RBUDP needs to send data
 681 to achieve the application's requested throughput:

$$B_{Send} = \frac{S_{obj}B_{obj}}{1 - RTT/D} \quad (11) \quad 682$$

683 Hence, using a graphics streaming application as an
 684 example: given that RTT is 100 ms, S_{obj} is $800 \times 600 \times$
 685 3 (assuming image resolution of 800×600 and 3 bytes
 686 color information for each pixel) if we want to achieve
 687 a frame rate B_{obj} of 20 frames/s, the maximum extra
 688 delay introduced will be 0.5 s, the sending rate needs
 689 to be at least 288 Mbps and each payload must encapsu-
 690 late 10 image frames. During IGrid 2002, Luc Ren-
 691 nambot applied Quanta's RBUDP to a parallel graph-
 692 ics streaming application called *Griz*. Using our ana-
 693 lytical model and the parameters from the above exam-
 694 ple, we were able to predict the number of animation
 695 frames that *Griz* had to package into a single payload
 696 to achieve full utilization of the Amsterdam–Chicago
 697 Starlight link [21].

698 **5. Conclusions**

699 We have described the overall architecture and capabilities of Quanta, a cross-platform C++ toolkit for building high performance networking applications. In particular we described in detail, an aggressive bulk data transfer scheme, called RBUDP, which is intended for either dedicated, or QoS-enabled high bandwidth networks. RBUDP eliminates TCP's slow-start and congestion control mechanisms, and aggregates acknowledgments so that the full bandwidth of a link is used for pure data delivery. For large bulk transfers, RBUDP can provide delivery at precise, user-specified sending rates. RBUDP performs at its best for large payloads, rather than smaller ones. This is because with smaller payloads, the time taken for completing the delivery approaches the time taken to acknowledge the payload.

715 We have provided an analytical model that gives a good prediction of RBUDP's performance. This prediction can be used as a rule of thumb in a manner similar to the *bandwidth* \times *delay* product for TCP. In addition, this prediction can be used to estimate how future ideas for improving the algorithm might impact RBUDP performance. Even though the initial application of RBUDP is bulk data transfer over high-speed networks, this protocol can also be extended for use in streaming applications. Here an application must make a tradeoff between latency and throughput. To achieve higher throughput, latency will increase because more data must be aggregated as a single transmission payload.

729 Through the combined use of PIN and Quanta, bandwidth intensive Optiputer applications will soon be able to allocate light paths between multiple photonic domains and make full use of the available bandwidth.

734 **Uncited references**

735 [\[11,15,22,23\]](#).

736 **Acknowledgements**

737 We would like to thank Cees de Laat at University
738 of Amsterdam for providing the endpoint at SARA

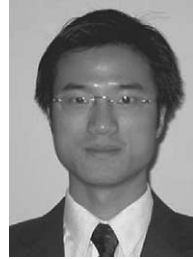
in Amsterdam to perform these experiments. The virtual reality and advanced networking research, collaborations, and outreach programs at the Electronic Visualization Laboratory (EVL) at the University of Illinois at Chicago are made possible by major funding from the National Science Foundation (NSF), awards EIA-9802090, EIA-0115809, ANI-9980480, ANI-0229642, ANI-9730202, ANI-0123399, ANI-0129527 and EAR-0218918, as well as the NSF Information Technology Research (ITR) cooperative agreement (ANI-0225642) to the University of California San Diego (UCSD) for "The Optiputer" and the NSF Partnerships for Advanced Computational Infrastructure (PACI) cooperative agreement (ACI-9619019) to the National Computational Science Alliance. EVL also receives funding from the US Department of Energy (DOE) ASCI VIEWS program. In addition, EVL receives funding from the State of Illinois, Microsoft Research, General Motors Research, and Pacific Interface on behalf of NTT Optical Network Systems Laboratory in Japan. The CAVE and ImmersaDesk are registered trademarks of the Board of Trustees of the University of Illinois. PARIS, Wanda, CAVELib and ElsieDesk are trademarks of the Board of Trustees of the University of Illinois. STAR TAP and Euro-Link are service marks of the Board of Trustees of the University of Illinois. StarLight is a service mark of the Board of Trustees of the University of Illinois at Chicago and the Board of Trustees of Northwestern University.

769 **References**

- [1] <http://www-fp.mcs.anl.gov/fl/accessgrid/>. 770
- [2] W. Allcock, J. Bester, J. Bresnahan, et al., Data management and transfer in high-performance computational grid environments, *Parallel Comput.* (2001). 771–773
- [3] C. Cruz-Neira, D. Sandin, T.A. DeFanti, Virtual reality: the design and implementation of the CAVE, in: *Proceedings of the SIGGRAPH'93 Computer Graphics Conference*, August 1993, ACM SIGGRAPH, pp. 135–142. 774–777
- [4] <http://www.evl.uic.edu/cavern/continuum/indexmain.html>. 778
- [5] R. Fang, D. Schonfeld, R. Ansari, J. Leigh, Forward Error Correction for Multimedia and Tele-immersion Streams, 2000. <http://www.startup.net/images/PDF/RayFangFEC1999.pdf>. 779–781
- [6] E. Gamma, E. Helm, R. Johnson, J. Vlissides, *Design Patterns—Elements of Reusable Object-Oriented Software*, Addison-Wesley, Reading, MA, 1995, pp. 293–303. 782–784
- [7] G. Gilder, Fiber Keeps Its Promise: Get Ready, Bandwidth will Triple Each Year for the Next 25, *Forbes*, April 7, 1997. 785–786

- 787 [8] S. Harinath, Data management support for distributed data
788 mining of large datasets over high speed wide area networks,
789 PhD Thesis, University of Illinois at Chicago, 2002.
- 790 [9] E. He, J. Leigh, O. Yu, T. DeFanti, Reliable blast UDP:
791 predictable high performance bulk data transfer, in: Proce-
792 dings of the IEEE Cluster Computing 2002, Chicago, IL,
793 September 2002.
- 794 [10] <http://dast.nlanr.net/Projects/Iperf/>.
- 795 [11] A. Johnson, M. Roussos, J. Leigh, C. Barnes, C. Vasilakis,
796 T. Moher, The NICE Project: learning together in a virtual
797 world, in: Proceedings of the IEEE VRAIS'98, Atlanta, GA,
798 March 14–18, 1998, pp. 176–183.
- 799 [12] J. Leigh, O. Yu, D. Schonfeld, R. Ansari, et al.,
800 Adaptive networking for tele-immersion, in: Proceedings of
801 the Immersive Projection Technology/Eurographics Virtual
802 Environments Workshop (IPT/EGVE), May 16–18, Stuttgart,
803 Germany, 2001.
- 804 [13] J. Leigh, G. Dawe, J. Talandis, E. He, S. Venkataraman, J.
805 Ge, D. Sandin, T.A. DeFanti, AGAVE: access grid augmented
806 virtual environment, in: Proceedings of the AccessGrid
807 Retreat, Argonne, IL, January 2001.
- 808 [14] J. Leigh, T.A. DeFanti, A. Johnson, M. Brown, D.
809 Sandin, Global tele-immersion: better than being there, in:
810 Proceedings of the Seventh International Conference on
811 Artificial Reality and Tele-Existence, Tokyo, Japan, December
812 1997, pp. 10–17.
- 813 [15] J. Leigh, A. Johnson, T.A. DeFanti, Issues in the design of
814 a flexible distributed architecture for supporting persistence
815 and interoperability in collaborative virtual environments,
816 in: Proceedings of the Supercomputing'97, San Jose, CA,
817 November 15–21, 1997.
- 818 [16] <http://netperf.org/netperf/NetperfPage.html>.
- 819 [17] [http://www.evl.uic.edu/activity/template_act_project.php?
820 indi=147](http://www.evl.uic.edu/activity/template_act_project.php?indi=147).
- 821 [18] <http://www.evl.uic.edu/cavern/optiputer>.
- 822 [19] K. Park, Y. Cho, N. Krishnaprasad, C. Scharver, M. Lewis,
823 J. Leigh, A. Johnson, CAVERNsoft G2: a toolkit for high
824 performance tele-immersive collaboration, in: Proceedings
825 of the ACM Symposium on Virtual Reality Software and
826 Technology 2000, Seoul, South Korea, October 22–25, 2000,
827 pp. 8–15.
- 828 [20] <http://www.evl.uic.edu/cavern/quanta/documentation.html>.
- 829 [21] L. Renambot, T.V.D. Schaaf, H.E. Bal, D. Germans, H.J.W.
830 Spoelder, Griz: experience with remote visualization over
831 optical grids, in: Proceedings of the IGrid 2002, October
832 2002.
- 833 [22] D.C. Schmidt, S.D. Huston, C++ Network Programming:
834 Mastering Complexity Using ACE and Patterns, Addison-
835 Wesley, Reading, MA, 2001.
- 836 [23] R. Singh, J. Leigh, T.A. DeFanti, F. Karayannis, TeraVision: a
837 high speed, high resolution graphics streaming device for the
838 OptIPuter, in: Proceedings of the IGrid 2002, October 2002.
- 839 [24] <http://www.startup.net/starlight>.
- 840 [25] <http://www.trecc.org>.
- 841 [26] W.R. Stevens, TCP/IP Illustrated, vol. 1, Addison-Wesley,
842 Reading, MA, 1994, pp. 344–350.
- 843 [27] <http://www.indiana.edu/~anml/anmlresearch.html>.

- [28] <http://www.web100.org>. 844
- [29] J.S. Olson, L. Covi, E. Rocco, W.J. Miller, P. Allie, A room of 845
your own: what would it take to help remote groups work as 846
well as collocated groups? in: Proceedings of the Conference 847
on Human Factors in Computing Systems (CHI'98), Short 848
Paper, 1998, pp. 279–280. 849



Eric He received his BS, MS degrees in Optical Engineering from Beijing Institute of Technology in 1994 and 1997, respectively. He is currently working toward a PhD degree in Computer Science at the Electronic Visualization Laboratory (EVL) at the University of Illinois at Chicago. His research interests include High Performance Network Transport Protocols, Optical Networking, Quality-of-Service, Cluster Computing and Virtual Reality.

850



Javid Alimohideen is a Master's student at Electronic Visualization Laboratory, UIC. His current research work involves development and maintenance of the Quanta toolkit.

851



Josh Eliason is an MFA student and Research Assistant at Electronic Visualization Laboratory, University of Illinois at Chicago. Research interests include human computer interaction with augmented and virtual reality systems, augmented and virtual reality tracking systems-tracking and controller devices, and distributed heterogeneous networking-mobile computing, device discovery, services, and synchronization.

852



Naveen K. Krishnaprasad received his Bachelor of Engineering (BE) in Instrumentation and Control from the University of Madras, India and his Master of Science (MS) in Computer Science from the University of Illinois at Chicago. Now he is a PhD student at the Electronic Visualization Laboratory (EVL) at University of Illinois at Chicago (UIC). His research interests include distributed com-

853

puting, large data visualization, and collaborative environments. Krishnaprasad worked at Christian Michelsen Research AS in Norway, on the design of a virtual reality kernel for tele-immersive applications. His prior work at EVL included development and maintenance of the CAVERNsoft/QUANTA networking library and design of a unified collaborative framework for network analysis.

network development. He is an Assistant Professor in the ECE department at the University of Illinois at Chicago. He has been working on research projects that are funded by the NSF, DOE and NTT Japan. His research interests are in the areas of intelligent optical network signaling, next generation Internet service provisioning, and 3G/4G wireless mobile networks.

854



Jason Leigh is an Associate Professor in the department of Computer Science at the University of Illinois at Chicago; and a Senior Scientist at the Electronic Visualization Laboratory (EVL). Leigh is a co-chair of the Global Grid Forum's Advanced Collaborative Environments research group; and a co-founder of the GeoWall Consortium.



Thomas A. DeFanti is the Director of the Electronic Visualization Laboratory (EVL), a distinguished professor in the department of Computer Science, and director of the Software Technologies Research Center at the University of Illinois at Chicago.

862

855



Oliver Yu received his Master's and PhD degrees in Electrical and Computer Engineering from the University of British Columbia. He has over 12 years of industry experiences in computer networking and telecommunications. He held technical leadership positions in Microtel Pacific Research, Nortel and Hughes. He last worked in Motorola as the section manager for the GPRS wireless mobile

UNCORRECTED PROOF