

# The *ns* Manual (formerly *ns* Notes and Documentation)<sup>1</sup>

The VINT Project

A Collaboration between researchers at  
UC Berkeley, LBL, USC/ISI, and Xerox PARC.

Kevin Fall <kfall@ee.lbl.gov>, Editor  
Kannan Varadhan <kannan@catarina.usc.edu>, Editor

August 24, 2006

*ns* © is LBNL's Network Simulator [24]. The simulator is written in C++; it uses OTcl as a command and configuration interface. *ns* v2 has three substantial changes from *ns* v1: (1) the more complex objects in *ns* v1 have been decomposed into simpler components for greater flexibility and composability; (2) the configuration interface is now OTcl, an object oriented version of Tcl; and (3) the interface code to the OTcl interpreter is separate from the main simulator.

Ns documentation is available in html, Postscript, and PDF formats. See <http://www.isi.edu/nsnam/ns/ns-documentation.html> for pointers to these.

---

<sup>1</sup>The VINT project is a joint effort by people from UC Berkeley, USC/ISI, LBL, and Xerox PARC. The project is supported by the Defense Advanced Research Projects Agency (DARPA) at LBL under DARPA grant DABT63-96-C-0105, at USC/ISI under DARPA grant ABT63-96-C-0054, at Xerox PARC under DARPA grant DABT63-96-C-0105. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the DARPA.

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Undocumented Facilities</b>	<b>6</b>
<b>I</b>	<b>Interface to the Interpreter</b>	<b>8</b>
<b>3</b>	<b>OTcl Linkage</b>	<b>9</b>
3.1	Concept Overview . . . . .	9
3.2	Code Overview . . . . .	10
3.3	Class Tcl . . . . .	10
3.3.1	Obtain a Reference to the class Tcl instance . . . . .	11
3.3.2	Invoking OTcl Procedures . . . . .	11
3.3.3	Passing Results to/from the Interpreter . . . . .	11
3.3.4	Error Reporting and Exit . . . . .	12
3.3.5	Hash Functions within the Interpreter . . . . .	12
3.3.6	Other Operations on the Interpreter . . . . .	13
3.4	Class TclObject . . . . .	13
3.4.1	Creating and Destroying TclObjects . . . . .	14
3.4.2	Variable Bindings . . . . .	15
3.4.3	Variable Tracing . . . . .	17
3.4.4	commandMethods: Definition and Invocation . . . . .	18
3.5	Class TclClass . . . . .	20
3.5.1	How to Bind Static C++ Class Member Variables . . . . .	21
3.6	Class TclCommand . . . . .	23
3.7	Class EmbeddedTcl . . . . .	24
3.8	Class InstVar . . . . .	25
<b>II</b>	<b>Simulator Basics</b>	<b>27</b>
<b>4</b>	<b>The Class Simulator</b>	<b>28</b>
4.1	Simulator Initialization . . . . .	28
4.2	Schedulers and Events . . . . .	28
4.2.1	The List Scheduler . . . . .	29
4.2.2	the heap scheduler . . . . .	30
4.2.3	The Calendar Queue Scheduler . . . . .	30
4.2.4	The Real-Time Scheduler . . . . .	30
4.2.5	Precision of the scheduler clock used in ns . . . . .	30
4.3	Other Methods . . . . .	30
4.4	Commands at a glance . . . . .	32

<b>5</b>	<b>Nodes and Packet Forwarding</b>	<b>35</b>
5.1	Node Basics . . . . .	35
5.2	Node Methods: Configuring the Node . . . . .	37
5.3	Node Configuration Interface . . . . .	39
5.4	The Classifier . . . . .	40
5.4.1	Address Classifiers . . . . .	43
5.4.2	Multicast Classifiers . . . . .	44
5.4.3	MultiPath Classifier . . . . .	45
5.4.4	Hash Classifier . . . . .	45
5.4.5	Replicator . . . . .	47
5.5	Routing Module and Classifier Organization . . . . .	48
5.5.1	Routing Module . . . . .	48
5.5.2	Node Interface . . . . .	50
5.6	Commands at a glance . . . . .	51
<b>6</b>	<b>Links: Simple Links</b>	<b>53</b>
6.1	Instance Procedures for Links and SimpleLinks . . . . .	54
6.2	Connectors . . . . .	56
6.3	Object hierarchy . . . . .	56
6.4	Commands at a glance . . . . .	57
<b>7</b>	<b>Queue Management and Packet Scheduling</b>	<b>60</b>
7.1	The C++ Queue Class . . . . .	60
7.1.1	Queue blocking . . . . .	61
7.1.2	PacketQueue Class . . . . .	62
7.2	Example: Drop Tail . . . . .	63
7.3	Different types of Queue objects . . . . .	64
7.4	Commands at a glance . . . . .	68
7.5	Queue/JoBS . . . . .	68
7.5.1	The JoBS algorithm . . . . .	68
7.5.2	Configuration . . . . .	70
7.5.3	Tracing . . . . .	71
7.5.4	Variables . . . . .	71
7.5.5	Commands at a glance . . . . .	72
<b>8</b>	<b>Delays and Links</b>	<b>75</b>
8.1	The LinkDelay Class . . . . .	75
8.2	Commands at a glance . . . . .	76
<b>9</b>	<b>Differentiated Services Module in ns</b>	<b>77</b>
9.1	Overview . . . . .	77
9.2	Implementation . . . . .	78
9.2.1	RED queue in DiffServ module . . . . .	78
9.2.2	Edge and core routers . . . . .	78
9.2.3	Policy . . . . .	79
9.3	Configuration . . . . .	80
9.4	Commands at a glance . . . . .	82
<b>10</b>	<b>Agents</b>	<b>86</b>
10.1	Agent state . . . . .	86
10.2	Agent methods . . . . .	86
10.3	Protocol Agents . . . . .	87
10.4	OTcl Linkage . . . . .	88
10.4.1	Creating and Manipulating Agents . . . . .	88
10.4.2	Default Values . . . . .	88

10.4.3	OTcl Methods	89
10.5	Examples: Tcp, TCP Sink Agents	89
10.5.1	Creating the Agent	89
10.5.2	Starting the Agent	90
10.5.3	Processing Input at Receiver	91
10.5.4	Processing Responses at the Sender	92
10.5.5	Implementing Timers	93
10.6	Creating a New Agent	93
10.6.1	Example: A “ping” requestor (Inheritance Structure)	93
10.6.2	The <code>recv()</code> and <code>timeout()</code> Methods	94
10.6.3	Linking the “ping” Agent with OTcl	94
10.6.4	Using the agent through OTcl	96
10.7	The Agent API	96
10.8	Different agent objects	96
10.9	Commands at a glance	99
<b>11</b>	<b>Timers</b>	<b>101</b>
11.1	C++ abstract base class <code>TimerHandler</code>	101
11.1.1	Definition of a new timer	102
11.1.2	Example: Tcp retransmission timer	102
11.2	OTcl Timer class	105
11.3	Commands at a glance	105
<b>12</b>	<b>Packet Headers and Formats</b>	<b>106</b>
12.1	A Protocol-Specific Packet Header	106
12.1.1	Adding a New Packet Header Type	108
12.1.2	Selectively Including Packet Headers in Your Simulation	108
12.2	Packet Classes	109
12.2.1	The Packet Class	109
12.2.2	<code>p_info</code> Class	112
12.2.3	The <code>hdr_cmn</code> Class	112
12.2.4	The <code>PacketHeaderManager</code> Class	113
12.3	Commands at a glance	114
<b>13</b>	<b>Error Model</b>	<b>116</b>
13.1	Implementation	116
13.2	Configuration	117
13.3	Multi-state error model	118
13.4	Commands at a glance	119
<b>14</b>	<b>Local Area Networks</b>	<b>121</b>
14.1	Tcl configuration	121
14.2	Components of a LAN	122
14.3	Channel Class	123
14.3.1	Channel State	123
14.3.2	Example: Channel and classifier of the physical layer	123
14.3.3	Channel Class in C++	123
14.4	<code>MacClassifier</code> Class	124
14.5	MAC Class	125
14.5.1	Mac State	125
14.5.2	Mac Methods	125
14.5.3	Mac Class in C++	125
14.5.4	CSMA-based MAC	126
14.6	LL (link-layer) Class	127
14.6.1	LL Class in C++	127

14.6.2 Example: Link Layer configuration . . . . .	127
14.7 LanRouterclass . . . . .	128
14.8 Other Components . . . . .	128
14.9 LANs and <i>ns</i> routing . . . . .	128
14.10 Commands at a glance . . . . .	130
<b>15 The (Revised) Addressing Structure in NS</b>	<b>131</b>
15.1 The Default Address Format . . . . .	131
15.2 The Hierarchical Address Format . . . . .	131
15.2.1 Default Hierarchical Setting . . . . .	132
15.2.2 Specific Hierarchical Setting . . . . .	132
15.3 The Expanded Node-Address Format . . . . .	132
15.4 Expanding port-id field . . . . .	132
15.5 Errors in setting address format . . . . .	133
15.6 Commands at a glance . . . . .	133
<b>16 Mobile Networking in ns</b>	<b>134</b>
16.1 The basic wireless model in ns . . . . .	134
16.1.1 Mobilenode: creating wireless topology . . . . .	134
16.1.2 Creating Node movements . . . . .	138
16.1.3 Network Components in a mobilenode . . . . .	139
16.1.4 Different MAC layer protocols for mobile networking . . . . .	142
16.1.5 Different types of Routing Agents in mobile networking . . . . .	143
16.1.6 Trace Support . . . . .	144
16.1.7 Revised format for wireless traces . . . . .	148
16.1.8 Generation of node-movement and traffic-connection for wireless scenarios . . . . .	150
16.2 Extensions made to CMU's wireless model . . . . .	151
16.2.1 wired-cum-wireless scenarios . . . . .	151
16.2.2 MobileIP . . . . .	152
16.3 Lists of changes for merging code developed in older version of ns (2.1b5 or later) into the current version (2.1b8) . . . . .	155
16.4 Commands at a glance . . . . .	157
<b>17 Satellite Networking in ns</b>	<b>160</b>
17.1 Overview of satellite models . . . . .	160
17.1.1 Geostationary satellites . . . . .	160
17.1.2 Low-earth-orbiting satellites . . . . .	161
17.2 Using the satellite extensions . . . . .	163
17.2.1 Nodes and node positions . . . . .	163
17.2.2 Satellite links . . . . .	164
17.2.3 Handoffs . . . . .	166
17.2.4 Routing . . . . .	167
17.2.5 Trace support . . . . .	168
17.2.6 Error models . . . . .	169
17.2.7 Other configuration options . . . . .	170
17.2.8 <i>nam</i> support . . . . .	170
17.2.9 Integration with wired and wireless code . . . . .	170
17.2.10 Example scripts . . . . .	171
17.3 Implementation . . . . .	171
17.3.1 Use of linked lists . . . . .	172
17.3.2 Node structure . . . . .	172
17.3.3 Detailed look at satellite links . . . . .	173
17.4 Commands at a glance . . . . .	175

<b>18</b>	<b>Radio Propagation Models</b>	<b>177</b>
18.1	Free space model . . . . .	177
18.2	Two-ray ground reflection model . . . . .	178
18.3	Shadowing model . . . . .	178
18.3.1	Backgroud . . . . .	178
18.3.2	Using shadowing model . . . . .	180
18.4	Communication range . . . . .	180
18.5	Commands at a glance . . . . .	181
<b>19</b>	<b>Energy Model in ns</b>	<b>182</b>
19.1	The C++ EnergyModel Class . . . . .	182
19.2	The OTcl interface . . . . .	183
<b>20</b>	<b>Directed Diffusion</b>	<b>184</b>
20.1	What is Directed Diffusion? . . . . .	184
20.2	The diffusion model in ns . . . . .	184
20.3	Some mac issues for diffusion in ns . . . . .	185
20.4	APIs for using fi lters in diffusion . . . . .	186
20.5	Ping: an example diffusion application implementation . . . . .	186
20.5.1	Ping Application as implemented in C++ . . . . .	186
20.5.2	Tcl APIs for the ping application . . . . .	187
20.6	Changes required to add yr diffusion application to ns . . . . .	187
20.7	Test-suites for diffusion . . . . .	189
20.8	Commands at a glance . . . . .	189
<b>21</b>	<b>XCP: eXplicit Congestion control Protocol</b>	<b>191</b>
21.1	What is XCP? . . . . .	191
21.2	Implementation of XCP in NS . . . . .	192
21.2.1	Endpoints in XCP . . . . .	192
21.2.2	XCP Router . . . . .	193
21.2.3	XCP queue . . . . .	193
21.3	XCP example script . . . . .	194
21.4	Test-suites for XCP . . . . .	197
21.5	Commands at a glance . . . . .	197
<b>22</b>	<b>DelayBox: Per-Flow Delay and Loss</b>	<b>198</b>
22.1	Implementation Details . . . . .	198
22.2	Example . . . . .	199
22.3	Commands at a Glance . . . . .	200
<b>III</b>	<b>Support</b>	<b>202</b>
<b>23</b>	<b>Debugging ns</b>	<b>203</b>
23.1	Tcl-level Debugging . . . . .	203
23.2	C++-Level Debugging . . . . .	203
23.3	Mixing Tcl and C debugging . . . . .	204
23.4	Memory Debugging . . . . .	205
23.4.1	Using dmalloc . . . . .	205
23.4.2	Memory Conservation Tips . . . . .	206
23.4.3	Some statistics collected by dmalloc . . . . .	206
23.5	Memory Leaks . . . . .	206
23.5.1	OTcl . . . . .	207
23.5.2	C/C++ . . . . .	207

<b>24 Mathematical Support</b>	<b>208</b>
24.1 Random Number Generation . . . . .	208
24.1.1 Seeding The RNG . . . . .	209
24.1.2 OTcl Support . . . . .	211
24.1.3 C++ Support . . . . .	212
24.2 Random Variables . . . . .	213
24.3 Integrals . . . . .	214
24.4 ns-random . . . . .	215
24.5 Some mathematical-support related objects . . . . .	216
24.6 Commands at a glance . . . . .	216
<b>25 Trace and Monitoring Support</b>	<b>218</b>
25.1 Trace Support . . . . .	218
25.1.1 OTcl Helper Functions . . . . .	219
25.2 Library support and examples . . . . .	220
25.3 The C++ Trace Class . . . . .	222
25.4 Trace File Format . . . . .	223
25.5 Packet Types . . . . .	225
25.6 Queue Monitoring . . . . .	226
25.7 Per-Flow Monitoring . . . . .	228
25.7.1 The Flow Monitor . . . . .	228
25.7.2 Flow Monitor Trace Format . . . . .	228
25.7.3 The Flow Class . . . . .	229
25.8 Commands at a glance . . . . .	229
<b>26 Test Suite Support</b>	<b>232</b>
26.1 Test Suite Components . . . . .	232
26.2 Write a Test Suite . . . . .	232
<b>27 ns Code Styles</b>	<b>235</b>
27.1 Indentation style . . . . .	235
27.2 Variable Naming Conventions . . . . .	235
27.3 Miscellaneous . . . . .	235
<b>IV Routing</b>	<b>237</b>
<b>28 Unicast Routing</b>	<b>238</b>
28.1 The Interface to the Simulation Operator (The API) . . . . .	238
28.2 Other Configuration Mechanisms for Specialised Routing . . . . .	239
28.3 Protocol Specific Configuration Parameters . . . . .	240
28.4 Internals and Architecture of Routing . . . . .	241
28.4.1 The classes . . . . .	241
28.4.2 Interface to Network Dynamics and Multicast . . . . .	245
28.5 Protocol Internals . . . . .	246
28.6 Unicast routing objects . . . . .	247
28.7 Commands at a glance . . . . .	247
<b>29 Multicast Routing</b>	<b>249</b>
29.1 Multicast API . . . . .	249
29.1.1 Multicast Behavior Monitor Configuration . . . . .	250
29.1.2 Protocol Specific configuration . . . . .	251
29.2 Internals of Multicast Routing . . . . .	252
29.2.1 The classes . . . . .	252
29.2.2 Extensions to other classes in <i>ns</i> . . . . .	254

29.2.3	Protocol Internals . . . . .	257
29.2.4	The internal variables . . . . .	259
29.3	Commands at a glance . . . . .	259
<b>30</b>	<b>Network Dynamics</b>	<b>262</b>
30.1	The user level API . . . . .	262
30.2	The Internal Architecture . . . . .	264
30.2.1	The class rtModel . . . . .	264
30.2.2	class rtQueue . . . . .	265
30.3	Interaction with Unicast Routing . . . . .	266
30.3.1	Extensions to Other Classes . . . . .	266
30.4	Deficiencies in the Current Network Dynamics API . . . . .	267
30.5	Commands at a glance . . . . .	267
<b>31</b>	<b>Hierarchical Routing</b>	<b>269</b>
31.1	Overview of Hierarchical Routing . . . . .	269
31.2	Usage of Hierarchical routing . . . . .	269
31.3	Creating large Hierarchical topologies . . . . .	271
31.4	Hierarchical Routing with SessionSim . . . . .	272
31.5	Commands at a glance . . . . .	272
<b>V</b>	<b>Transport</b>	<b>273</b>
<b>32</b>	<b>UDP Agents</b>	<b>274</b>
32.1	UDP Agents . . . . .	274
32.2	Commands at a glance . . . . .	275
<b>33</b>	<b>TCP Agents</b>	<b>276</b>
33.1	One-Way TCP Senders . . . . .	277
33.1.1	The Base TCP Sender (Tahoe TCP) . . . . .	277
33.1.2	Configuration . . . . .	277
33.1.3	Simple Configuration . . . . .	277
33.1.4	Other Configuration Parameters . . . . .	278
33.1.5	Other One-Way TCP Senders . . . . .	279
33.2	TCP Receivers (sinks) . . . . .	279
33.2.1	The Base TCP Sink . . . . .	280
33.2.2	Delayed-ACK TCP Sink . . . . .	280
33.2.3	Sack TCP Sink . . . . .	280
33.3	Two-Way TCP Agents (FullTcp) . . . . .	280
33.3.1	Simple Configuration . . . . .	281
33.3.2	BayFullTcp . . . . .	282
33.4	Architecture and Internals . . . . .	282
33.5	Tracing TCP Dynamics . . . . .	284
33.6	One-Way Trace TCP Trace Dynamics . . . . .	284
33.7	One-Way Trace TCP Trace Dynamics . . . . .	284
33.8	Commands at a glance . . . . .	284
<b>34</b>	<b>SCTP Agents</b>	<b>286</b>
34.1	The Base SCTP Agent . . . . .	286
34.1.1	Configuration Parameters . . . . .	287
34.1.2	Commands . . . . .	289
34.2	Extensions . . . . .	290
34.2.1	HbAfterRto SCTP . . . . .	290
34.2.2	MultipleFastRtx SCTP . . . . .	290



34.2.3	Timestamp Sctp	291
34.2.4	MfrHbAfterRto Sctp	291
34.2.5	MfrHbAfterRto Sctp	291
34.3	Tracing Sctp Dynamics	291
34.4	Sctp Applications	292
34.5	Example Scripts	293
34.5.1	Singled Homed Example	293
34.5.2	Multihomed Example	294
<b>35</b>	<b>Agent/SRM</b>	<b>296</b>
35.1	Confi guration	296
35.1.1	Trivial Confi guration	296
35.1.2	Other Confi guration Parameters	298
35.1.3	Statistics	299
35.1.4	Tracing	300
35.2	Architecture and Internals	302
35.3	Packet Handling: Processing received messages	302
35.4	Loss Detection—The Class SRMinfo	304
35.5	Loss Recovery Objects	304
35.6	Session Objects	306
35.7	Extending the Base Class Agent	307
35.7.1	Fixed Timers	307
35.7.2	Adaptive Timers	307
35.8	SRM objects	308
35.9	Commands at a glance	309
<b>36</b>	<b>PLM</b>	<b>311</b>
36.1	Confi guration	311
36.2	The Packet Pair Source Generator	313
36.3	Architecture of the PLM Protocol	314
36.3.1	Instantiation of a PLM Source	314
36.3.2	Instantiation of a PLM Receiver	314
36.3.3	Reception of a Packet	315
36.3.4	Detection of a Loss	316
36.3.5	Joining or Leaving a Layer	316
36.4	Commands at a Glance	316
<b>VI</b>	<b>Application</b>	<b>318</b>
<b>37</b>	<b>Applications and transport agent API</b>	<b>319</b>
37.1	The class Application	319
37.2	The transport agent API	320
37.2.1	Attaching transport agents to nodes	320
37.2.2	Attaching applications to agents	321
37.2.3	Using transport agents via system calls	321
37.2.4	Agent upcalls to applications	321
37.2.5	An example	322
37.3	The class Traffi cGenerator	323
37.3.1	An example	325
37.4	Simulated applications: Telnet and FTP	326
37.5	Applications objects	326
37.6	Commands at a glance	328

<b>38 Web cache as an application</b>	<b>329</b>
38.1 Using application-level data in <i>ns</i>	329
38.1.1 ADU	329
38.1.2 Passing data between applications	330
38.1.3 Transmitting user data over UDP	331
38.1.4 Transmitting user data over TCP	332
38.1.5 Class hierarchy related to user data handling	333
38.2 Overview of web cache classes	333
38.2.1 Managing HTTP connections	333
38.2.2 Managing web pages	334
38.2.3 Debugging	335
38.3 Representing web pages	335
38.4 Page pools	336
38.4.1 PagePool/Math	336
38.4.2 PagePool/CompMath	337
38.4.3 PagePool/ProxyTrace	337
38.4.4 PagePool/Client	338
38.4.5 PagePool/WebTraf	338
38.5 Web client	340
38.6 Web server	341
38.7 Web cache	342
38.7.1 Http/Cache	342
38.8 Putting together: a simple example	343
38.9 Http trace format	345
38.10 Commands at a glance	346
<b>39 Worm Model</b>	<b>348</b>
39.1 Overview	348
39.2 Configuration	349
39.3 Commands at a glance	349
<b>40 PackMime-HTTP: Web Traffic Generation in NS-2</b>	<b>351</b>
40.1 Implementation Details	351
40.1.1 PackMimeHTTP Client Application	352
40.1.2 PackMimeHTTP Server Application	353
40.2 PackMimeHTTP Random Variables	353
40.3 Use of DelayBox with PackMime-HTTP	353
40.4 Example	354
40.5 Commands at a Glance	355
<b>VII Scale</b>	<b>358</b>
<b>41 Session-level Packet Distribution</b>	<b>359</b>
41.1 Configuration	359
41.1.1 Basic Configuration	359
41.1.2 Inserting a Loss Module	361
41.2 Architecture	361
41.3 Internals	362
41.3.1 Object Linkage	362
41.3.2 Packet Forwarding	363
41.4 Commands at a glance	364
<b>42 Asim: approximate analytical simulation</b>	<b>365</b>

<b>VIII</b>	<b>Emulation</b>	<b>369</b>
<b>43</b>	<b>Emulation</b>	<b>370</b>
43.1	Introduction . . . . .	370
43.2	Real-Time Scheduler . . . . .	371
43.3	Tap Agents . . . . .	371
43.4	Network Objects . . . . .	372
43.4.1	Pcap/BPF Network Objects . . . . .	372
43.4.2	IP Network Objects . . . . .	373
43.4.3	IP/UDP Network Objects . . . . .	373
43.5	An Example . . . . .	374
43.6	Commands at a glance . . . . .	375
<b>IX</b>	<b>Visualization with Nam - The Network Animator</b>	<b>376</b>
<b>44</b>	<b>Nam</b>	<b>377</b>
44.1	Introduction . . . . .	377
44.2	Nam Command Line Options . . . . .	377
44.3	User Interface . . . . .	378
44.4	Keyboard Commands . . . . .	379
44.5	Generating External Animations from Nam . . . . .	380
44.6	Network Layout . . . . .	380
44.7	Animation Objects . . . . .	381
<b>45</b>	<b>Nam Trace</b>	<b>382</b>
45.1	Nam Trace Format . . . . .	382
45.1.1	Initialization Events . . . . .	383
45.1.2	Nodes . . . . .	384
45.1.3	Links . . . . .	384
45.1.4	Queues . . . . .	385
45.1.5	Packets . . . . .	385
45.1.6	Node Marking . . . . .	386
45.1.7	Agent Tracing . . . . .	387
45.1.8	Variable Tracing . . . . .	387
45.1.9	Executing Tcl Procedures and External Code from within Nam . . . . .	387
45.1.10	Using Streams for Realtime Applications . . . . .	389
45.1.11	Nam Trace File Format Lookup Table . . . . .	392
45.2	Ns commands for creating and controlling nam animations . . . . .	398
45.2.1	Node . . . . .	398
45.2.2	Link/Queue . . . . .	398
45.2.3	Agent and Features . . . . .	399
45.2.4	Some Generic Commands . . . . .	399
<b>X</b>	<b>Other</b>	<b>400</b>
<b>46</b>	<b>Educational use of NS and NAM</b>	<b>401</b>
46.1	Using NS for educational purposes . . . . .	401
46.1.1	Installing/building/running <i>ns</i> . . . . .	401
46.1.2	The educational scripts' inventory page: . . . . .	401
46.2	Using NAM for educational purposes . . . . .	402

# Chapter 1

## Introduction

Let's start at the very beginning,  
a very nice place to start,  
when you sing, you begin with A, B, C,  
when you simulate, you begin with the topology,<sup>1</sup>  
...

This document (*ns Notes and Documentation*) provides reference documentation for ns. Although we begin with a simple simulation script, resources like Marc Greis's tutorial web pages (originally at his web site, now at <http://www.isi.edu/nsnam/ns/tutorial/>) or the slides from one of the ns tutorials are probably better places to begin for the ns novice.

We first begin by showing a simple simulation script. This script is also available in the sources in `~ns/tcl/ex/simple.tcl`.

This script defines a simple topology of four nodes, and two agents, a UDP agent with a CBR traffic generator, and a TCP agent. The simulation runs for 3s. The output is two trace files, `out.tr` and `out.nam`. When the simulation completes at the end of 3s, it will attempt to run a nam visualisation of the simulation on your screen.

```
# The preamble
set ns [new Simulator]                                ; # initialise the simulation

# Predefine tracing
set f [open out.tr w]
$ns trace-all $f
set nf [open out.nam w]
$ns namtrace-all $nf
```

---

<sup>1</sup>with apologies to Rodgers and Hammerstein

*# so, we lied. now, we define the topology*

```
#
#
#      n0
#      |
#      5Mb \
#      2ms  \
#            \
#            n2 ----- n3
#            /
#      1.5Mb /
#      5Mb  /
#      2ms  /
#            /
#            n1
```

```
#
set n0 [$ns node]
set n1 [$ns node]
set n2 [$ns node]
set n3 [$ns node]
```

```
$ns duplex-link $n0 $n2 5Mb 2ms DropTail
$ns duplex-link $n1 $n2 5Mb 2ms DropTail
$ns duplex-link $n2 $n3 1.5Mb 10ms DropTail
```

*# Some agents.*

```
set udp0 [new Agent/UDP] ;# A UDP agent
$ns attach-agent $n0 $udp0 ;# on node $n0
set cbr0 [new Application/Traffic/CBR] ;# A CBR traffic generator agent
$cbr0 attach-agent $udp0 ;# attached to the UDP agent
$udp0 set class_ 0 ;# actually, the default, but...
```

```
set null0 [new Agent/Null] ;# Its sink
$ns attach-agent $n3 $null0 ;# on node $n3
```

```
$ns connect $udp0 $null0
$ns at 1.0 "$cbr0 start"
```

```
puts [$cbr0 set packetSize_]
puts [$cbr0 set interval_]
```

*# A FTP over TCP/Tahoe from \$n1 to \$n3, flowid 2*

```
set tcp [new Agent/TCP]
$tcp set class_ 1
$ns attach-agent $n1 $tcp
```

```
set sink [new Agent/TCPSink]
$ns attach-agent $n3 $sink
```

```
set ftp [new Application/FTP] ;# TCP does not generate its own traffic
$ftp attach-agent $tcp
$ns at 1.2 "$ftp start"
```

```
$ns connect $tcp $sink
$ns at 1.35 "$ns detach-agent $n0 $tcp ; $ns detach-agent $n3 $sink"
```



*# The simulation runs for 3s.*  
*# The simulation comes to an end when the scheduler invokes the finish{} procedure below.*  
*# This procedure closes all trace files, and invokes nam visualization on one of the trace files.*

```
$ns at 3.0 "finish"  
proc finish {} {  
    global ns f nf  
    $ns flush-trace  
    close $f  
    close $nf  
  
    puts "running nam..."  
    exec nam out.nam &  
    exit 0  
}
```

*# Finally, start the simulation.*  
\$ns run

## Chapter 2

# Undocumented Facilities

Ns is often growing to include new protocols. Unfortunately the documentation doesn't grow quite as often. This section lists what remains to be documented, or what needs to be improved.

(The documentation is in the doc subdirectory of the ns source code if you want to add to it. :-)

**Interface to the Interpreter**     • nothing currently

**Simulator Basics**     • LANs need to be updated for new wired/wireless support (Yuri updated this?)

- wireless support needs to be added (done)
- should explicitly list queueing options in the queue mgt chapter?

**Support**     • should pick a single list mgt package and document it

- should document the trace-post-processing utilities in bin

**Routing**     • The usage and design of link state and MPLS routing modules are not documented at all. (Note: link state and MPLS appeared only in daily snapshots and releases after 09/14/2000.)

- need to document hierarchical routing/addressing (Padma has done)
- need a chapter on supported ad-hoc routing protocols

**Queueing**     • CBQ needs documentation (can maybe build off of `ftp://ftp.ee.lbl.gov/papers/cbqsims.ps.Z`?)

**Transport**     • need to document MFTP

- need to document RTP (session-rtp.cc, etc.)
- need to document multicast building blocks
- should repair and document snoop and tcp-int

**Traffic and scenarios** (new section)

- should add a description of how to drive the simulator from traces
- should add discussion of the scenario generator
- should add discussion of http traffic sources

**Application**     • is the non-Haobo http stuff documented? no.



- Scale**     • should add disucssion of mixed mode (pending)
- Emulation**     • nothing currently
- Other**     • should document admission control policies?
- should add a validation chapter and snarf up the contents of ns-tests.html
  - should snarf up Marc Greis' tutorial rather than just referring to it?

## **Part I**

# **Interface to the Interpreter**

## Chapter 3

# OTcl Linkage

*ns* is an object oriented simulator, written in C++, with an OTcl interpreter as a frontend. The simulator supports a class hierarchy in C++ (also called the compiled hierarchy in this document), and a similar class hierarchy within the OTcl interpreter (also called the interpreted hierarchy in this document). The two hierarchies are closely related to each other; from the user's perspective, there is a one-to-one correspondence between a class in the interpreted hierarchy and one in the compiled hierarchy. The root of this hierarchy is the class `TclObject`. Users create new simulator objects through the interpreter; these objects are instantiated within the interpreter, and are closely mirrored by a corresponding object in the compiled hierarchy. The interpreted class hierarchy is automatically established through methods defined in the class `TclClass`. user instantiated objects are mirrored through methods defined in the class `TclObject`. There are other hierarchies in the C++ code and OTcl scripts; these other hierarchies are not mirrored in the manner of `TclObject`.

### 3.1 Concept Overview

*Why two languages?* *ns* uses two languages because simulator has two different kinds of things it needs to do. On one hand, detailed simulations of protocols requires a systems programming language which can efficiently manipulate bytes, packet headers, and implement algorithms that run over large data sets. For these tasks run-time speed is important and turn-around time (run simulation, find bug, fix bug, recompile, re-run) is less important.

On the other hand, a large part of network research involves slightly varying parameters or configurations, or quickly exploring a number of scenarios. In these cases, iteration time (change the model and re-run) is more important. Since configuration runs once (at the beginning of the simulation), run-time of this part of the task is less important.

*ns* meets both of these needs with two languages, C++ and OTcl. C++ is fast to run but slower to change, making it suitable for detailed protocol implementation. OTcl runs much slower but can be changed very quickly (and interactively), making it ideal for simulation configuration. *ns* (via `tclcl`) provides glue to make objects and variables appear on both languages.

For more information about the idea of scripting languages and split-language programming, see Ousterhout's article in IEEE Computer [26]. For more information about split level programming for network simulation, see the *ns* paper [2].

*Which language for what?* Having two languages raises the question of which language should be used for what purpose.

Our basic advice is to use OTcl:

- for configuration, setup, and "one-time" stuff

- if you can do what you want by manipulating existing C++ objects

and use C++:

- if you are doing *anything* that requires processing each packet of a flow
- if you have to change the behavior of an existing C++ class in ways that weren't anticipated

For example, links are OTcl objects that assemble delay, queueing, and possibly loss modules. If your experiment can be done with those pieces, great. If instead you want to do something fancier (a special queueing discipline or model of loss), then you'll need a new C++ object.

There are certainly grey areas in this spectrum: most routing is done in OTcl (although the core Dijkstra algorithm is in C++). We've had HTTP simulations where each flow was started in OTcl and per-packet processing was all in C++. This approach worked OK until we had 100s of flows starting per second of simulated time. In general, if you're ever having to invoke Tcl many times per second, you probably should move that code to C++.

## 3.2 Code Overview

In this document, we use the term “interpreter” to be synonymous with the OTcl interpreter. The code to interface with the interpreter resides in a separate directory, `tclcl`. The rest of the simulator code resides in the directory, `ns-2`. We will use the notation `~tclcl/<file>` to refer to a particular `<file>` in the `Tcl` directory. Similarly, we will use the notation, `~ns/<file>` to refer to a particular `<file>` in the `ns-2` directory.

There are a number of classes defined in `~tclcl/`. We only focus on the six that are used in `ns`: The class `Tcl` (Section 3.3) contains the methods that C++ code will use to access the interpreter. The class `TclObject` (Section 3.4) is the base class for all simulator objects that are also mirrored in the compiled hierarchy. The class `TclClass` (Section 3.5) defines the interpreted class hierarchy, and the methods to permit the user to instantiate `TclObjects`. The class `TclCommand` (Section 3.6) is used to define simple global interpreter commands. The class `EmbeddedTcl` (Section 3.7) contains the methods to load higher level builtin commands that make configuring simulations easier. Finally, the class `InstVar` (Section 3.8) contains methods to access C++ member variables as OTcl instance variables.

The procedures and functions described in this chapter can be found in `~tclcl/Tcl.{cc, h}`, `~tclcl/Tcl2.cc`, `~tclcl/tcl-object.tcl`, and, `~tclcl/tracedvar.{cc, h}`. The file `~tclcl/tcl2c++.c` is used in building `ns`, and is mentioned briefly in this chapter.

## 3.3 Class Tcl

The class `Tcl` encapsulates the actual instance of the OTcl interpreter, and provides the methods to access and communicate with that interpreter. The methods described in this section are relevant to the `ns` programmer who is writing C++ code. The class provides methods for the following operations:

- obtain a reference to the `Tcl` instance;
- invoke OTcl procedures through the interpreter;
- retrieve, or pass back results to the interpreter;
- report error situations and exit in an uniform manner; and

- store and lookup “TclObjects”.
- acquire direct access to the interpreter.

We describe each of the methods in the following subsections.

### 3.3.1 Obtain a Reference to the class Tcl instance

A single instance of the class is declared in *~tclcl/Tcl.cc* as a static member variable; the programmer must obtain a reference to this instance to access other methods described in this section. The statement required to access this instance is:

```
Tcl& tcl = Tcl::instance();
```

### 3.3.2 Invoking OTcl Procedures

There are four different methods to invoke an OTcl command through the instance, `tcl`. They differ essentially in their calling arguments. Each function passes a string to the interpreter, that then evaluates the string in a global context. These methods will return to the caller if the interpreter returns `TCL_OK`. On the other hand, if the interpreter returns `TCL_ERROR`, the methods will call `tkerror{}`. The user can overload this procedure to selectively disregard certain types of errors. Such intricacies of OTcl programming are outside the scope of this document. The next section (Section 3.3.3) describes methods to access the result returned by the interpreter.

- `tcl.eval(char* s)` invokes `Tcl_GlobalEval()` to execute *s* through the interpreter.
- `tcl.evalc(const char* s)` preserves the argument string *s*. It copies the string *s* into its internal buffer; it then invokes the previous `eval(char* s)` on the internal buffer.
- `tcl.eval()` assumes that the command is already stored in the class’ internal `bp_`; it directly invokes `tcl.eval(char* bp_)`. A handle to the buffer itself is available through the method `tcl.buffer(void)`.
- `tcl.evalf(const char* s, ...)` is a `Printf(3)` like equivalent. It uses `vsprintf(3)` internally to create the input string.

As an example, here are some of the ways of using the above methods:

```
Tcl& tcl = Tcl::instance();
char wrk[128];
strcpy(wrk, "Simulator set NumberInterfaces_ 1");
tcl.eval(wrk);

sprintf(tcl.buffer(), "Agent/SRM set requestFunction_ %s", "Fixed");
tcl.eval();

tcl.evalc("puts stdout hello world");
tcl.evalf("%s request %d %d", name_, sender, msgid);
```

### 3.3.3 Passing Results to/from the Interpreter

When the interpreter invokes a C++ method, it expects the result back in the private member variable, `tcl_>result`. Two methods are available to set this variable.

- `tcl.result(const char* s)`  
Pass the result string *s* back to the interpreter.
- `tcl.resultf(const char* fmt, ...)`  
`varargs(3)` variant of above to format the result using `vsprintf(3)`, pass the result string back to the interpreter.

```
if (strcmp(argv[1], "now") == 0) {
    tcl.resultf("%.17g", clock());
    return TCL_OK;
}
tcl.result("Invalid operation specified");
return TCL_ERROR;
```

Likewise, when a C++ method invokes an OTcl command, the interpreter returns the result in `tcl_>result`.

- `tcl.result(void)` must be used to retrieve the result. Note that the result is a string, that must be converted into an internal format appropriate to the type of result.

```
tcl.evalc("Simulator set NumberInterfaces_");
char* ni = tcl.result();
if (atoi(ni) != 1)
    tcl.evalc("Simulator set NumberInterfaces_ 1");
```

### 3.3.4 Error Reporting and Exit

This method provides a uniform way to report errors in the compiled code.

- `tcl.error(const char* s)` performs the following functions: write *s* to stdout; write `tcl_>result` to stdout; exit with error code 1.

```
tcl.resultf("cmd = %s", cmd);
tcl.error("invalid command specified");
/*NOTREACHED*/
```

Note that there are minor differences between returning `TCL_ERROR` as we did in the previous subsection (Section 3.3.3), and calling `Tcl::error()`. The former generates an exception within the interpreter; the user can trap the exception and possibly recover from the error. If the user has not specified any traps, the interpreter will print a stack trace and exit. However, if the code invokes `error()`, then the simulation user cannot trap the error; in addition, *ns* will not print any stack trace.

### 3.3.5 Hash Functions within the Interpreter

*ns* stores a reference to every `TclObject` in the compiled hierarchy in a hash table; this permits quick access to the objects. The hash table is internal to the interpreter. *ns* uses the name of the `TclObject` as the key to enter, lookup, or delete the `TclObject` in the hash table.

- `tcl.enter(TclObject* o)` will insert a pointer to the `TclObject o` into the hash table.  
It is used by `TclClass::create_shadow()` to insert an object into the table, when that object is created.
- `tcl.lookup(char* s)` will retrieve the `TclObject` with the name `s`.  
It is used by `TclObject::lookup()`.
- `tcl.remove(TclObject* o)` will delete references to the `TclObject o` from the hash table.  
It is used by `TclClass::delete_shadow()` to remove an existing entry from the hash table, when that object is deleted.

These functions are used internally by the class `TclObject` and class `TclClass`.

### 3.3.6 Other Operations on the Interpreter

If the above methods are not sufficient, then we must acquire the handle to the interpreter, and write our own functions.

- `tcl.interp(void)` returns the handle to the interpreter that is stored within the class `Tcl`.

## 3.4 Class `TclObject`

class `TclObject` is the base class for most of the other classes in the interpreted and compiled hierarchies. Every object in the class `TclObject` is created by the user from within the interpreter. An equivalent shadow object is created in the compiled hierarchy. The two objects are closely associated with each other. The class `TclClass`, described in the next section, contains the mechanisms that perform this shadowing.

In the rest of this document, we often refer to an object as a `TclObject`<sup>1</sup>. By this, we refer to a particular object that is either in the class `TclObject`, or in a class that is derived from the class `TclObject`. If it is necessary, we will explicitly qualify whether that object is an object within the interpreter, or an object within the compiled code. In such cases, we will use the abbreviations “interpreted object”, and “compiled object” to distinguish the two. and within the compiled code respectively.

**Differences from *ns v1*** Unlike *ns v1*, the class `TclObject` subsumes the earlier functions of the `NsObject` class. It therefore stores the interface variable bindings (Section 3.4.2) that tie OTcl instance variables in the interpreted object to corresponding C++ member variables in the compiled object. The binding is stronger than in *ns v1* in that any changes to the OTcl variables are trapped, and the current C++ and OTcl values are made consistent after each access through the interpreter. The consistency is done through the class `InstVar` (Section 3.8). Also unlike *ns v1*, objects in the class `TclObject` are no longer stored as a global link list. Instead, they are stored in a hash table in the class `Tcl` (Section 3.3.5).

**Example configuration of a `TclObject`** The following example illustrates the configuration of an SRM agent (class `Agent/SRM/Adaptive`).

```
set srm [new Agent/SRM/Adaptive]
$srms set packetSize_ 1024
$srms traffic-source $s0
```

---

<sup>1</sup>In the latest release of *ns* and *ns/tclcl*, this object has been renamed to `SplitObject`, which more accurately reflects its nature of existence. However, for the moment, we will continue to use the term `TclObject` to refer to these objects and this class.

By convention in *ns*, the class Agent/SRM/Adaptive is a subclass of Agent/SRM, is a subclass of Agent, is a subclass of TclObject. The corresponding compiled class hierarchy is the ASRMAgent, derived from SRMAgent, derived from Agent, derived from TclObject respectively. The first line of the above example shows how a TclObject is created (or destroyed) (Section 3.4.1); the next line configures a bound variable (Section 3.4.2); and finally, the last line illustrates the interpreted object invoking a C++ method as if they were an instance procedure (Section 3.4.4).

### 3.4.1 Creating and Destroying TclObjects

When the user creates a new TclObject, using the procedures `new{}` and `delete{}`; these procedures are defined in `~tclcl/tcl-object.tcl`. They can be used to create and destroy objects in all classes, including TclObjects.<sup>2</sup> In this section, we describe the internal actions executed when a TclObject is created.

**Creating TclObjects** By using `new{}`, the user creates an interpreted TclObject. the interpreter will execute the constructor for that object, `init{}`, passing it any arguments provided by the user. *ns* is responsible for automatically creating the compiled object. The shadow object gets created by the base class TclObject's constructor. Therefore, the constructor for the new TclObject must call the parent class constructor first. `new{}` returns a handle to the object, that can then be used for further operations upon that object.

The following example illustrates the Agent/SRM/Adaptive constructor:

```
Agent/SRM/Adaptive instproc init args {
    eval $self next $args
    $self array set closest_ "requestor 0 repairor 0"
    $self set eps_ [$class set eps_]
}
```

The following sequence of actions are performed by the interpreter as part of instantiating a new TclObject. For ease of exposition, we describe the steps that are executed to create an Agent/SRM/Adaptive object. The steps are:

1. Obtain an unique handle for the new object from the TclObject name space. The handle is returned to the user. Most handles in *ns* have the form `_o<NNN>`, where `<NNN>` is an integer. This handle is created by `getid{}`. It can be retrieved from C++ with the `name( ) {}` method.
2. Execute the constructor for the new object. Any user-specified arguments are passed as arguments to the constructor. This constructor must invoke the constructor associated with its parent class.

In our example above, the Agent/SRM/Adaptive calls its parent class in the very first line.

Note that each constructor, in turn invokes its parent class' constructor *ad nauseum*. The last constructor in *ns* is the TclObject constructor. This constructor is responsible for setting up the shadow object, and performing other initializations and bindings, as we explain below. *It is preferable to call the parent constructors first before performing the initializations required in this class.* This allows the shadow objects to be set up, and the variable bindings established.

3. The TclObject constructor invokes the instance procedure `create-shadow{}` for the class Agent/SRM/Adaptive.
4. When the shadow object is created, *ns* calls all of the constructors for the compiled object, each of which may establish variable bindings for objects in that class, and perform other necessary initializations. Hence our earlier injunction that it is preferable to invoke the parent constructors prior to performing the class initializations.
5. After the shadow object is successfully created, `create_shadow(void)`

---

<sup>2</sup>As an example, the classes Simulator, Node, Link, or rtObject, are classes that are *not* derived from the class TclObject. Objects in these classes are not, therefore, TclObjects. However, a Simulator, Node, Link, or route Object is also instantiated using the `new` procedure in *ns*.



- (a) adds the new object to hash table of TclObjects described earlier (Section 3.3.5).
- (b) makes `cmd{ }` an instance procedure of the newly created interpreted object. This instance procedure invokes the `command()` method of the compiled object. In a later subsection (Section 3.4.4), we describe how the `command` method is defined, and invoked.

Note that all of the above shadowing mechanisms only work when the user creates a new TclObject through the interpreter. It will not work if the programmer creates a compiled TclObject unilaterally. Therefore, the programmer is enjoined not to use the C++ new method to create compiled objects directly.

**Deletion of TclObjects** The `delete` operation destroys the interpreted object, and the corresponding shadow object. For example, `use-scheduler{<scheduler>}` uses the `delete` procedure to remove the default list scheduler, and instantiate an alternate scheduler in its place.

```

Simulator instproc use-scheduler type {
    $self instvar scheduler_

    delete scheduler_                                ; # first delete the existing list scheduler
    set scheduler_ [new Scheduler/$type]
}

```

As with the constructor, the object destructor must call the destructor for the parent class explicitly as the very last statement of the destructor. The TclObject destructor will invoke the instance procedure `delete-shadow`, that in turn invokes the equivalent compiled method to destroy the shadow object. The interpreter itself will destroy the interpreted object.

### 3.4.2 Variable Bindings

In most cases, access to compiled member variables is restricted to compiled code, and access to interpreted member variables is likewise confined to access via interpreted code; however, it is possible to establish bi-directional bindings such that both the interpreted member variable and the compiled member variable access the same data, and changing the value of either variable changes the value of the corresponding paired variable to same value.

The binding is established by the compiled constructor when that object is instantiated; it is automatically accessible by the interpreted object as an instance variable. *ns* supports five different data types: reals, bandwidth valued variables, time valued variables, integers, and booleans. The syntax of how these values can be specified in OTcl is different for each variable type.

- Real and Integer valued variables are specified in the “normal” form. For example,

```

$object set realvar 1.2e3
$object set intvar 12

```

- Bandwidth is specified as a real value, optionally suffixed by a ‘k’ or ‘K’ to mean kilo-quantities, or ‘m’ or ‘M’ to mean mega-quantities. A final optional suffix of ‘B’ indicates that the quantity expressed is in Bytes per second. The default is bandwidth expressed in bits per second. For example, all of the following are equivalent:

```

$object set bwvar 1.5m
$object set bwvar 1.5mb
$object set bwvar 1500k

```

```

$object set bwvar 1500kb
$object set bwvar .1875MB
$object set bwvar 187.5kB
$object set bwvar 1.5e6

```

- Time is specified as a real value, optionally suffixed by a ‘m’ to express time in milli-seconds, ‘n’ to express time in nano-seconds, or ‘p’ to express time in pico-seconds. The default is time expressed in seconds. For example, all of the following are equivalent:

```

$object set timevar 1500m
$object set timevar 1.5
$object set timevar 1.5e9n
$object set timevar 1500e9p

```

Note that we can also safely add a *s* to reflect the time unit of seconds. *ns* will ignore anything other than a valid real number specification, or a trailing ‘m’, ‘n’, or ‘p’.

- Booleans can be expressed either as an integer, or as ‘T’ or ‘t’ for true. Subsequent characters after the first letter are ignored. If the value is neither an integer, nor a true value, then it is assumed to be false. For example,

```

$object set boolvar t                                ; # set to true
$object set boolvar true
$object set boolvar 1                                ; # or any non-zero value

$object set boolvar false                            ; # set to false
$object set boolvar junk
$object set boolvar 0

```

The following example shows the constructor for the ASRMAgent<sup>3</sup>.

```

ASRMAgent::ASRMAgent() {
    bind("pdistance_", &pdistance_);           /* real variable */
    bind("requestor_", &requestor_);           /* integer variable */
    bind_time("lastSent_", &lastSessSent_);    /* time variable */
    bind_bw("ctrlLimit_", &ctrlBWLlimit_);     /* bandwidth variable */
    bind_bool("running_", &running_);          /* boolean variable */
}

```

Note that all of the functions above take two arguments, the name of an OTcl variable, and the address of the corresponding compiled member variable that is linked. While it is often the case that these bindings are established by the constructor of the object, it need not always be done in this manner. We will discuss such alternate methods when we describe the class InstVar (Section 3.8) in detail later.

Each of the variables that is bound is automatically initialised with default values when the object is created. The default values are specified as interpreted class variables. This initialisation is done by the routine `init-instvar{}`, invoked by methods in the class Instvar, described later (Section 3.8). `init-instvar{}` checks the class of the interpreted object, and all of the parent class of that object, to find the first class in which the variable is defined. It uses the value of the variable in that class to initialise the object. Most of the bind initialisation values are defined in `~ns/tcl/lib/ns-default.tcl`.

For example, if the following class variables are defined for the ASRMAgent:

---

<sup>3</sup>Note that this constructor is embellished to illustrate the features of the variable binding mechanism.

```

Agent/SRM/Adaptive set pdistance_ 15.0
Agent/SRM set pdistance_ 10.0
Agent/SRM set lastSent_ 8.345m
Agent set ctrlLimit_ 1.44M
Agent/SRM/Adaptive set running_ f

```

Therefore, every new Agent/SRM/Adaptive object will have `pdistance_` set to 15.0; `lastSent_` is set to 8.345m from the setting of the class variable of the parent class; `ctrlLimit_` is set to 1.44M using the class variable of the parent class twice removed; `running` is set to false; the instance variable `pdistance_` is not initialised, because no class variable exists in any of the class hierarchy of the interpreted object. In such instance, `init-instvar{ }` will invoke `warn-instvar{ }`, to print out a warning about such a variable. The user can selectively override this procedure in their simulation scripts, to elide this warning.

Note that the actual binding is done by instantiating objects in the class `InstVar`. Each object in the class `InstVar` binds one compiled member variable to one interpreted member variable. A `TclObject` stores a list of `InstVar` objects corresponding to each of its member variable that is bound in this fashion. The head of this list is stored in its member variable `instvar_` of the `TclObject`.

One last point to consider is that *ns* will guarantee that the actual values of the variable, both in the interpreted object and the compiled object, will be identical at all times. However, if there are methods and other variables of the compiled object that track the value of this variable, they must be explicitly invoked or changed whenever the value of this variable is changed. This usually requires additional primitives that the user should invoke. One way of providing such primitives in *ns* is through the `command()` method described in the next section.

### 3.4.3 Variable Tracing

In addition to variable bindings, `TclObject` also supports tracing of both C++ and Tcl instance variables. A traced variable can be created and configured either in C++ or Tcl. To establish variable tracing at the Tcl level, the variable must be visible in Tcl, which means that it must be a bounded C++/Tcl or a pure Tcl instance variable. In addition, the object that owns the traced variable is also required to establish tracing using the Tcl `trace` method of `TclObject`. The first argument to the `trace` method must be the name of the variable. The optional second argument specifies the trace object that is responsible for tracing that variable. If the trace object is not specified, the object that owns the variable is responsible for tracing it.

For a `TclObject` to trace variables, it must extend the C++ `trace` method that is virtually defined in `TclObject`. The `Trace` class implements a simple `trace` method, thereby, it can act as a generic tracer for variables.

```

class Trace : public Connector {
    ...
    virtual void trace(TracedVar*);
};

```

Below is a simple example for setting up variable tracing in Tcl:

```

# \tcp tracing its own variable cwnd_
\tcp trace cwnd_

# the variable ssthresh_ of \tcp is traced by a generic \stracer
set tracer [new Trace/Var]
\tcp trace ssthresh_ \stracer

```

For a C++ variable to be traceable, it must belong to a class that derives from TracedVar. The virtual base class TracedVar keeps track of the variable's name, owner, and tracer. Classes that derives from TracedVar must implement the virtual method value, that takes a character buffer as an argument and writes the value of the variable into that buffer.

```
class TracedVar {
    ...
    virtual char* value(char* buf) = 0;
protected:
    TracedVar(const char* name);
    const char* name_;      // name of the variable
    TclObject* owner_;      // the object that owns this variable
    TclObject* tracer_;     // callback when the variable is changed
    ...
};
```

The TclCL library exports two classes of TracedVar: TracedInt and TracedDouble. These classes can be used in place of the basic type int and double respectively. Both TracedInt and TracedDouble overload all the operators that can change the value of the variable such as assignment, increment, and decrement. These overloaded operators use the assign method to assign the new value to the variable and call the tracer if the new value is different from the old one. TracedInt and TracedDouble also implement their value methods that output the value of the variable into string. The width and precision of the output can be pre-specified.

### 3.4.4 command Methods: Definition and Invocation

For every TclObject that is created, *ns* establishes the instance procedure, cmd{ }, as a hook to executing methods through the compiled shadow object. The procedure cmd{ } invokes the method command() of the shadow object automatically, passing the arguments to cmd{ } as an argument vector to the command() method.

The user can invoke the cmd{ } method in one of two ways: by explicitly invoking the procedure, specifying the desired operation as the first argument, or implicitly, as if there were an instance procedure of the same name as the desired operation. Most simulation scripts will use the latter form, hence, we will describe that mode of invocation first.

Consider the that the distance computation in SRM is done by the compiled object; however, it is often used by the interpreted object. It is usually invoked as:

```
$srmObject distance? <agentAddress>
```

If there is no instance procedure called distance?, the interpreter will invoke the instance procedure unknown{ }, defined in the base class TclObject. The unknown procedure then invokes

```
$srmObject cmd distance? <agentAddress>
```

to execute the operation through the compiled object's command() procedure.

Ofcourse, the user could explicitly invoke the operation directly. One reason for this might be to overload the operation by using an instance procedure of the same name. For example,

```
Agent/SRM/Adaptive instproc distance? addr {
```

```

$self instvar distanceCache_
if ![info exists distanceCache_($addr)] {
    set distanceCache_($addr) [$self cmd distance? $addr]
}
set distanceCache_($addr)
}

```

We now illustrate how the `command()` method using `ASRMAgent::command()` as an example.

```

int ASRMAgent::command(int argc, const char*const*argv) {
    Tcl& tcl = Tcl::instance();
    if (argc == 3) {
        if (strcmp(argv[1], "distance?") == 0) {
            int sender = atoi(argv[2]);
            SRMInfo* sp = get_state(sender);
            tcl.resultf("%f", sp->distance_);
            return TCL_OK;
        }
    }
    return (SRMAgent::command(argc, argv));
}

```

We can make the following observations from this piece of code:

- The function is called with two arguments:  
 The first argument (`argc`) indicates the number of arguments specified in the command line to the interpreter.  
 The command line arguments vector (`argv`) consists of  
 — `argv[0]` contains the name of the method, “cmd”.  
 — `argv[1]` specifies the desired operation.  
 — If the user specified any arguments, then they are placed in `argv[2... (argc - 1)]`.  
 The arguments are passed as strings; they must be converted to the appropriate data type.
- If the operation is successfully matched, the match should return the result of the operation using methods described earlier (Section 3.3.3).
- `command()` itself must return either `TCL_OK` or `TCL_ERROR` to indicate success or failure as its return code.
- If the operation is not matched in this method, it must invoke its parent’s `command` method, and return the corresponding result.  
 This permits the user to conceive of operations as having the same inheritance properties as instance procedures or compiled methods.  
 In the event that this `command` method is defined for a class with multiple inheritance, the programmer has the liberty to choose one of two implementations:
  - 1) Either they can invoke one of the parent’s `command` method, and return the result of that invocation, or
  - 2) They can each of the parent’s `command` methods in some sequence, and return the result of the first invocation that is successful. If none of them are successful, then they should return an error.

In our document, we call operations executed through the `command()` *instproc*-likes. This reflects the usage of these operations as if they were OTcl instance procedures of an object, but can be very subtly different in their realisation and usage.

### 3.5 Class TclClass

This compiled class (`class TclClass`) is a pure virtual class. Classes derived from this base class provide two functions: construct the interpreted class hierarchy to mirror the compiled class hierarchy; and provide methods to instantiate new `TclObjects`. Each such derived class is associated with a particular compiled class in the compiled class hierarchy, and can instantiate new objects in the associated class.

As an example, consider a class such as the class `RenoTcpClass`. It is derived from class `TclClass`, and is associated with the class `RenoTcpAgent`. It will instantiate new objects in the class `RenoTcpAgent`. The compiled class hierarchy for `RenoTcpAgent` is that it derives from `TcpAgent`, that in turn derives from `Agent`, that in turn derives (roughly) from `TclObject`. `RenoTcpClass` is defined as

```
static class RenoTcpClass: public TclClass {
public:
    RenoTcpClass() : TclClass("Agent/TCP/Reno") {}
    TclObject* create(int argc, const char*const* argv) {
        return (new RenoTcpAgent());
    }
} class_reno;
```

We can make the following observations from this definition:

1. The class defines only the constructor, and one additional method, to create instances of the associated `TclObject`.
2. *ns* will execute the `RenoTcpClass` constructor for the static variable `class_reno`, when it is first started. This sets up the appropriate methods and the interpreted class hierarchy.
3. The constructor specifies the interpreted class explicitly as `Agent/TCP/Reno`. This also specifies the interpreted class hierarchy implicitly.

Recall that the convention in *ns* is to use the character slash (`'/'`) as a separator. For any given class `A/B/C/D`, the class `A/B/C/D` is a sub-class of `A/B/C`, that is itself a sub-class of `A/B`, that, in turn, is a sub-class of `A`. `A` itself is a sub-class of `TclObject`.

In our case above, the `TclClass` constructor creates three classes, `Agent/TCP/Reno` sub-class of `Agent/TCP` sub-class of `Agent` sub-class of `TclObject`.

4. This class is associated with the class `RenoTcpAgent`; it creates new objects in this associated class.
5. The `RenoTcpClass::create` method returns `TclObjects` in the class `RenoTcpAgent`.
6. When the user specifies new `Agent/TCP/Reno`, the routine `RenoTcpClass::create` is invoked.
7. The arguments vector (`argv`) consists of
  - `argv[0]` contains the name of the object.
  - `argv[1...3]` contain `$self`, `$class`, and `$proc`. Since `create` is called through the instance procedure `create-shadow`, `argv[3]` contains `create-shadow`.
  - `argv[4]` contain any additional arguments (passed as a string) provided by the user.

The class `Trace` illustrates argument handling by `TclClass` methods.

```
class TraceClass : public TclClass {
public:
```

```

        TraceClass() : TclClass("Trace") {}
        TclObject* create(int args, const char*const* argv) {
            if (args >= 5)
                return (new Trace(*argv[4]));
            else
                return NULL;
        }
    } trace_class;

```

A new Trace object is created as

```
new Trace "X"
```

Finally, the nitty-gritty details of how the interpreted class hierarchy is constructed:

1. The object constructor is executed when *ns* first starts.
2. This constructor calls the `TclClass` constructor with the name of the interpreted class as its argument.
3. The `TclClass` constructor stores the name of the class, and inserts this object into a linked list of the `TclClass` objects.
4. During initialization of the simulator, `Tcl_AppInit(void)` invokes `TclClass::bind(void)`
5. For each object in the list of `TclClass` objects, `bind()` invokes `register{}`, specifying the name of the interpreted class as its argument.
6. `register{}` establishes the class hierarchy, creating the classes that are required, and not yet created.
7. Finally, `bind()` defines instance procedures `create-shadow` and `delete-shadow` for this new class.

### 3.5.1 How to Bind Static C++ Class Member Variables

In Section 3.4, we have seen how to expose member variables of a C++ object into OTcl space. This, however, does not apply to static member variables of a C++ class. Of course, one may create an OTcl variable for the static member variable of every C++ object; obviously this defeats the whole meaning of static members.

We cannot solve this binding problem using a similar solution as binding in `TclObject`, which is based on `InstVar`, because `InstVars` in `TclCL` require the presence of a `TclObject`. However, we can create a method of the corresponding `TclClass` and access static members of a C++ class through the methods of its corresponding `TclClass`. The procedure is as follows:

1. Create your own derived `TclClass` as described above;
2. Declare methods `bind()` and `method()` in your derived class;
3. Create your binding methods in the implementation of your `bind()` with `add_method("your_method")`, then implement the handler in `method()` in a similar way as you would do in `TclObject::command()`. Notice that the number of arguments passed to `TclClass::method()` are different from those passed to `TclObject::command()`. The former has two more arguments in the front.

As an example, we show a simplified version of `PacketHeaderClass` in `~ns/packet.cc`. Suppose we have the following class `Packet` which has a static variable `hdrlen_` that we want to access from OTcl:

```

class Packet {
    .....
    static int hdrlen_;
};

```

Then we do the following to construct an accessor for this variable:

```

class PacketHeaderClass : public TclClass {
protected:
    PacketHeaderClass(const char* classname, int hdrsize);
    TclObject* create(int argc, const char*const* argv);
                                /* These two implements OTcl class access methods */
    virtual void bind();
    virtual int method(int argc, const char*const* argv);
};

void PacketHeaderClass::bind()
{
                                /* Call to base class bind() must precede add_method() */
    TclClass::bind();
    add_method("hdrlen");
}

int PacketHeaderClass::method(int ac, const char*const* av)
{
    Tcl& tcl = Tcl::instance();
                                /* Notice this argument translation; we can then handle them as if in TclObject::command() */
    int argc = ac - 2;
    const char*const* argv = av + 2;
    if (argc == 2) {
        if (strcmp(argv[1], "hdrlen") == 0) {
            tcl.resultf("%d", Packet::hdrlen_);
            return (TCL_OK);
        }
    } else if (argc == 3) {
        if (strcmp(argv[1], "hdrlen") == 0) {
            Packet::hdrlen_ = atoi(argv[2]);
            return (TCL_OK);
        }
    }
    return TclClass::method(ac, av);
}

```

After this, we can then use the following OTcl command to access and change values of `Packet::hdrlen_`:

```

PacketHeader hdrlen 120
set i [PacketHeader hdrlen]

```



## 3.6 Class TclCommand

This class (`class TclCommand`) provides just the mechanism for *ns* to export simple commands to the interpreter, that can then be executed within a global context by the interpreter. There are two functions defined in `~ns/misc.cc`: `ns-random` and `ns-version`. These two functions are initialized by the function `init_misc(void)`, defined in `~ns/misc.cc`; `init_misc` is invoked by `Tcl_AppInit(void)` during startup.

- `class VersionCommand` defines the command `ns-version`. It takes no argument, and returns the current *ns* version string.

```
% ns-version                                ; # get the current version
2.0a12
```

- `class RandomCommand` defines the command `ns-random`. With no argument, `ns-random` returns an integer, uniformly distributed in the interval  $[0, 2^{31} - 1]$ .

When specified an argument, it takes that argument as the seed. If this seed value is 0, the command uses a heuristic seed value; otherwise, it sets the seed for the random number generator to the specified value.

```
% ns-random                                ; # return a random number
2078917053
% ns-random 0                              ; # set the seed heuristically
858190129
% ns-random 23786                          ; # set seed to specified value
23786
```

*Note that, it is generally not advisable to construct top-level commands that are available to the user. We now describe how to define a new command using the example `class say_hello`. The example defines the command `hi`, to print the string “hello world”, followed by any command line arguments specified by the user. For example,*

```
% hi this is ns [ns-version]
hello world, this is ns 2.0a12
```

1. The command must be defined within a class derived from the `class TclCommand`. The class definition is:

```
class say_hello : public TclCommand {
public:
    say_hello();
    int command(int argc, const char*const* argv);
};
```

2. The constructor for the class must invoke the `TclCommand` constructor with the command as argument; *i.e.*,

```
say_hello() : TclCommand("hi") {}
```

The `TclCommand` constructor sets up “hi” as a global procedure that invokes `TclCommand::dispatch_cmd()`.

3. The method `command()` must perform the desired action.

The method is passed two arguments. The first argument, `argc`, contains the number of actual arguments passed by the user.

The actual arguments passed by the user are passed as an argument vector (`argv`) and contains the following:

— `argv[0]` contains the name of the command (`hi`).

— `argv[1... (argc - 1)]` contains additional arguments specified on the command line by the user.

`command()` is invoked by `dispatch_cmd()`.

```
#include <streams.h>                                     /* because we are using stream I/O */

int say_hello::command(int argc, const char*const* argv) {
    cout << "hello world:";
    for (int i = 1; i < argc; i++)
        cout << ' ' << argv[i];
    cout << '\n';
    return TCL_OK;
}
```

4. Finally, we require an instance of this class. `TclCommand` instances are created in the routine `init_misc(void)`.

```
new say_hello;
```

Note that there used to be more functions such as `ns-at` and `ns-now` that were accessible in this manner. Most of these functions have been subsumed into existing classes. In particular, `ns-at` and `ns-now` are accessible through the scheduler `TclObject`. These functions are defined in `~ns/tcl/lib/ns-lib.tcl`.

```
% set ns [new Simulator]                                ; # get new instance of simulator
_ol
% $ns now                                                ; # query simulator for current time
0
% $ns at ...                                            ; # specify at operations for simulator
...
```

### 3.7 Class `EmbeddedTcl`

`ns` permits the development of functionality in either compiled code, or through interpreter code, that is evaluated at initialization. For example, the scripts `~tcl/tcl-object.tcl` or the scripts in `~ns/tcl/lib`. Such loading and evaluation of scripts is done through objects in the class `EmbeddedTcl`.

The easiest way to extend `ns` is to add OTcl code to either `~tcl/tcl-object.tcl` or through scripts in the `~ns/tcl/lib` directory. Note that, in the latter case, `ns` sources `~ns/tcl/lib/ns-lib.tcl` automatically, and hence the programmer must add a couple of lines to this file so that their script will also get automatically sourced by `ns` at startup. As an example, the file `~ns/tcl/mcast/srm.tcl` defines some of the instance procedures to run SRM. In `~ns/tcl/lib/ns-lib.tcl`, we have the lines:

```
source tcl/mcast/srm.tcl
```

to automatically get `srm.tcl` sourced by `ns` at startup.

Three points to note with `EmbeddedTcl` code are that firstly, if the code has an error that is caught during the eval, then `ns` will not run. Secondly, the user can explicitly override any of the code in the scripts. In particular, they can re-source the entire

script after making their own changes. Finally, after adding the scripts to `~ns/tcl/lib/ns-lib.tcl`, and every time thereafter that they change their script, the user must recompile *ns* for their changes to take effect. Of course, in most cases<sup>4</sup>, the user can source their script to override the embedded code.

The rest of this subsection illustrate how to integrate individual scripts directly into *ns*. The first step is convert the script into an EmbeddedTcl object. The lines below expand `ns-lib.tcl` and create the EmbeddedTcl object instance called `et_ns_lib`:

```
tclsh bin/tcl-expand.tcl tcl/lib/ns-lib.tcl | \
    ../Tcl/tcl2c++ et_ns_lib > gen/ns_tcl.cc
```

The script, `~ns/bin/tcl-expand.tcl` expands `ns-lib.tcl` by replacing all source lines with the corresponding source files. The program, `~tcl/tcl2c.c`, converts the OTcl code into an equivalent EmbeddedTcl object, `et_ns_lib`.

During initialization, invoking the method `EmbeddedTcl::load` explicitly evaluates the array.

- `~tcl/tcl-object.tcl` is evaluated by the method `Tcl::init(void)`; `Tcl_AppInit()` invokes `Tcl::Init()`. The exact command syntax for the load is:

```
et_tclobject.load();
```

- Similarly, `~ns/tcl/lib/ns-lib.tcl` is evaluated directly by `Tcl_AppInit` in `~ns/ns_tclsh.cc`.

```
et_ns_lib.load();
```

## 3.8 Class InstVar

This section describes the internals of the class `InstVar`. This class defines the methods and mechanisms to bind a C++ member variable in the compiled shadow object to a specified OTcl instance variable in the equivalent interpreted object. The binding is set up such that the value of the variable can be set or accessed either from within the interpreter, or from within the compiled code at all times.

There are five instance variable classes: `class InstVarReal`, `class InstVarTime`, `class InstVarBandwidth`, `class InstVarInt`, and `class InstVarBool`, corresponding to bindings for real, time, bandwidth, integer, and boolean valued variables respectively.

We now describe the mechanism by which instance variables are set up. We use the class `InstVarReal` to illustrate the concept. However, this mechanism is applicable to all five types of instance variables.

When setting up an interpreted variable to access a member variable, the member functions of the class `InstVar` assume that they are executing in the appropriate method execution context; therefore, they do not query the interpreter to determine the context in which this variable must exist.

In order to guarantee the correct method execution context, a variable must only be bound if its class is already established within the interpreter, and the interpreter is currently operating on an object in that class. Note that the former requires that when a method in a given class is going to make its variables accessible via the interpreter, there must be an associated

---

<sup>4</sup>The few places where this might not work are when certain variables might have to be defined or undefined, or otherwise the script contains code other than procedure and variable definitions and executes actions directly that might not be reversible.

class `TclClass` (Section 3.5) defined that identifies the appropriate class hierarchy to the interpreter. The appropriate method execution context can therefore be created in one of two ways.

An implicit solution occurs whenever a new `TclObject` is created within the interpreter. This sets up the method execution context within the interpreter. When the compiled shadow object of the interpreted `TclObject` is created, the constructor for that compiled object can bind its member variables of that object to interpreted instance variables in the context of the newly created interpreted object.

An explicit solution is to define a `bind-variables` operation within a `command` function, that can then be invoked via the `cmd` method. The correct method execution context is established in order to execute the `cmd` method. Likewise, the compiled code is now operating on the appropriate shadow object, and can therefore safely bind the required member variables.

An instance variable is created by specifying the name of the interpreted variable, and the address of the member variable in the compiled object. The constructor for the base class `InstVar` creates an instance of the variable in the interpreter, and then sets up a trap routine to catch all accesses to the variable through the interpreter.

Whenever the variable is read through the interpreter, the trap routine is invoked just prior to the occurrence of the read. The routine invokes the appropriate `get` function that returns the current value of the variable. This value is then used to set the value of the interpreted variable that is then read by the interpreter.

Likewise, whenever the variable is set through the interpreter, the trap routine is invoked just after the write is completed. The routine gets the current value set by the interpreter, and invokes the appropriate `set` function that sets the value of the compiled member to the current value set within the interpreter.

## **Part II**

# **Simulator Basics**

## Chapter 4

# The Class Simulator

The overall simulator is described by a Tcl class `Simulator`. It provides a set of interfaces for configuring a simulation and for choosing the type of event scheduler used to drive the simulation. A simulation script generally begins by creating an instance of this class and calling various methods to create nodes, topologies, and configure other aspects of the simulation. A subclass of `Simulator` called `OldSim` is used to support *ns v1* backward compatibility.

The procedures and functions described in this chapter can be found in `~ns/tcl/lib/ns-lib.tcl`, `~ns/scheduler.{cc,h}`, and, `~ns/heap.h`.

### 4.1 Simulator Initialization

When a new simulation object is created in tcl, the initialization procedure performs the following operations:

- initialize the packet format (calls `create_packetformat`)
- create a scheduler (defaults to a calendar scheduler)
- create a “null agent” (a discard sink used in various places)

The packet format initialization sets up field offsets within packets used by the entire simulation. It is described in more detail in the following chapter on packets (Chapter 12). The scheduler runs the simulation in an event-driven manner and may be replaced by alternative schedulers which provide somewhat different semantics (see the following section for more detail). The null agent is created with the following call:

```
set nullAgent_ [new Agent/Null]
```

This agent is generally useful as a sink for dropped packets or as a destination for packets that are not counted or recorded.

### 4.2 Schedulers and Events

The simulator is an event-driven simulator. There are presently four schedulers available in the simulator, each of which is implemented using a different data structure: a simple linked-list, heap, calendar queue (default), and a special type called

“real-time”. Each of these are described below. The scheduler runs by selecting the next earliest event, executing it to completion, and returning to execute the next event. Unit of time used by scheduler is seconds. Presently, the simulator is single-threaded, and only one event in execution at any given time. If more than one event are scheduled to execute at the same time, their execution is performed on the first scheduled – first dispatched manner. Simultaneous events are not re-ordered anymore by schedulers (as it was in earlier versions) and all schedulers should yield the same order of dispatching given the same input.

No partial execution of events or pre-emption is supported.

An *event* generally comprises a “fi ring time” and a handler function. The actual defi nition of an event is found in `~ns/scheduler.h`:

```
class Event {
public:
    Event* next_;                                /* event list */
    Handler* handler_;                            /* handler to call when event ready */
    double time_;                                /* time at which event is ready */
    int uid_;                                    /* unique ID */
    Event() : time_(0), uid_(0) {}
};
/*
 * The base class for all event handlers. When an event's scheduled
 * time arrives, it is passed to handle which must consume it.
 * i.e., if it needs to be freed it, it must be freed by the handler.
 */
class Handler {
public:
    virtual void handle(Event* event);
};
```

Two types of objects are derived from the base class `Event`: packets and “at-events”. Packets are described in detail in the next chapter (Chapter 12.2.1). An at-event is a tcl procedure execution scheduled to occur at a particular time. This is frequently used in simulation scripts. A simple example of how it is used is as follows:

```
...
set ns_ [new Simulator]
$ns_ use-scheduler Heap
$ns_ at 300.5 "$self complete_sim"
...
```

This tcl code fragment first creates a simulation object, then changes the default scheduler implementation to be heap-based (see below), and finally schedules the function `$self complete_sim` to be executed at time 300.5 (seconds)(Note that this particular code fragment expects to be encapsulated in an object instance procedure, where the appropriate reference to `$self` is correctly defi ned.). At-events are implemented as events where the handler is effectively an execution of the tcl interpreter.

## 4.2.1 The List Scheduler

The list scheduler (`class Scheduler/List`) implements the scheduler using a simple linked-list structure. The list is kept in time-order (earliest to latest), so event insertion and deletion require scanning the list to find the appropriate entry. Choosing the next event for execution requires trimming the first entry off the head of the list. This implementation preserves event execution in a FIFO manner for simultaneous events.

### 4.2.2 the heap scheduler

The heap scheduler (`class Scheduler/Heap`) implements the scheduler using a heap structure. This structure is superior to the list structure for a large number of events, as insertion and deletion times are in  $O(\log n)$  for  $n$  events. This implementation in *ns* v2 is borrowed from the MaRS-2.0 simulator [1]; it is believed that MaRS itself borrowed the code from NetSim [14], although this lineage has not been completely verified.

### 4.2.3 The Calendar Queue Scheduler

The calendar queue scheduler (`class Scheduler/Calendar`) uses a data structure analogous to a one-year desk calendar, in which events on the same month/day of multiple years can be recorded in one day. It is formally described in [6], and informally described in Jain (p. 410) [15]. The implementation of Calendar queues in *ns* v2 was contributed by David Wetherall (presently at MIT/LCS).

### 4.2.4 The Real-Time Scheduler

The real-time scheduler (`class Scheduler/RealTime`) attempts to synchronize the execution of events with real-time. It is currently implemented as a subclass of the list scheduler. The real-time capability in *ns* is still under development, but is used to introduce an *ns* simulated network into a real-world topology to experiment with easily-configured network topologies, cross-traffic, etc. This only works for relatively slow network traffic data rates, as the simulator must be able to keep pace with the real-world packet arrival rate, and this synchronization is not presently enforced.

### 4.2.5 Precision of the scheduler clock used in ns

Precision of the scheduler clock can be defined as the smallest time-scale of the simulator that can be correctly represented. The clock variable for *ns* is represented by a double. As per the IEEE std for floating numbers, a double, consisting of 64 bits must allocate the following bits between its sign, exponent and mantissa fields.

sign	exponent	mantissa
1 bit	11 bits	52 bits

Any floating number can be represented in the form  $(X * 2^n)$  where  $X$  is the mantissa and  $n$  is the exponent. Thus the precision of timeclock in *ns* can be defined as  $(1/2^{52})$ . As simulation runs for longer times the number of remaining bits to represent the time reduces thus reducing the accuracy. Given 52 bits we can safely say time upto around  $(2^{40})$  can be represented with considerable accuracy. Anything greater than that might not be very accurate as you have remaining 12 bits to represent the time change. However  $(2^{40})$  is a very large number and we do not anticipate any problem regarding precision of time in *ns*.

## 4.3 Other Methods

The `Simulator` class provides a number of methods used to set up the simulation. They generally fall into three categories: methods to create and manage the topology (which in turn consists of managing the nodes (Chapter 5) and managing the links (Chapter 6)), methods to perform tracing (Chapter 25), and helper functions to deal with the scheduler. The following is a list of the non-topology related simulator methods:



```

Simulator instproc now                ; # return scheduler's notion of current time
Simulator instproc at args            ; # schedule execution of code at specified time
Simulator instproc cancel args        ; # cancel event
Simulator instproc run args           ; # start scheduler
Simulator instproc halt               ; # stop (pause) the scheduler
Simulator instproc flush-trace        ; # flush all trace object write buffers
Simulator instproc create-trace type files src dst ; # create trace object
Simulator instproc create_packetformat ; # set up the simulator's packet format

```

## 4.4 Commands at a glance

### Synopsis:

```
ns <otclfile> <arg> <arg>..
```

### Description:

Basic command to run a simulation script in ns.

The simulator (ns) is invoked via the ns interpreter, an extension of the vanilla otclsh command shell. A simulation is defined by a OTcl script (file). Several examples of OTcl scripts can be found under *ns/tcl/ex* directory.

The following is a list of simulator commands commonly used in simulation scripts:

```
set ns_ [new Simulator]
```

This command creates an instance of the simulator object.

```
set now [$ns_ now]
```

The scheduler keeps track of time in a simulation. This returns scheduler's notion of current time.

```
$ns_ halt
```

This stops or pauses the scheduler.

```
$ns_ run
```

This starts the scheduler.

```
$ns_ at <time> <event>
```

This schedules an <event> (which is normally a piece of code) to be executed at the specified <time>.

e.g `$ns_ at $opt(stop) "puts NS EXITING.." ; $ns_ halt"`

or, `$ns_ at 10.0 "$ftp start"`

`$ns_ cancel <event>`

Cancels the event. In effect, event is removed from scheduler's list of ready to run events.

`$ns_ create-trace <type> <file> <src> <dst> <optional arg: op>`

This creates a trace-object of type <type> between <src> and <dst> objects and attaches trace-object to <file> for writing trace-outputs. If op is defined as "nam", this creates nam tracefiles; otherwise if op is not defined, ns tracefiles are created on default.

`$ns_ flush-trace`

Flushes all trace object write buffers.

`$ns_ gen-map`

This dumps information like nodes, node components, links etc created for a given simulation. This may be broken for some scenarios (like wireless).

`$ns_ at-now <args>`

This is in effect like command "`$ns_ at $now $args`". Note that this function may not work because of tcl's string number resolution.

These are additional simulator (internal) helper functions (normally used for developing/changing the ns core code) :

`$ns_ use-scheduler <type>`

Used to specify the type of scheduler to be used for simulation. The different types of scheduler available are List, Calendar, Heap and RealTime. Currently Calendar is used as default.

`$ns_ after <delay> <event>`

Scheduling an <event> to be executed after the lapse of time <delay>.

`$ns_ clearMemTrace`

Used for memory debugging purposes.

`$ns_ is-started`

This returns true if simulator has started to run and false if not.

`$ns_ dumpq`

Command for dumping events queued in scheduler while scheduler is halted.

`$ns_ create_packetformat`

This sets up simulator's packet format.

## Chapter 5

# Nodes and Packet Forwarding

This chapter describes one aspect of creating a topology in *ns*, *i.e.*, creating the nodes. In the next chapter (Chapter 6), we will describe second aspect of creating the topology, *i.e.*, connecting the nodes to form links.

Recall that each simulation requires a single instance of the class `Simulator` to control and operate that simulation. The class provides instance procedures to create and manage the topology, and internally stores references to each element of the topology. We begin by describing the procedures in the class `Simulator` (Section 5.1). We then describe the instance procedures in the class `Node` (Section 5.2) to access and operate on individual nodes. We conclude with detailed descriptions of the `Classifier` (Section 5.4) from which the more complex node objects are formed.

The procedures and functions described in this chapter can be found in `~ns/tcl/lib/ns-lib.tcl`, `~ns/tcl/lib/ns-node.tcl`, `~ns/tcl/lib/ns-rtmodule.tcl`, `~ns/rtmodule.{cc,h}`, `~ns/classifier.{cc,h}`, `~ns/classifier-addr.cc`, `~ns/classifier-mcast.cc`, `~ns/classifier-mpath.cc`, and, `~ns/replicator.cc`.

### 5.1 Node Basics

The basic primitive for creating a node is

```
set ns [new Simulator]
$ns node
```

The instance procedure `node` constructs a node out of more simple classifier objects (Section 5.4). The Node itself is a standalone class in OTcl. However, most of the components of the node are themselves TclObjects. The typical structure of a (unicast) node is as shown in Figure 5.1. This simple structure consists of two TclObjects: an address classifier (`classifier_`) and a port classifier (`dmux_`). The function of these classifiers is to distribute incoming packets to the correct agent or outgoing link.

All nodes contain at least the following components:

- an address or `id_`, monotonically increasing by 1 (from initial value 0) across the simulation namespace as nodes are created,
- a list of neighbors (`neighbor_`),

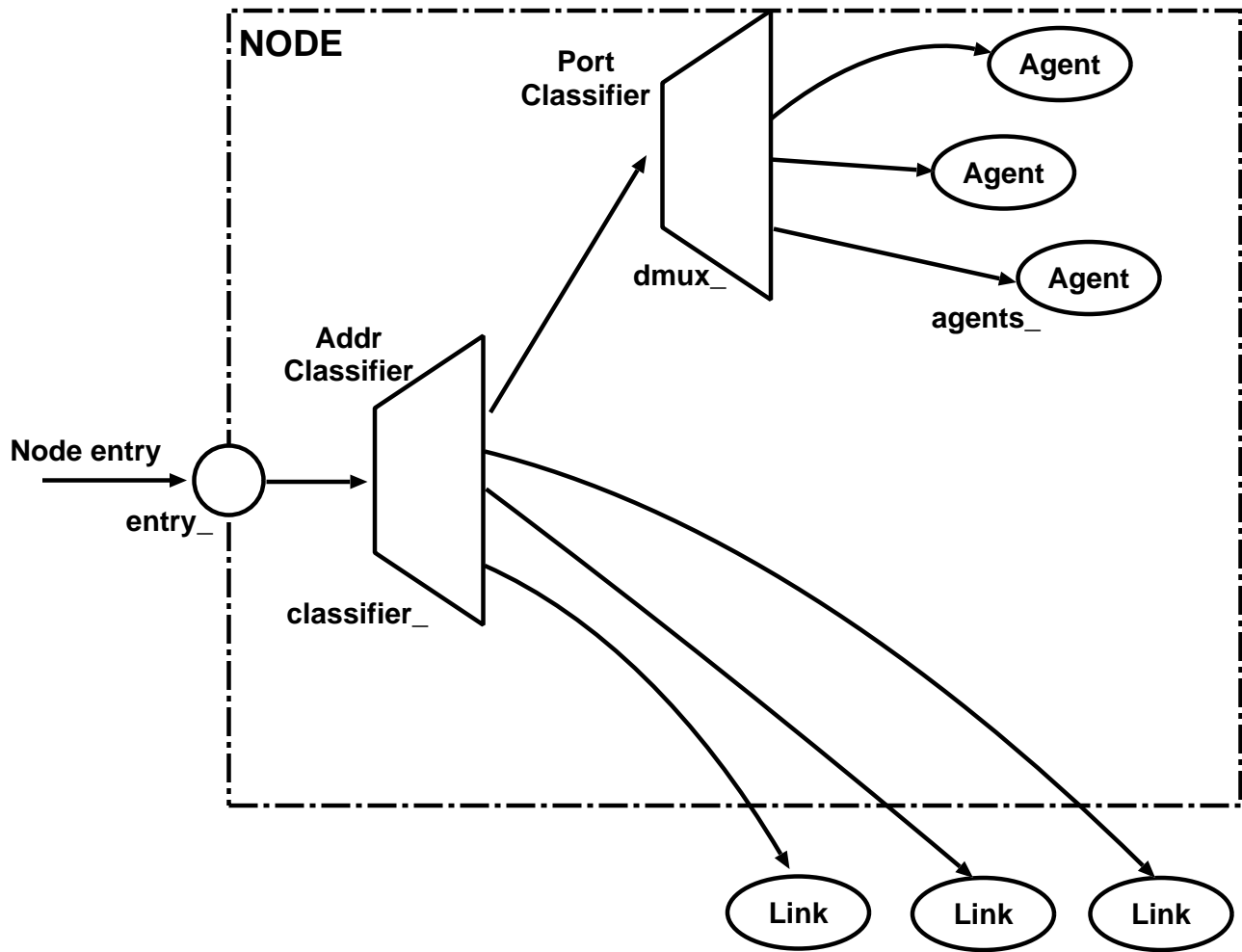


Figure 5.1: Structure of a Unicast Node. Notice that `entry_` is simply a label variable instead of a real object, e.g., the `classifier_`.

- a list of agents (`agent_`),
- a node type identifier (`nodetype_`), and
- a routing module (described in Section 5.5 below)

By default, nodes in *ns* are constructed for unicast simulations. In order to enable multicast simulation, the simulation should be created with an option “-multicast on”, e.g.:

```
set ns [new Simulator -multicast on]
```

The internal structure of a typical multicast node is shown in Figure 5.2.

When a simulation uses multicast routing, the highest bit of the address indicates whether the particular address is a multicast address or an unicast address. If the bit is 0, the address represents a unicast address, else the address represents a multicast address.

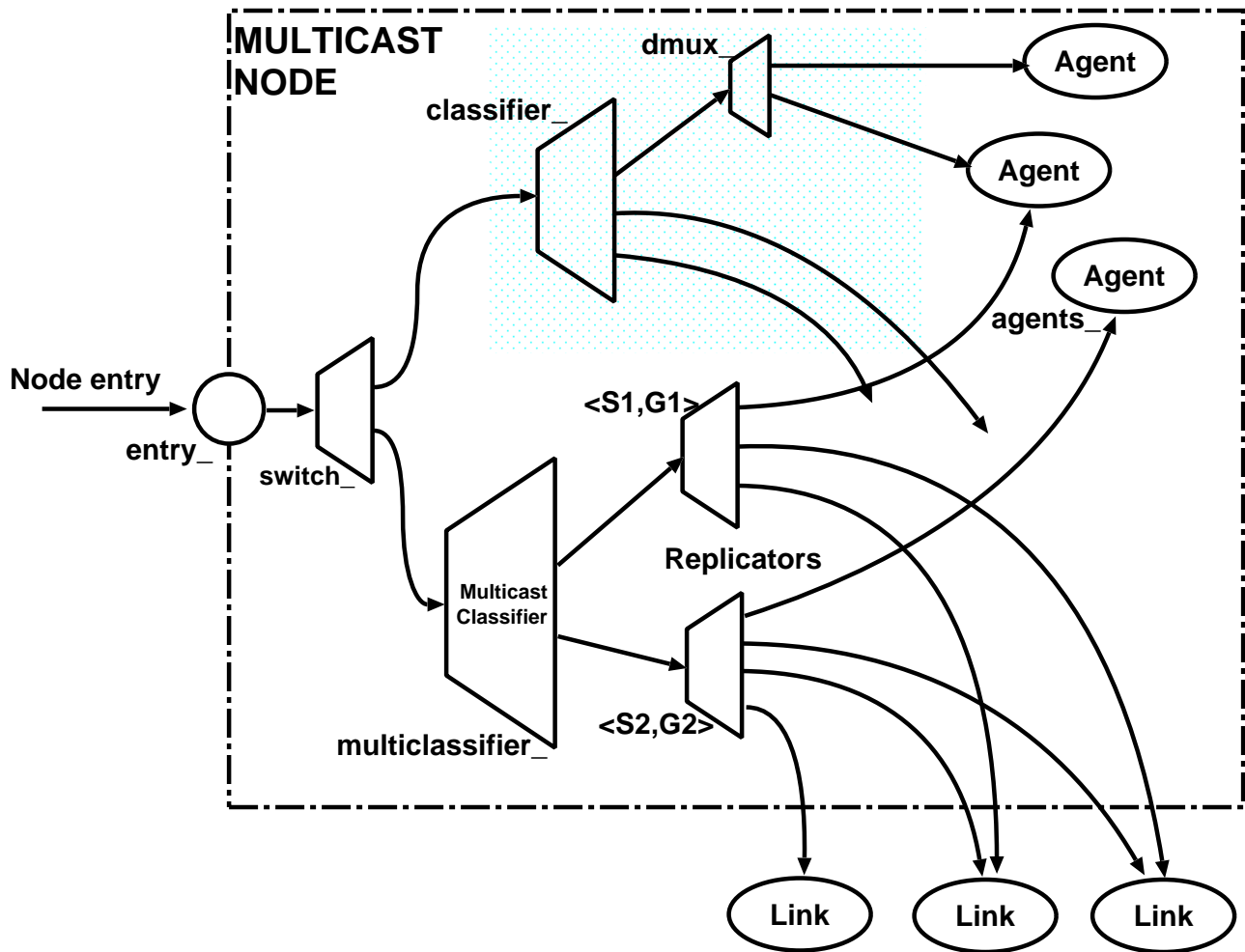


Figure 5.2: Internal Structure of a Multicast Node.

## 5.2 Node Methods: Configuring the Node

Procedures to configure an individual node can be classified into:

- Control functions
- Address and Port number management, unicast routing functions
- Agent management
- Adding neighbors

We describe each of the functions in the following paragraphs.

### Control functions

1. `$node entry` returns the entry point for a node. This is the first element which will handle packets arriving at that node.

The Node instance variable, `entry_`, stores the reference this element. For unicast nodes, this is the address classifier that looks at the higher bits of the destination address. The instance variable, `classifier_` contains the reference to this classifier. However, for multicast nodes, the entry point is the `switch_` which looks at the first bit to decide whether it should forward the packet to the unicast classifier, or the multicast classifier as appropriate.

2. `$node reset` will reset all agents at the node.

**Address and Port number management** The procedure `$node id` returns the node number of the node. This number is automatically incremented and assigned to each node at creation by the class Simulator method, `$ns node`. The class Simulator also stores an instance variable array<sup>1</sup>, `Node_`, indexed by the node id, and contains a reference to the node with that id.

The procedure `$node agent <port>` returns the handle of the agent at the specified port. If no agent at the specified port number is available, the procedure returns the null string.

The procedure `alloc-port` returns the next available port number. It uses an instance variable, `np_`, to track the next unallocated port number.

The procedures, `add-route` and `add-routes`, are used by unicast routing (Chapter 28) to add routes to populate the `classifier_`. The usage syntax is `$node add-route <destination id> <TclObject>`. `TclObject` is the entry of `dmux_`, the port demultiplexer at the node, if the destination id is the same as this node's id, it is often the head of a link to send packets for that destination to, but could also be the the entry for other classifiers or types of classifiers.

`$node add-routes <destination id> <TclObjects>` is used to add multiple routes to the same destination that must be used simultaneously in round robin manner to spread the bandwidth used to reach that destination across all links equally. It is used only if the instance variable `multiPath_` is set to 1, and detailed dynamic routing strategies are in effect, and requires the use of a `multiPath` classifier. We describe the implementation of the `multiPath` classifier later in this chapter (Section 5.4); however, we defer the discussion of multipath routing (Chapter 28) to the chapter on unicast routing.

The dual of `add-routes{}` is `delete-routes{}`. It takes the id, a list of `TclObjects`, and a reference to the simulator's `nullagent`. It removes the `TclObjects` in the list from the installed routes in the `multiPath` classifier. If the route entry in the classifier does not point to a `multiPath` classifier, the routine simply clears the entry from `classifier_`, and installs the `nullagent` in its place.

Detailed dynamic routing also uses two additional methods: the instance procedure `init-routing{}` sets the instance variable `multiPath_` to be equal to the class variable of the same name. It also adds a reference to the route controller object at that node in the instance variable, `rtObject_`. The procedure `rtObject?{}` returns the handle for the route object at the node.

Finally, the procedure `intf-changed{}` is invoked by the network dynamics code if a link incident on the node changes state. Additional details on how this procedure is used are discussed later in the chapter on network dynamics (Chapter 30).

**Agent management** Given an `<agent>`, the procedure `attach{}` will add the agent to its list of `agents_`, assign a port number the agent and set its source address, set the target of the agent to be its (*i.e.*, the node's) `entry{}`, and add a pointer to the port demultiplexer at the node (`dmux_`) to the agent at the corresponding slot in the `dmux_` classifier.

Conversely, `detach{}` will remove the agent from `agents_`, and point the agent's target, and the entry in the node `dmux_` to `nullagent`.

---

<sup>1</sup>*i.e.*, an instance variable of a class that is also an array variable



**Tracking Neighbors** Each node keeps a list of its adjacent neighbors in its instance variable, `neighbor_`. The procedure `add-neighbor{}` adds a neighbor to the list. The procedure `neighbors{}` returns this list.

## 5.3 Node Configuration Interface

**NOTE:** This API, especially its internal implementation which is messy at this point, is still a moving target. It may undergo significant changes in the near future. However, we will do our best to maintain the same interface as described in this chapter. In addition, this API currently does not cover all existing nodes in the old format, namely, nodes built using inheritance, and parts of mobile IP. It is principally oriented towards wireless and satellite simulation. [Sep 15, 2000; updated June 2001].

`Simulator::node-config{}` accommodates flexible and modular construction of different node definitions within the same base Node class. For instance, to create a mobile node capable of wireless communication, one no longer needs a specialized node creation command, e.g., `dsdv-create-mobile-node{}`; instead, one changes default configuration parameters, such as

```
$ns node-config -adhocRouting dsdv
```

before actually creating the node with the command: `$ns node`. Together with routing modules, this allows one to combine “arbitrary” routing functionalities within a single node without resorting to multiple inheritance and other fancy object gimmicks. We will describe this in more detail in Section 5.5. The functions and procedures relevant to the new node APIs may be found in `~ns/tcl/lib/ns-node.tcl`.

The node configuration interface consists of two parts. The first part deals with node configuration, while the second part actually creates nodes of the specified type. We have already seen the latter in Section 5.1, in this section we will describe the configuration part.

Node configuration essentially consists of defining the different node characteristics before creating them. They may consist of the type of addressing structure used in the simulation, defining the network components for mobilenodes, turning on or off the trace options at Agent/Router/MAC levels, selecting the type of adhoc routing protocol for wireless nodes or defining their energy model.

As an example, node-configuration for a wireless, mobile node that runs AODV as its adhoc routing protocol in a hierarchical topology would be as shown below. We decide to turn tracing on at the agent and router level only. Also we assume a topology has been instantiated with `"set topo [new Topography]"`. The node-config command would look like the following:

```
$ns_ node-config -addressType hierarchical \  
                -adhocRouting AODV \  
                -llType LL \  
                -macType Mac/802_11 \  
                -ifqType Queue/DropTail/PriQueue \  
                -ifqLen 50 \  
                -antType Antenna/OmniAntenna \  
                -propType Propagation/TwoRayGround \  
                -phyType Phy/WirelessPhy \  
                -topologyInstance $topo \  
                -channel Channel/WirelessChannel \  
                -agentTrace ON \  
                -routerTrace ON \  
                -macTrace OFF \  
                -movementTrace OFF
```

The default values for all the above options are NULL except `-addressingType` whose default value is flat. The option `-reset` can be used to reset all node-configuration parameters to their default value.

Note that the `config` command can be broken down into separate lines like

```
$ns_ node-config -addressingType hier
$ns_ node-config -macTrace ON
```

The options that need to be changed may only be called. For example after configuring for AODV mobilenodes as shown above (and after creating AODV mobilenodes), we may configure for AODV base-station nodes in the following way:

```
$ns_ node-config -wiredRouting ON
```

While all other features for base-station nodes and mobilenodes are same, the base-station nodes are capable of wired routing, while mobilenodes are not. In this way we can change node-configuration only when it is required.

All node instances created after a given node-configuration command will have the same property unless a part or all of the node-configuration command is executed with different parameter values. And all parameter values remain unchanged unless they are explicitly changed. So after creation of the AODV base-station and mobilenodes, if we want to create simple nodes, we will use the following node-configuration command:

```
$ns_ node-config -reset
```

This will set all parameter values to their default setting which basically defines configuration of a simple node.

Currently, this type of node configuration is oriented towards wireless and satellite nodes. Table 5.1 lists the available options for these kinds of nodes. The example scripts `~ns/tcl/ex/simple-wireless.tcl` and `~ns/tcl/ex/sat-mixed.tcl` provide usage examples.

## 5.4 The Classifier

The function of a node when it receives a packet is to examine the packet's fields, usually its destination address, and on occasion, its source address. It should then map the values to an outgoing interface object that is the next downstream recipient of this packet.

In *ns*, this task is performed by a simple *classifier* object. Multiple classifier objects, each looking at a specific portion of the packet forward the packet through the node. A node in *ns* uses many different types of classifiers for different purposes. This section describes some of the more common, or simpler, classifier objects in *ns*.

We begin with a description of the base class in this section. The next subsections describe the address classifier (Section 5.4.1), the multicast classifier (Section 5.4.2), the multipath classifier (Section 5.4.3), the hash classifier (Section 5.4.4), and finally, the replicator (Section 5.4.5).

A classifier provides a way to match a packet against some logical criteria and retrieve a reference to another simulation object based on the match results. Each classifier contains a table of simulation objects indexed by *slot number*. The job of a classifier is to determine the slot number associated with a received packet and forward that packet to the object referenced by that particular slot. The C++ class `Classifier` (defined in `~ns/classifier.h`) provides a base class from which other classifiers are derived.

option	available values	default
<b>general</b>		
addressType	flat, hierarchical	flat
MPLS	ON, OFF	OFF
<b>both satellite- and wireless-oriented</b>		
wiredRouting	ON, OFF	OFF
llType	LL, LL/Sat	""
macType	Mac/802_11, Mac/Csma/Ca, Mac/Sat, Mac/Sat/UnslottedAloha, Mac/Tdma	""
ifqType	Queue/DropTail, Queue/DropTail/PriQueue	""
phyType	Phy/WirelessPhy, Phy/Sat	""
<b>wireless-oriented</b>		
adhocRouting	DIFFUSION/RATE, DIFFUSION/PROB, DSDV, DSR, FLOODING, OMNIMCAST, AODV, TORA	""
propType	Propagation/TwoRayGround, Propagation/Shadowing	""
propInstance	Propagation/TwoRayGround, Propagation/Shadowing	""
antType	Antenna/OmniAntenna	""
channel	Channel/WirelessChannel, Channel/Sat	""
topoInstance	<topology file>	""
mobileIP	ON, OFF	OFF
energyModel	EnergyModel	""
initialEnergy	<value in Joules>	""
rxPower	<value in W>	""
txPower	<value in W>	""
idlePower	<value in W>	""
agentTrace	ON, OFF	OFF
routerTrace	ON, OFF	OFF
macTrace	ON, OFF	OFF
movementTrace	ON, OFF	OFF
errProc	UniformErrorProc	""
FECProc	?	?
toraDebug	ON, OFF	OFF
<b>satellite-oriented</b>		
satNodeType	polar, geo, terminal, geo-repeater	""
downlinkBW	<bandwidth value, e.g. "2Mb">	""

Table 5.1: Available options for node configuration (see tcl/lib/ns-lib.tcl).

```

class Classifier : public NsObject {
public:
    ~Classifier();
    void recv(Packet*, Handler* h = 0);
protected:
    Classifier();
    void install(int slot, NsObject*);
    void clear(int slot);
    virtual int command(int argc, const char*const* argv);
    virtual int classify(Packet *const) = 0;
    void alloc(int);
    NsObject** slot_; /* table that maps slot number to a NsObject */
    int nslot_;
    int maxslot_;
};

```

The `classify()` method is pure virtual, indicating the class `Classifier` is to be used only as a base class. The `alloc()` method dynamically allocates enough space in the table to hold the specified number of slots. The `install()` and `clear()` methods add or remove objects from the table. The `recv()` method and the OTcl interface are implemented as follows in `~ns/classifier.cc`:

```

/*
 * objects only ever see "packet" events, which come either
 * from an incoming link or a local agent (i.e., packet source).
 */
void Classifier::recv(Packet* p, Handler*)
{
    NsObject* node;
    int cl = classify(p);
    if (cl < 0 || cl >= nslot_ || (node = slot_[cl]) == 0) {
        Tcl::instance().evalf("%s no-slot %d", name(), cl);
        Packet::free(p);
        return;
    }
    node->recv(p);
}

int Classifier::command(int argc, const char*const* argv)
{
    Tcl& tcl = Tcl::instance();
    if (argc == 3) {
        /*
         * $classifier clear $slot
         */
        if (strcmp(argv[1], "clear") == 0) {
            int slot = atoi(argv[2]);
            clear(slot);
            return (TCL_OK);
        }
        /*
         * $classifier installNext $node
         */
        if (strcmp(argv[1], "installNext") == 0) {
            int slot = maxslot_ + 1;
            NsObject* node = (NsObject*)TclObject::lookup(argv[2]);
            install(slot, node);
            tcl.resultf("%u", slot);
            return TCL_OK;
        }
        if (strcmp(argv[1], "slot") == 0) {
            int slot = atoi(argv[2]);
            if ((slot >= 0) || (slot < nslot_)) {
                tcl.resultf("%s", slot_[slot]->name());
                return TCL_OK;
            }
            tcl.resultf("Classifier: no object at slot %d", slot);
            return (TCL_ERROR);
        }
    }
    else if (argc == 4) {
        /*

```

```

        * $classifier install $slot $node
        */
    if (strcmp(argv[1], "install") == 0) {
        int slot = atoi(argv[2]);
        NsObject* node = (NsObject*)TclObject::lookup(argv[3]);
        install(slot, node);
        return (TCL_OK);
    }
}
return (NsObject::command(argc, argv));
}

```

When a classifier `recv()`'s a packet, it hands it to the `classify()` method. This is defined differently in each type of classifier derived from the base class. The usual format is for the `classify()` method to determine and return a slot index into the table of slots. If the index is valid, and points to a valid `TclObject`, the classifier will hand the packet to that object using that object's `recv()` method. If the index is not valid, the classifier will invoke the instance procedure `no-slot{}` to attempt to populate the table correctly. However, in the base class `Classifier::no-slot{}` prints an error message and terminates execution.

The `command()` method provides the following instproc-like to the interpreter:

- `clear{<slot>}` clears the entry in a particular slot.
- `installNext{<object>}` installs the object in the next available slot, and returns the slot number.  
Note that this instproc-like is overloaded by an instance procedure of the same name that stores a reference to the object stored. This then helps quick query of the objects installed in the classifier from OTcl.
- `slot{<index>}` returns the object stored in the specified slot.
- `install{<index>, <object>}` installs the specified <object> at the slot <index>.  
Note that this instproc-like too is overloaded by an instance procedure of the same name that stores a reference to the object stored. This is also to quickly query of the objects installed in the classifier from OTcl.

### 5.4.1 Address Classifiers

An address classifier is used in supporting unicast packet forwarding. It applies a bitwise shift and mask operation to a packet's destination address to produce a slot number. The slot number is returned from the `classify()` method. The class `AddressClassifier` (defined in `~ns/classifier-addr.cc`) is defined as follows:

```

class AddressClassifier : public Classifier {
public:
    AddressClassifier() : mask_(~0), shift_(0) {
        bind("mask_", (int*)&mask_);
        bind("shift_", &shift_);
    }
protected:
    int classify(Packet *const p) {
        IPHeader *h = IPHeader::access(p->bits());
        return ((h->dst() >> shift_) & mask_);
    }
    nsaddr_t mask_;
    int shift_;
}

```

```
};
```

The class imposes no direct semantic meaning on a packet's destination address field. Rather, it returns some number of bits from the packet's `dst_` field as the slot number used in the `Classifier::recv()` method. The `mask_` and `shift_` values are set through OTcl.

## 5.4.2 Multicast Classifiers

The multicast classifier classifies packets according to both source and destination (group) addresses. It maintains a (chained hash) table mapping source/group pairs to slot numbers. When a packet arrives containing a source/group unknown to the classifier, it invokes an OTcl procedure `Node::new-group{}` to add an entry to its table. This OTcl procedure may use the method `set-hash` to add new (source, group, slot) 3-tuples to the classifier's table. The multicast classifier is defined in `~ns/classifier-mcast.cc` as follows:

```
static class MCastClassifierClass : public TclClass {
public:
    MCastClassifierClass() : TclClass("Classifier/Multicast") {}
    TclObject* create(int argc, const char*const* argv) {
        return (new MCastClassifier());
    }
} class_mcast_classifier;

class MCastClassifier : public Classifier {
public:
    MCastClassifier();
    ~MCastClassifier();
protected:
    int command(int argc, const char*const* argv);
    int classify(Packet *const p);
    int findslot();
    void set_hash(nsaddr_t src, nsaddr_t dst, int slot);
    int hash(nsaddr_t src, nsaddr_t dst) const {
        u_int32_t s = src ^ dst;
        s ^= s >> 16;
        s ^= s >> 8;
        return (s & 0xff);
    }
    struct hashnode {
        int slot;
        nsaddr_t src;
        nsaddr_t dst;
        hashnode* next;
    };
    hashnode* ht_[256];
    const hashnode* lookup(nsaddr_t src, nsaddr_t dst) const;
};

int MCastClassifier::classify(Packet *const pkt)
{
    IPHeader *h = IPHeader::access(pkt->bits());
    nsaddr_t src = h->src() >> 8; /*XXX*/
```

```

nsaddr_t dst = h->dst();
const hashnode* p = lookup(src, dst);
if (p == 0) {
    /*
     * Didn't find an entry.
     * Call tcl exactly once to install one.
     * If tcl doesn't come through then fail.
     */
    Tcl::instance().evalf("%s new-group %u %u", name(), src, dst);
    p = lookup(src, dst);
    if (p == 0)
        return (-1);
}
return (p->slot);
}

```

The class `MCastClassifier` implements a chained hash table and applies a hash function on both the packet source and destination addresses. The hash function returns the slot number to index the `slot_` table in the underlying object. A hash miss implies packet delivery to a previously-unknown group; OTcl is called to handle the situation. The OTcl code is expected to insert an appropriate entry into the hash table.

### 5.4.3 MultiPath Classifier

This object is devised to support equal cost multipath forwarding, where the node has multiple equal cost routes to the same destination, and would like to use all of them simultaneously. This object does not look at any field in the packet. With every succeeding packet, it simply returns the next filled slot in round robin fashion. The definitions for this classifier are in `~ns/classifier-mpath.cc`, and are shown below:

```

class MultiPathForwarder : public Classifier {
public:
    MultiPathForwarder() : ns_(0), Classifier() {}
    virtual int classify(Packet* const) {
        int cl;
        int fail = ns_;
        do {
            cl = ns_++;
            ns_ %= (maxslot_ + 1);
        } while (slot_[cl] == 0 && ns_ != fail);
        return cl;
    }
private:
    int ns_;
};
/* next slot to be used. Probably a misnomer? */

```

### 5.4.4 Hash Classifier

This object is used to classify a packet as a member of a particular *flow*. As their name indicates, hash classifiers use a hash table internally to assign packets to flows. These objects are used where flow-level information is required (e.g. in

flow-specific queuing disciplines and statistics collection). Several “flow granularities” are available. In particular, packets may be assigned to flows based on flow ID, destination address, source/destination addresses, or the combination of source/destination addresses plus flow ID. The fields accessed by the hash classifier are limited to the ip header: `src()`, `dst()`, `flowid()` (see `ip.h`).

The hash classifier is created with an integer argument specifying the initial size of its hash table. The current hash table size may be subsequently altered with the `resize` method (see below). When created, the instance variables `shift_` and `mask_` are initialized with the simulator’s current `NodeShift` and `NodeMask` values, respectively. These values are retrieved from the `AddrParams` object when the hash classifier is instantiated. The hash classifier will fail to operate properly if the `AddrParams` structure is not initialized. The following constructors are used for the various hash classifiers:

```
Classifier/Hash/SrcDest
Classifier/Hash/Dst
Classifier/Hash/Fid
Classifier/Hash/SrcDestFid
```

The hash classifier receives packets, classifies them according to their flow criteria, and retrieves the classifier *slot* indicating the next node that should receive the packet. In several circumstances with hash classifiers, most packets should be associated with a single slot, while only a few flows should be directed elsewhere. The hash classifier includes a `default_` instance variable indicating which slot is to be used for packets that do not match any of the per-flow criteria. The `default_` may be set optionally.

The methods for a hash classifier are as follows:

```
$hashcl set-hash buck src dst fid slot
$hashcl lookup buck src dst fid
$hashcl del-hash src dst fid
$hashcl resize nbuck
```

The `set-hash()` method inserts a new entry into the hash table within the hash classifier. The `buck` argument specifies the hash table bucket number to use for the insertion of this entry. When the bucket number is not known, `buck` may be specified as `auto`. The `src`, `dst` and `fid` arguments specify the IP source, destination, and flow IDs to be matched for flow classification. Fields not used by a particular classifier (e.g. specifying `src` for a flow-id classifier) is ignored. The `slot` argument indicates the index into the underlying slot table in the base `Classifier` object from which the hash classifier is derived. The `lookup` function returns the name of the object associated with the given `buck/src/dst/fid` tuple. The `buck` argument may be `auto`, as for `set-hash`. The `del-hash` function removes the specified entry from the hash table. Currently, this is done by simply marking the entry as inactive, so it is possible to populate the hash table with unused entries. The `resize` function resizes the hash table to include the number of buckets specified by the argument `nbuck`.

Provided no default is defined, a hash classifier will perform a call into `OTcl` when it receives a packet which matches no flow criteria. The call takes the following form:

```
$obj unknown-flow src dst flowid buck
```

Thus, when a packet matching no flow criteria is received, the method `unknown-flow` of the instantiated hash classifier object is invoked with the source, destination, and flow id fields from the packet. In addition, the `buck` field indicates the hash bucket which should contain this flow if it were inserted using `set-hash`. This arrangement avoids another hash lookup when performing insertions into the classifier when the bucket is already known.



### 5.4.5 Replicator

The replicator is different from the other classifiers we have described earlier, in that it does not use the `classify` function. Rather, it simply uses the classifier as a table of  $n$  slots; it overloads the `recv()` method to produce  $n$  copies of a packet, that are delivered to all  $n$  objects referenced in the table.

To support multicast packet forwarding, a classifier receiving a multicast packet from source  $S$  destined for group  $G$  computes a hash function  $h(S, G)$  giving a “slot number” in the classifier’s object table. In multicast delivery, the packet must be copied once for each link leading to nodes subscribed to  $G$  minus one. Production of additional copies of the packet is performed by a `Replicator` class, defined in `replicator.cc`:

```
/*
 * A replicator is not really a packet classifier but
 * we simply find convenience in leveraging its slot table.
 * (this object used to implement fan-out on a multicast
 * router as well as broadcast LANs)
 */
class Replicator : public Classifier {
public:
    Replicator();
    void recv(Packet*, Handler* h = 0);
    virtual int classify(Packet* const) {};
protected:
    int ignore_;
};

void Replicator::recv(Packet* p, Handler*)
{
    IPHeader *iph = IPHeader::access(p->bits());
    if (maxslot_ < 0) {
        if (!ignore_)
            Tcl::instance().evalf("%s drop %u %u", name(),
                                  iph->src(), iph->dst());
        Packet::free(p);
        return;
    }
    for (int i = 0; i < maxslot_; ++i) {
        NsObject* o = slot_[i];
        if (o != 0)
            o->recv(p->copy());
    }
    /* we know that maxslot is non-null */
    slot_[maxslot_]->recv(p);
}
```

As we can see from the code, this class does not really classify packets. Rather, it replicates a packet, one for each entry in its table, and delivers the copies to each of the nodes listed in the table. The last entry in the table gets the “original” packet. Since the `classify()` method is pure virtual in the base class, the replicator defines an empty `classify()` method.

## 5.5 Routing Module and Classifier Organization

As we have seen, a *ns* node is essentially a collection of classifiers. The simplest node (unicast) contains only one address classifier and one port classifier, as shown in Figure 5.1. When one extends the functionality of the node, more classifiers are added into the base node, for instance, the multicast node shown in Figure 5.2. As more function blocks are added, and each of these blocks requires its own classifier(s), it becomes important for the node to provide a *uniform* interface to organize these classifiers and to bridge these classifiers to the route computation blocks.

The classical method to handle this case is through class inheritance. For instance, if one wants a node that supports hierarchical routing, one simply derive a *Node/Hier* from the base node and override the classifier setup methods to insert hierarchical classifiers. This method works well when the new function blocks are independent and cannot be “arbitrarily” mixed. For instance, both hierarchical routing and ad hoc routing use their own set of classifiers. Inheritance would require that we have *Node/Hier* that supports the former, and *Node/Mobile* for the latter. This becomes slightly problematic when one wants an ad hoc routing node that supports hierarchical routing. In this simple case one may use multiple inheritance to solve the problem, but this quickly becomes infeasible as the number of such function blocks increases.

The only method to solve this problem is object composition. The base node needs to define a set of interfaces for classifier access and organization. These interfaces should

- allow individual routing modules that implement their own classifiers to insert their classifiers into the node;
- allow route computation blocks to populate routes to classifiers in all routing modules that need this information,
- provide a single point of management for existing routing modules.

In addition, we should also define a uniform interface for routing modules to connect to the node interfaces, so as to provide a systematic approach to extending node functionality. In this section we will describe the design of routing modules as well as that of the corresponding node interfaces.

### 5.5.1 Routing Module

In general, every routing implementation in *ns* consists of three function blocks:

- *Routing agent* exchanges routing packet with neighbors,
- *Route logic* uses the information gathered by routing agents (or the global topology database in the case of static routing) to perform the actual route computation,
- *Classifiers* sit inside a Node. They use the computed routing table to perform packet forwarding.

Notice that when implementing a new routing protocol, one does not necessarily implement all of these three blocks. For instance, when one implements a link state routing protocol, one simply implement a routing agent that exchanges information in the link state manner, and a route logic that does Dijkstra on the resulting topology database. It can then use the same classifiers as other unicast routing protocols.

When a new routing protocol implementation includes more than one function blocks, especially when it contains its own classifier, it is desirable to have another object, which we call a *routing module*, that manages all these function blocks and to interface with node to organize its classifiers. Figure 5.3 shows functional relation among these objects. Notice that routing modules may have direct relationship with route computation blocks, i.e., route logic and/or routing agents. However, route computation MAY not install their routes directly through a routing module, because there may exist other modules that

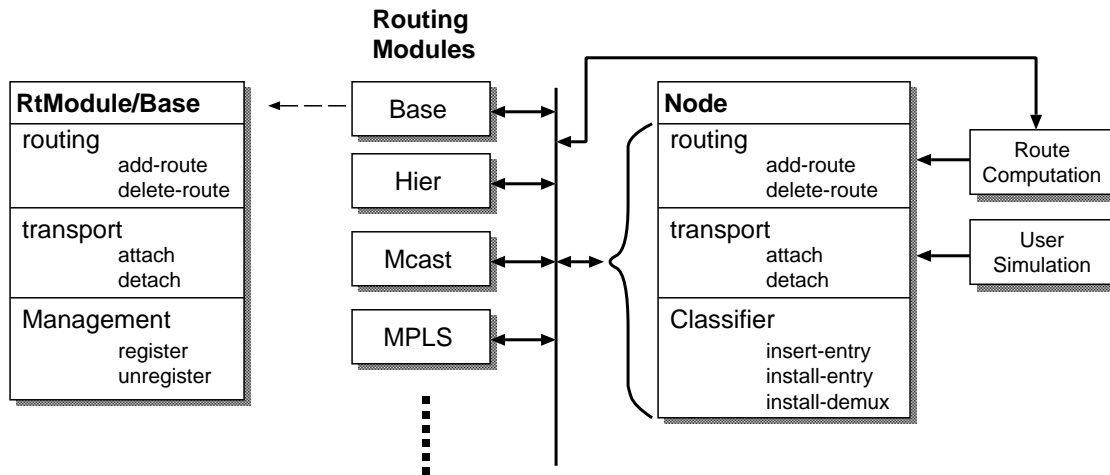


Figure 5.3: Interaction among node, routing module, and routing. The dashed line shows the details of one routing module.

are interested in learning about the new routes. This is not a requirement, however, because it is possible that some route computation is specific to one particular routing module, for instance, label installation in the MPLS module.

A routing module contains three major functionalities:

1. A routing module initializes its connection to a node through `register{}`, and tears the connection down via `unregister{}`. Usually, in `register{}` a routing module (1) tells the node whether it interests in knowing route updates and transport agent attachments, and (2) creates its classifiers and install them in the node (details described in the next subsection). In `unregister{}` a routing module does the exact opposite: it deletes its classifiers and removes its hooks on routing update in the node.
2. If a routing module is interested in knowing routing updates, the node will inform the module via `RtModule::add-route{dst, target}` and `RtModule::delete-route{dst, nullagent}`.
3. If a routing module is interested in learning about transport agent attachment and detachment in a node, the node will inform the module via `RtModule::attach{agent, port}` and `RtModule::detach{agent, nullagent}`.

There are two steps to write your own routing module:

1. You need to declare the C++ part of your routing module (see `~ns/rmodule.{cc,h}`). For many modules this only means to declare a virtual method `name()` which returns a string descriptor of the module. However, you are free to implement as much functionality as you like in C++; if necessary you may later move functionality from OTcl into C++ for better performance.
2. You need to look at the above interfaces implemented in the base routing module (see `~ns/tcl/lib/ns-rmodule.tcl`) and decide which one you'll inherit, which one you'll override, and put them in OTcl interfaces of your own module.

There are several derived routing module examples in `~ns/tcl/lib/ns-rmodule.tcl`, which may serve as templates for your modules.

Currently, there are six routing modules implemented in `ns`:

Module Name	Functionality
RtModule/Base	Interface to unicast routing protocols. Provide basic functionality to add/delete route and attach/detach agents.
RtModule/Mcast	Interface to multicast routing protocols. Its only purpose is establishes multicast classifiers. All other multicast functionalities are implemented as instprocs of Node. This should be converted in the future.
RtModule/Hier	Hierarchical routing. It's a wrapper for managing hierarchical classifiers and route installation. Can be combined with other routing protocols, e.g., ad hoc routing.
RtModule/Manual	Manual routing.
RtModule/VC	Uses virtual classifier instead of vanilla classifier.
RtModule/MPLS	Implements MPLS functionality. This is the only existing module that is completely self-contained and does not pollute the Node namespace.

Table 5.2: Available routing modules

## 5.5.2 Node Interface

To connect to the above interfaces of routing module, a node provides a similar set of interfaces:

- In order to know which module to register during creation, the Node class keeps a list of modules as a class variable. The default value of this list contains only the base routing module. The Node class provides the following two *procs* to manipulate this module list:

`Node::enable-module{[name]}` If module `RtModule/[name]` exists, this proc puts `[name]` into the module list.

`Node::disable-module{[name]}` If `[name]` is in the module list, remove it from the list.

When a node is created, it goes through the module list of the Node class, creates all modules included in the list, and register these modules at the node.

After a node is created, one may use the following instprocs to list modules registered at the node, or to get a handle of a module with a particular name:

`Node::list-modules{}` Return a list of the handles (shadow objects) of all registered modules.

`Node::get-module{[name]}` Return a handle of the registered module whose name matches the given one. Notice that any routing module can only have a single instance registered at any node.

- To allow routing modules register their interests of routing updates, a node object provide the following instprocs:

`Node::route-notify{module}` Add module into route update notification list.

`Node::unreg-route-notify{module}` Remove module from route update notification list.

Similarly, the following instprocs provide hooks on the attachment of transport agents:

`Node::port-notify{module}` Add module into agent attachment notification list.

`Node::unreg-port-notify{module}` Remove module from agent attachment notification list.

Notice that in all of these instprocs, parameter `module` should be a module handle instead of a module name.

- Node provides the following instprocs to manipulate its address and port classifiers:
  - `Node::insert-entry{module, clsfr, hook}` inserts classifier `clsfr` into the entry point of the node. It also associates the new classifier with `module` so that if this classifier is removed later, `module` will be unregistered. If `hook` is specified as a number, the existing classifier will be inserted into slot `hook` of the new classifier. In this way, one may establish a “chain” of classifiers; see Figure 5.2 for an example. **NOTE:** `clsfr` needs NOT

to be a classifier. In some cases one may want to put an agent, or any class derived from Connector, at the entry point of a node. In such cases, one simply supplies `target` to parameter `hook`.

- `Node::install-entry{module, clsfr, hook}` differs from `Node::insert-entry` in that it deletes the existing classifier at the node entry point, unregisters any associated routing module, and installs the new classifier at that point. If `hook` is given, and the old classifier is connected into a classifier chain, it will connect the chain into slot `hook` of the new classifier. As above, if `hook` equals to `target`, `clsfr` will be treated as an object derived from Connector instead of a classifier.
- `Node::install-demux{demux, port}` places the given classifier `demux` as the default demultiplexer. If `port` is given, it plugs the existing demultiplexer into slot `port` of the new one. Notice that in either case it does not delete the existing demultiplexer.

## 5.6 Commands at a glance

Following is a list of common node commands used in simulation scripts:

```
$ns_ node [<hier_addr>]
```

Command to create and return a node instance. If `<hier_addr>` is given, assign the node address to be `<hier_addr>`. Note that the latter MUST only be used when hierarchical addressing is enabled via either `set-address-format hierarchical{}` or `node-config -addressType hierarchical{}`.

```
$ns_ node-config -<config-parameter> <optional-val>
```

This command is used to configure nodes. The different config-parameters are `addressingType`, different type of the network stack components, whether tracing will be turned on or not, `mobileIP` flag is turned on or not, energy model is being used or not etc. An option `-reset` maybe used to set the node configuration to its default state. The default setting of node-config, i.e if no values are specified, creates a simple node (base class Node) with flat addressing/routing. For the syntax details see Section 5.3.

```
$node id
```

Returns the id number of the node.

```
$node node-addr
```

Returns the address of the node. In case of flat addressing, the node address is same as its node-id. In case of hierarchical addressing, the node address in the form of a string (viz. "1.4.3") is returned.

```
$node reset
```

Resets all agent attached to this node.

```
$node agent <port_num>
```

Returns the handle of the agent at the specified port. If no agent is found at the given port, a null string is returned.

```
$node entry
```

Returns the entry point for the node. This is first object that handles packet receiving at this node.

```
$node attach <agent> <optional:port_num>
```

Attaches the `<agent>` to this node. In case no specific port number is passed, the node allocates a port number and binds the agent to this port. Thus once the agent is attached, it receives packets destined for this host (node) and port.

```
$node detach <agent> <null_agent>
```

This is the dual of "attach" described above. It detaches the agent from this node and installs a null-agent to the port this agent was attached. This is done to handle transit packets that may be destined to the detached agent. These on-the-fly packets are then sinked at the null-agent.

`$node neighbors`

This returns the list of neighbors for the node.

`$node add-neighbor <neighbor_node>`

This is a command to add `<neighbor_node>` to the list of neighbors maintained by the node.

---

Following is a list of internal node methods:

`$node add-route <destination_id> <target>`

This is used in unicast routing to populate the classifier. The target is a Tcl object, which may be the entry of `dmux_` (port demultiplexer in the node) in case the `<destination_id>` is same as this node-id. Otherwise it is usually the head of the link for that destination. It could also be the entry for other classifiers.

`$node alloc-port <null_agent>`

This returns the next available port number.

`$node incr-rtgtable-size`

The instance variable `rtsize_` is used to keep track of size of routing-table in each node. This command is used to increase the routing-table size every time an routing-entry is added to the classifiers.

There are other node commands that supports hierarchical routing, detailed dynamic routing, equal cost multipath routing, manual routing, and energy model for mobile nodes. These and other methods described earlier can be found in `~ns/tcl/lib/ns-node.tcl` and `~ns/tcl/lib/ns-mobilenode.tcl`.

## Chapter 6

# Links: Simple Links

This is the second aspect of defining the topology. In the previous chapter (Chapter 5), we had described how to create the nodes in the topology in *ns*. We now describe how to create the links to connect the nodes and complete the topology. In this chapter, we restrict ourselves to describing the simple point to point links. *ns* supports a variety of other media, including an emulation of a multi-access LAN using a mesh of simple links, and other true simulation of wireless and broadcast media. They will be described in a separate chapter. The CBQlink is derived from simple links and is a considerably more complex form of link that is also not described in this chapter.

We begin by describing the commands to create a link in this section. As with the node being composed of classifiers, a simple link is built up from a sequence of connectors. We also briefly describe some of the connectors in a simple link. We then describe the instance procedures that operate on the various components of defined by some of these connectors (Section 6.1). We conclude the chapter with a description the connector object (Section 6.2), including brief descriptions of the common link connectors.

The class `Link` is a standalone class in OTcl, that provides a few simple primitives. The class `SimpleLink` provides the ability to connect two nodes with a point to point link. *ns* provides the instance procedure `simplex-link{}` to form a unidirectional link from one node to another. The link is in the class `SimpleLink`. The following describes the syntax of the simplex link:

```
set ns [new Simulator]
$ns simplex-link <node0> <node1> <bandwidth> <delay> <queue_type>
```

The command creates a link from `<node0>` to `<node1>`, with specified `<bandwidth>` and `<delay>` characteristics. The link uses a queue of type `<queue_type>`. The procedure also adds a TTL checker to the link. Five instance variables define the link:

<code>head_</code>	Entry point to the link, it points to the first object in the link.
<code>queue_</code>	Reference to the main queue element of the link. Simple links usually have one queue per link. Other more complex types of links may have multiple queue elements in the link.
<code>link_</code>	A reference to the element that actually models the link, in terms of the delay and bandwidth characteristics of the link.
<code>ttl_</code>	Reference to the element that manipulates the ttl in every packet.
<code>drophead_</code>	Reference to an object that is the head of a queue of elements that process link drops.

In addition, if the simulator instance variable, `$traceAllFile_`, is defined, the procedure will add trace elements that

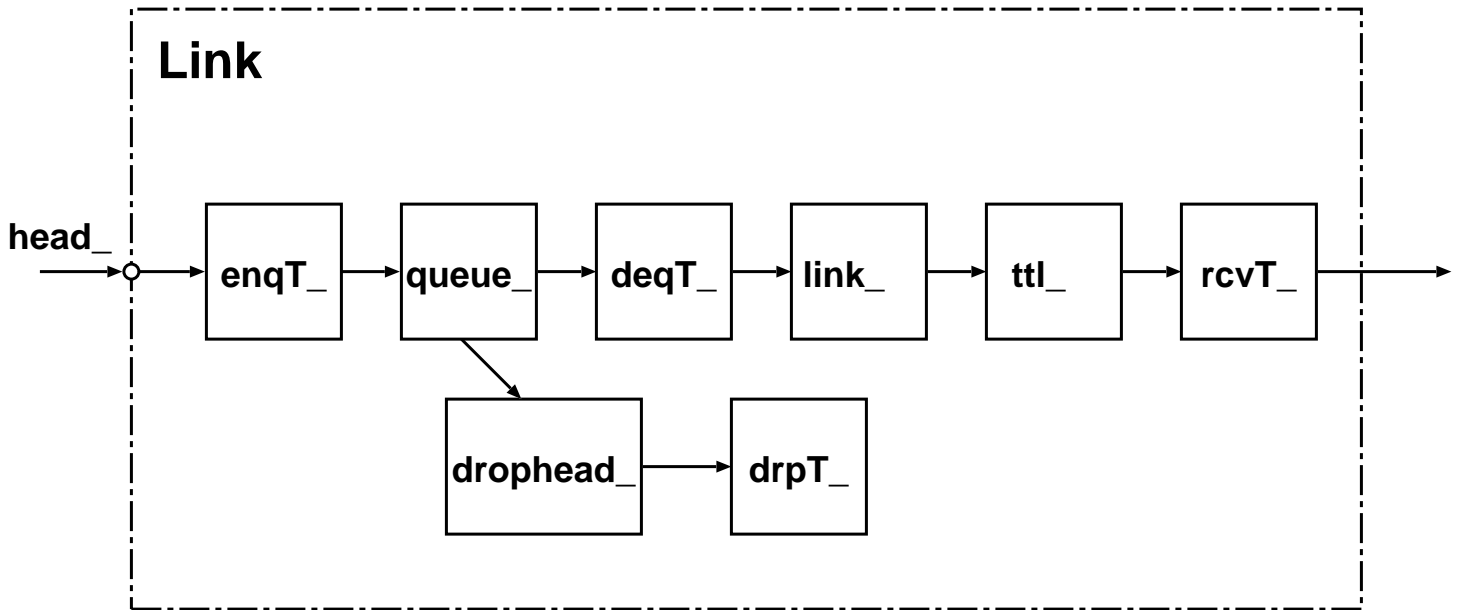


Figure 6.1: Composite Construction of a Unidirectional Link

track when a packet is enqueued and dequeued from `queue_`. Furthermore, tracing interposes a drop trace element after the `drophead_`. The following instance variables track the trace elements:

- `enqT_`    Reference to the element that traces packets entering `queue_`.
- `deqT_`    Reference to the element that traces packets leaving `queue_`.
- `drpT_`    Reference to the element that traces packets dropped from `queue_`.
- `rcvT_`    Reference to the element that traces packets received by the next node.

Note however, that if the user enable tracing multiple times on the link, these instance variables will only store a reference to the last elements inserted.

Other configuration mechanisms that add components to a simple link are network interfaces (used in multicast routing), link dynamics models, and tracing and monitors. We give a brief overview of the related objects at the end of this chapter (Section 6.2), and discuss their functionality/implementation in other chapters.

The instance procedure `duplex-link{}` constructs a bi-directional link from two simplex links.

## 6.1 Instance Procedures for Links and SimpleLinks

**Link procedures** The class `Link` is implemented entirely in Otcl. The Otcl `SimpleLink` class uses the C++ `LinkDelay` class to simulate packet delivery delays. The instance procedures in the class `Link` are:



`head{}` returns the handle for `head_`.  
`queue{}` returns the handle for `queue_`.  
`link{}` returns the handle for the delay element, `link_`.  
`up{}` set link status to “up” in the `dynamics_` element. Also, writes out a trace line to each file specified through the procedure `trace-dynamics{}`.  
`down{}` As with `up{}`, set link status to “down” in the `dynamics_` element. Also, writes out a trace line to each file specified through the procedure `trace-dynamics{}`.  
`up?{}` returns status of the link. Status is “up” or “down”; status is “up” if link dynamics is not enabled.  
`all-connectors{}` Apply specified operation to all connectors on the link.p An example of such usage is `link all-connectors reset`.  
`cost{}` set link cost to value specified.  
`cost?{}` returns the cost of the link. Default cost of link is 1, if no cost has been specified earlier.

**SimpleLink Procedures** The Otcl class `SimpleLink` implements a simple point-to-point link with an associated queue and delay<sup>1</sup>. It is derived from the base Otcl class `Link` as follows:

```

Class SimpleLink -superclass Link
SimpleLink instproc init { src dst bw delay q { lltype "DelayLink" } } {
    $self next $src $dst
    $self instvar link_ queue_ head_ toNode_ ttl_
    ...
    set queue_ $q
    set link_ [new Delay/Link]
    $link_ set bandwidth_ $bw
    $link_ set delay_ $delay

    $queue_ target $link_
    $link_ target [$toNode_ entry]

    ...
    # XXX
    # put the ttl checker after the delay
    # so we don't have to worry about accounting
    # for ttl-drops within the trace and/or monitor
    # fabric
    #
    set ttl_ [new TTLChecker]
    $ttl_ target [$link_ target]
    $link_ target $ttl_
}

```

Notice that when a `SimpleLink` object is created, new `Delay/Link` and `TTLChecker` objects are also created. Note also that, the `Queue` object must have already been created.

There are two additional methods implemented (in OTcl) as part of the `SimpleLink` class: `trace` and `init-monitor`. These functions are described in further detail in the section on tracing (Chapter 25).

---

<sup>1</sup>The current version also includes an object to examine the network layer “ttl” field and discard packets if the field reaches zero.

## 6.2 Connectors

Connectors, unlike classifiers, only generate data for one recipient; either the packet is delivered to the `target_neighbor`, or it is sent to `drop-target`.

A connector will receive a packet, perform some function, and deliver the packet to its neighbor, or drop the packet. There are a number of different types of connectors in *ns*. Each connector performs a different function.

networkinterface	labels packets with incoming interface identifier—it is used by some multicast routing protocols. The class variable “Simulator NumberInterfaces_1” tells <i>ns</i> to add these interfaces, and then, it is added to either end of the simplex link. Multicast routing protocols are discussed in a separate chapter (Chapter 29).
DynaLink	Object that gates traffic depending on whether the link is up or down. It expects to be at the head of the link, and is inserted on the link just prior to simulation start. Its <code>status_</code> variable control whether the link is up or down. The description of how the DynaLink object is used is in a separate chapter (Chapter 30).
DelayLink	Object that models the link’s delay and bandwidth characteristics. If the link is not dynamic, then this object simply schedules receive events for the downstream object for each packet it receives at the appropriate time for that packet. However, if the link is dynamic, then it queues the packets internally, and schedules one receives event for itself for the next packet that must be delivered. Thus, if the link goes down at some point, this object’s <code>reset()</code> method is invoked, and the object will drop all packets in transit at the instant of link failure. We discuss the specifics of this class in another chapter (Chapter 8).
Queues	model the output buffers attached to a link in a “real” router in a network. In <i>ns</i> , they are attached to, and are considered as part of the link. We discuss the details of queues and different types of queues in <i>ns</i> in another chapter (Chapter 7).
TTLChecker	will decrement the ttl in each packet that it receives. If that ttl then has a positive value, the packet is forwarded to the next element on the link. In the simple links, TTLCheckers are automatically added, and are placed as the last element on the link, between the delay element, and the entry for the next node.

## 6.3 Object hierarchy

The base class used to represent links is called `Link`. Methods for this class are listed in the next section. Other link objects derived from the base class are given as follows:

- **SimpleLink Object** A SimpleLink object is used to represent a simple unidirectional link. There are no state variables or configuration parameters associated with this object. Methods for this class are: `$simplelink enable-mcast <src> <dst>`

This turns on multicast for the link by creating an incoming network interface for the destination and adds an outgoing interface for the source.

```
$simplelink trace <ns> <file> <optional:op>
```

Build trace objects for this link and update object linkage. If `op` is specified as "nam" create nam trace files.

```
$simplelink nam-trace <ns> <file>
```

Sets up nam tracing in the link.

```
$simplelink trace-dynamics <ns> <file> <optional:op>
```

This sets up tracing specially for dynamic links. `<op>` allows setting up of nam tracing as well.

```
$simplelink init-monitor <ns> <qtrace> <sampleInterval>
```

Insert objects that allow us to monitor the queue size of this link. Return the name of the object that can be queried to determine the average queue size.

```
$simplelink attach-monitors <insnoop> <outsnoop> <dropsnoop> <qmon>
```

This is similar to init-monitor, but allows for specification of more of the items.

```
$simplelink dynamic
```

Sets up the dynamic flag for this link.

```
$simplelink errormodule <args>
```

Inserts an error module before the queue.

```
$simplelink insert-linkloss <args>
```

Inserts the error module after the queue.

//Other link objects derived from class SimpleLink are FQLink, CBQLink and IntServLink.

Configuration parameters for FQLink are:

**queueManagement\_** The type of queue management used in the link. Default value is DropTail.

No configuration parameters are specified for CBQLink and IntServLink objects.

- **DelayLink Object** The DelayLink Objects determine the amount of time required for a packet to traverse a link. This is defined to be  $\text{size}/\text{bw} + \text{delay}$  where size is the packet size, bw is the link bandwidth and delay is the link propagation delay. There are no methods or state variables associated with this object.

Configuration Parameters are:

**bandwidth\_** Link bandwidth in bits per second.

**delay\_** Link propagation delay in seconds.

## 6.4 Commands at a glance

Following is a list of common link commands used in simulation scripts:

```
$ns_ simplex-link <node1> <node2> <bw> <delay> <qtype> <args>
```

This command creates an unidirectional link between node1 and node2 with specified bandwidth (BW) and delay characteristics. The link uses a queue type of <qtype> and depending on the queue type different arguments are passed through <args>.

```
$ns_ duplex-link <node1> <node2> <bw> <delay> <qtype> <args>
```

This creates a bi-directional link between node1 and node2. This procedure essentially creates a duplex-link from two simplex links, one from node1 to node2 and the other from node2 to node1. The syntax for duplex-link is same as that of simplex-link described above.

```
$ns_ duplex-intserv-link <n1> <n2> <bw> <dly> <sched> <signal> <adc> <args>
```

This creates a duplex-link between n1 and n2 with queue type of intserv, with specified BW and delay. This type of queue implements a scheduler with two level services priority. The type of intserv queue is given by <sched>, with admission control unit type of <adc> and signal module of type <signal>.

```
$ns_ simplex-link-op <n1> <n2> <op> <args>
```

This is used to set attributes for a simplex link. The attributes may be the orientation, color, label, or queue-position.

```
$ns_ duplex-link-op <n1> <n2> <op> <args>
```

This command is used to set link attributes (like orientation of the links, color, label, or queue-position) for duplex links.

`$ns_ link-lossmodel <lossobj> <from> <to>`

This function generates losses (using the loss model <lossobj> inserted in the link between <from> node and <to> node) in the link that can be visualized by nam.

`$ns_ lossmodel <lossobj> <from> <to>`

This is used to insert a loss module in regular links.

Following is a list of internal link-related procedures:

`$ns_ register-nam-linkconfig <link>`

This is an internal procedure used by "`$link orient`" to register/update the order in which links should be created in nam.

`$ns_ remove-nam-linkconfig <id1> <id2>`

This procedure is used to remove any duplicate links (duplicate links may be created by GT-ITM topology generator).

`$link head`

Returns the instance variable `head_` for the link. The `head_` is the entry point to the link and it points to the first object in the link.

`$link add-to-head <connector>`

This allows the <connector> object to be now pointed by the `head_` element in the link, i.e, <connector> now becomes the first object in the link.

`$link link`

Returns the instance variable `link_`. The `link_` is the element in the link that actually models the link in terms of delay and bandwidth characteristics of the link.

`$link queue`

Returns the instance variable `queue_`. `queue_` is queue element in the link. There may be one or more queue elements in a particular link.

`$link cost <c>`

This sets a link cost of <c>.

`$link cost?`

Returns the cost value for the link. Default cost of link is set to 1.

`$link if-label?`

Returns the network interfaces associated with the link (for multicast routing).

`$link up`

This sets the link status to "up". This command is a part of network dynamics support in ns.

`$link down`

Similar to up, this command marks the link status as "down".

`$link up?`

Returns the link status. The status is always "up" as default, if link dynamics is not enabled.

`$link all-connectors op`

This command applies the specified operation <op> to all connectors in the link. Like, `$link all-connectors reset` or `$link all-connectors isDynamic`.

```
$link install-error <errmodel>
```

This installs an error module after the `link_` element.

In addition to the Link and link-related commands listed above, there are other procedures to support the specific requirements of different types of links derived from the base class "Link" like simple-link (SimpleLink), integrated service (IntServLink), class-based queue (CBQLink), fair queue (FQLink) and procedures to support multicast routing, sessionsim, nam etc. These and the above procedures may be found in *ns/tcl/lib*(*ns-lib.tcl*, *ns-link.tcl*, *ns-intserv.tcl*, *ns-namsupp.tcl*, *ns-queue.tcl*), *ns/tcl/mcast*(*McastMonitor.tcl*, *ns-mcast.tcl*), *ns/tcl/session/session.tcl*.

## Chapter 7

# Queue Management and Packet Scheduling

Queues represent locations where packets may be held (or dropped). Packet scheduling refers to the decision process used to choose which packets should be serviced or dropped. Buffer management refers to any particular discipline used to regulate the occupancy of a particular queue. At present, support is included for drop-tail (FIFO) queueing, RED buffer management, CBQ (including a priority and round-robin scheduler), and variants of Fair Queueing including, Fair Queueing (FQ), Stochastic Fair Queueing (SFQ), and Deficit Round-Robin (DRR). In the common case where a *delay* element is downstream from a queue, the queue may be *blocked* until it is re-enabled by its downstream neighbor. This is the mechanism by which transmission delay is simulated. In addition, queues may be forcibly blocked or unblocked at arbitrary times by their neighbors (which is used to implement multi-queue aggregate queues with inter-queue flow control). Packet drops are implemented in such a way that queues contain a “drop destination”; that is, an object that receives all packets dropped by a queue. This can be useful to (for example) keep statistics on dropped packets.

### 7.1 The C++ Queue Class

The `Queue` class is derived from a `Connector` base class. It provides a base class used by particular types of (derived) queue classes, as well as a call-back function to implement blocking (see next section). The following definitions are provided in `queue.h`:

```
class Queue : public Connector {
public:
    virtual void enqueue(Packet*) = 0;
    virtual Packet* deque() = 0;
    void recv(Packet*, Handler*);
    void resume();
    int blocked();
    void unblock();
    void block();
protected:
    Queue();
    int command(int argc, const char*const* argv);
    int qlim_; /* maximum allowed pkts in queue */
    int blocked_;
    int unblock_on_resume_; /* unblock q on idle? */
    QueueHandler qh_;
```

```
};
```

The enqueue and dequeue functions are pure virtual, indicating the Queue class is to be used as a base class; particular queues are derived from Queue and implement these two functions as necessary. Particular queues do not, in general, override the recv function because it invokes the particular enqueue and dequeue.

The Queue class does not contain much internal state. Often these are special monitoring objects (Chapter 25). The qlim\_ member is constructed to dictate a bound on the maximum queue occupancy, but this is not enforced by the Queue class itself; it must be used by the particular queue subclasses if they need this value. The blocked\_ member is a boolean indicating whether the queue is able to send a packet immediately to its downstream neighbor. When a queue is blocked, it is able to enqueue packets but not send them.

### 7.1.1 Queue blocking

A queue may be either blocked or unblocked at any given time. Generally, a queue is blocked when a packet is in transit between it and its downstream neighbor (most of the time if the queue is occupied). A blocked queue will remain blocked as long as its downstream link is busy and the queue has at least one packet to send. A queue becomes unblocked only when its resume function is invoked (by means of a downstream neighbor scheduling it via a callback), usually when no packets are queued. The callback is implemented by using the following class and methods:

```
class QueueHandler : public Handler {
public:
    inline QueueHandler(Queue& q) : queue_(q) {}
    void handle(Event*); /* calls queue_.resume() */
private:
    Queue& queue_;
};

void QueueHandler::handle(Event*)
{
    queue_.resume();
}

Queue::Queue() : drop_(0), blocked_(0), qh_(*this)
{
    Tcl& tcl = Tcl::instance();
    bind("limit_", &qlim_);
}

void Queue::recv(Packet* p, Handler*)
{
    enqueue(p);
    if (!blocked_) {
        /*
         * We're not block. Get a packet and send it on.
         * We perform an extra check because the queue
         * might drop the packet even if it was
         * previously empty! (e.g., RED can do this.)
         */
        p = dequeue();
        if (p != 0) {
            blocked_ = 1;
            target_->recv(p, &qh_);
        }
    }
}
```

```

    }
}
void Queue::resume()
{
    Packet* p = deque();
    if (p != 0)
        target_>recv(p, &qh_);
    else {
        if (unblock_on_resume_)
            blocked_ = 0;
        else
            blocked_ = 1;
    }
}

```

The handler management here is somewhat subtle. When a new Queue object is created, it includes a QueueHandler object (qh\_) which is initialized to contain a reference to the new Queue object (Queue& QueueHandler::queue\_). This is performed by the Queue constructor using the expression qh\_( \*this ). When a Queue receives a packet it calls the subclass (i.e. queueing discipline-specific) version of the enqueue function with the packet. If the queue is not blocked, it is allowed to send a packet and calls the specific deque function which determines which packet to send, blocks the queue (because a packet is now in transit), and sends the packet to the queue's downstream neighbor. Note that any future packets received from upstream neighbors will arrive to a blocked queue. When a downstream neighbor wishes to cause the queue to become unblocked it schedules the QueueHandler's handle function by passing &qh\_ to the simulator scheduler. The handle function invokes resume, which will send the next-scheduled packet downstream (and leave the queue blocked), or unblock the queue when no packet is ready to be sent. This process is made more clear by also referring to the LinkDelay::recv() method (Section 8.1).

### 7.1.2 PacketQueue Class

The Queue class may implement buffer management and scheduling but do not implement the low-level operations on a particular queue. The PacketQueue class is used for this purpose, and is defined as follows (see queue.h):

```

class PacketQueue {
public:
    PacketQueue();
    int length(); /* queue length in packets */
    void enqueue(Packet* p);
    Packet* deque();
    Packet* lookup(int n);
    /* remove a specific packet, which must be in the queue */
    void remove(Packet*);
protected:
    Packet* head_;
    Packet** tail_;
    int len_; /* packet count */
};

```

This class maintains a linked-list of packets, and is commonly used by particular scheduling and buffer management disciplines to hold an ordered set of packets. Particular scheduling or buffer management schemes may make use of several



PacketQueue objects. The PacketQueue class maintains current counts of the number of packets held in the queue which is returned by the length() method. The enqueue function places the specified packet at the end of the queue and updates the len\_ member variable. The dequeue function returns the packet at the head of the queue and removes it from the queue (and updates the counters), or returns NULL if the queue is empty. The lookup function returns the *n*th packet from the head of the queue, or NULL otherwise. The remove function deletes the packet stored in the given address from the queue (and updates the counters). It causes an abnormal program termination if the packet does not exist.

## 7.2 Example: Drop Tail

The following example illustrates the implementation of the Queue/DropTail object, which implements FIFO scheduling and drop-on-overflow buffer management typical of most present-day Internet routers. The following definitions declare the class and its OTcl linkage:

```
/*
 * A bounded, drop-tail queue
 */
class DropTail : public Queue {
protected:
    void enqueue(Packet*);
    Packet* dequeue();
    PacketQueue q_;
};
```

The base class Queue, from which DropTail is derived, provides most of the needed functionality. The drop-tail queue maintains exactly one FIFO queue, implemented by including an object of the PacketQueue class. Drop-tail implements its own versions of enqueue and dequeue as follows:

```
/*
 * drop-tail
 */
void DropTail::enqueue(Packet* p)
{
    q_.enqueue(p);
    if (q_.length() >= qlim_) {
        q_.remove(p);
        drop(p);
    }
}

Packet* DropTail::dequeue()
{
    return (q_.dequeue());
}
```

Here, the enqueue function first stores the packet in the internal packet queue (which has no size restrictions), and then checks the size of the packet queue versus qlim\_. Drop-on-overflow is implemented by dropping the packet most recently added to the packet queue if the limit is reached or exceeded. *Note:* in the implementation of enqueue above, setting qlim\_ to *n* actually means a queue size of *n*−1. Simple FIFO scheduling is implemented in the dequeue function by always returning the first packet in the packet queue.

## 7.3 Different types of Queue objects

A queue object is a general class of object capable of holding and possibly marking or discarding packets as they travel through the simulated topology. Configuration Parameters used for queue objects are:

**limit\_** The queue size in packets.

**blocked\_** Set to false by default, this is true if the queue is blocked (unable to send a packet to its downstream neighbor).

**unblock\_on\_resume\_** Set to true by default, indicates a queue should unblock itself at the time the last packet packet sent has been transmitted (but not necessarily received).

Other queue objects derived from the base class Queue are drop-tail, FQ, SFQ, DRR, RED and CBQ queue objects. Each are described as follows:

- Drop-tail objects: Drop-tail objects are a subclass of Queue objects that implement simple FIFO queue. There are no methods, configuration parameter, or state variables that are specific to drop-tail objects.
- FQ objects: FQ objects are a subclass of Queue objects that implement Fair queuing. There are no methods that are specific to FQ objects. Configuration Parameters are:

**secsPerByte\_**

There are no state variables associated with this object.

- SFQ objects: SFQ objects are a subclass of Queue objects that implement Stochastic Fair queuing. There are no methods that are specific to SFQ objects. Configuration Parameters are:

**maxqueue\_**

**buckets\_**

There are no state variables associated with this object.

- DRR objects: DRR objects are a subclass of Queue objects that implement deficit round robin scheduling. These objects implement deficit round robin scheduling amongst different flows ( A particular flow is one which has packets with the same node and port id OR packets which have the same node id alone). Also unlike other multi-queue objects, this queue object implements a single shared buffer space for its different flows. Configuration Parameters are:

**buckets\_** Indicates the total number of buckets to be used for hashing each of the flows.

**blimit\_** Indicates the shared buffer size in bytes.

**quantum\_** Indicates (in bytes) how much each flow can send during its turn.

**mask\_** mask\_, when set to 1, means that a particular flow consists of packets having the same node id (and possibly different port ids), otherwise a flow consists of packets having the same node and port ids.

- RED objects: RED objects are a subclass of Queue objects that implement random early-detection gateways. The object can be configured to either drop or “mark” packets. There are no methods that are specific to RED objects. Configuration Parameters are:

**bytes\_** Set to "true" to enable “byte-mode” RED, where the size of arriving packets affect the likelihood of marking (dropping) packets.

**queue-in-bytes\_** Set to "true" to measure the average queue size in bytes rather than packets. Enabling this option also causes thresh\_ and maxthresh\_ to be automatically scaled by mean\_pktsize\_ (see below).

**thresh\_** The minimum threshold for the average queue size in packets.

**maxthresh\_** The maximum threshold for the average queue size in packets.

**mean\_pktsize\_** A rough estimate of the average packet size in bytes. Used in updating the calculated average queue size after an idle period.

**q\_weight\_** The queue weight, used in the exponential-weighted moving average for calculating the average queue size.

**wait\_** Set to true to maintain an interval between dropped packets.

**linterm\_** As the average queue size varies between "thresh\_" and "maxthresh\_", the packet dropping probability varies between 0 and "1/linterm".

**setbit\_** Set to "true" to mark packets by setting the congestion indication bit in packet headers rather than drop packets.

**drop-tail\_** Set to true to use drop-tail rather than randomdrop when the queue overflows or the average queue size exceeds "maxthresh\_". For a further explanation of these variables, see [2].

None of the state variables of the RED implementation are accessible.

- CBQ objects: CBQ objects are a subclass of Queue objects that implement class-based queueing.

```
$cbq insert <class>
```

Insert traffic class class into the link-sharing structure associated with link object cbq.

```
$cbq bind <cbqclass> <id1> [$id2]
```

Cause packets containing flow id id1 (or those in the range id1 to id2 inclusive) to be associated with the traffic class cbqclass.

```
$cbq algorithm <alg>
```

Select the CBQ internal algorithm. <alg> may be set to one of: "ancestor-only", "top-level", or "formal".

- CBQ/WRR objects: CBQ/WRR objects are a subclass of CBQ objects that implement weighted round-robin scheduling among classes of the same priority level. In contrast, CBQ objects implement packet-by-packet round-robin scheduling among classes of the same priority level. Configuration Parameters are:

**maxpkt\_** The maximum size of a packet in bytes. This is used only by CBQ/WRR objects in computing maximum bandwidth allocations for the weighted round-robin scheduler.

## CBQCLASS OBJECTS

CBQClass objects implement the traffic classes associated with CBQ objects.

```
$cbqclass setparams <parent> <okborrow> <allot> <maxidle> <prio> <level>
```

Sets several of the configuration parameters for the CBQ traffic class (see below).

```
$cbqclass parent <cbqcl|none>
```

specify the parent of this class in the link-sharing tree. The parent may be specified as "none" to indicate this class is a root.

```
$cbqclass newallot <a>
```

Change the link allocation of this class to the specified amount (in range 0.0 to 1.0). Note that only the specified class is affected.

```
$cbqclass install-queue <q>
```

Install a Queue object into the compound CBQ or CBQ/WRR link structure. When a CBQ object is initially created, it includes no internal queue (only a packet classifier and scheduler).

Configuration Parameters are:

**okborrow\_** is a boolean indicating the class is permitted to borrow bandwidth from its parent.

**allot\_** is the maximum fraction of link bandwidth allocated to the class expressed as a real number between 0.0 and 1.0.

**maxidle\_** is the maximum amount of time a class may be required to have its packets queued before they are permitted to be forwarded

**priority\_** is the class' priority level with respect to other classes. This value may range from 0 to 10, and more than one class may exist at the same priority. Priority 0 is the highest priority.

**level\_** is the level of this class in the link-sharing tree. Leaf nodes in the tree are considered to be at level 1; their parents are at level 2, etc.

**extradelat\_** increase the delay experienced by a delayed class by the specified time

#### QUEUE-MONITOR OBJECTS

QueueMonitor Objects are used to monitor a set of packet and byte arrival, departure and drop counters. It also includes support for aggregate statistics such as average queue size, etc.

`$queuemonitor`

reset all the cumulative counters described below (arrivals, departures, and drops) to zero. Also, reset the integrators and delay sampler, if defined.

`$queuemonitor set-delay-samples <delaySamp_>`

Set up the Samples object `delaySamp_` to record statistics about queue delays. `delaySamp_` is a handle to a Samples object i.e the Samples object should have already been created.

`$queuemonitor get-bytes-integrator`

Returns an Integrator object that can be used to find the integral of the queue size in bytes.

`$queuemonitor get-pkts-integrator`

Returns an Integrator object that can be used to find the integral of the queue size in packets.

`$queuemonitor get-delay-samples`

Returns a Samples object `delaySamp_` to record statistics about queue delays.

There are no configuration parameters specific to this object.

State Variables are:

**size\_** Instantaneous queue size in bytes.

**pkts\_** Instantaneous queue size in packets.

**parrivals\_** Running total of packets that have arrived.

**barrivals\_** Running total of bytes contained in packets that have arrived.

**pdepartures\_** Running total of packets that have departed (not dropped).

**bdepartures\_** Running total of bytes contained in packets that have departed (not dropped).

**pdrops\_** Total number of packets dropped.

**bdrops\_** Total number of bytes dropped.

**bytesInt\_** Integrator object that computes the integral of the queue size in bytes. The `sum_` variable of this object has the running sum (integral) of the queue size in bytes.

**pktsInt\_** Integrator object that computes the integral of the queue size in packets. The `sum_` variable of this object has the running sum (integral) of the queue size in packets.

## QUEUEMONITOR/ED OBJECTS

This derived object is capable of differentiating regular packet drops from early drops. Some queues distinguish regular drops (e.g. drops due to buffer exhaustion) from other drops (e.g. random drops in RED queues). Under some circumstances, it is useful to distinguish these two types of drops.

State Variables are:

**epdrops\_** The number of packets that have been dropped “early”.

**ebdrops\_** The number of bytes comprising packets that have been dropped “early”.

Note: because this class is a subclass of QueueMonitor, objects of this type also have fields such as pdrops\_ and bdrops\_. These fields describe the total number of dropped packets and bytes, including both early and non-early drops.

## QUEUEMONITOR/ED/FLOWMON OBJECTS

These objects may be used in the place of a conventional QueueMonitor object when wishing to collect per-flow counts and statistics in addition to the aggregate counts and statistics provided by the basic QueueMonitor.

`$fmon classifier <cl>`

This inserts (read) the specified classifier into (from) the flow monitor object. This is used to map incoming packets to which flows they are associated with.

`$fmon dump`

Dump the current per-flow counters and statistics to the I/O channel specified in a previous attach operation.

`$fmon flows`

Return a character string containing the names of all flow objects known by this flow monitor. Each of these objects are of type QueueMonitor/ED/Flow.

`$fmon attach <chan>`

Attach a tcl I/O channel to the flow monitor. Flow statistics are written to the channel when the dump operation is executed.

Configuration Parameters are:

**enable\_in\_** Set to true by default, indicates that per-flow arrival state should be kept by the flow monitor. If set to false, only the aggregate arrival information is kept.

**enable\_out\_** Set to true by default, indicates that per-flow departure state should be kept by the flow monitor. If set to false, only the aggregate departure information is kept.

**enable\_drop\_** Set to true by default, indicates that per-flow drop state should be kept by the flow monitor. If set to false, only the aggregate drop information is kept.

**enable\_edrop\_** Set to true by default, indicates that per-flow early drop state should be kept by the flow monitor. If set to false, only the aggregate early drop information is kept.

## QUEUEMONITOR/ED/FLOW OBJECTS

These objects contain per-flow counts and statistics managed by a QueueMonitor/ED/Flowmon object. They are generally created in an OTcl callback procedure when a flow monitor is given a packet it cannot map on to a known flow. Note that the flow monitor’s classifier is responsible for mapping packets to flows in some arbitrary way. Thus, depending on the type of classifier used, not all of the state variables may be relevant (e.g. one may classify packets based only on flow id, in which case the source and destination addresses may not be significant). State Variables are:

**src\_** The source address of packets belonging to this flow.

**dst\_** The destination address of packets belonging to this flow.

**flowid\_** The flow id of packets belonging to this flow.

## 7.4 Commands at a glance

Following is a list of queue commands used in simulation scripts:

```
$ns_ queue-limit <n1> <n2> <limit>
```

This sets a limit on the maximum buffer size of the queue in the link between nodes <n1> and <n2>.

```
$ns_ trace-queue <n1> <n2> <optional:file>
```

This sets up trace objects to log events in the queue. If tracefile is not passed, it uses traceAllFile\_ to write the events.

```
$ns_ namtrace-queue <n1> <n2> <optional:file>
```

Similar to trace-queue above, this sets up nam-tracing in the queue.

```
$ns_ monitor-queue <n1> <n2> <optional:qtrace> <optional:sampleinterval>
```

This command inserts objects that allows us to monitor the queue size. This returns a handle to the object that may be queried to determine the average queue size. The default value for sampleinterval is 0.1.

## 7.5 Queue/JoBS

*JoBS is developed and contributed by Nicolas Christin <nicolas@cs.virginia.edu>*

This chapter describes the implementation of the Joint Buffer Management and Scheduling (JoBS) algorithm in *ns*. This chapter is in three parts. The first part summarizes the objectives of the JoBS algorithm. The second part explains how to configure a JoBS queue in *ns*. The third part focuses on the tracing mechanisms implemented for JoBS.

The procedures and functions described in this chapter can be found in *ns/jobs.{cc, h}*, *ns/marker.{cc, h}*, *ns/demarker.{cc, h}*. Example scripts can be found in *ns/tcl/ex/jobs-{lossdel, cn2002}.tcl*.

Additional information can be found at <http://qosbox.cs.virginia.edu>.

### 7.5.1 The JoBS algorithm

This section gives an overview of the objectives the JoBS algorithm aims at achieving, and of the mechanisms employed to reach these objectives. The original JoBS algorithm, as described in [20], was using the solution to a non-linear optimization problem. This *ns-2* implementation uses the feedback-control based heuristic described in [8].

*Important Note:* This *ns-2* implementation results from the merge between old code for *ns-2.1b5*, and code derived from the BSD kernel-level implementation of the JoBS algorithm. **It is still considered experimental.** Due to the absence of binding facilities for arrays between Tcl and C++ in *tclcl* at the moment, *the number of traffic classes is statically set to 4 and cannot be changed without modifying the C++ code.*

## Objective

The objective of the JoBS algorithm is to provide absolute and relative (proportional) loss and delay differentiation independently at each node for *classes* of traffic. JoBS therefore provides service guarantees on a *per-hop* basis. The set of performance requirements are specified to the algorithm as a set of per-class Quality of Service (QoS) constraints. As an example, for three classes, the QoS constraints could be of the form:

- Class-1 Delay  $\approx 2 \cdot$  Class-2 Delay,
- Class-2 Loss Rate  $\approx 10^{-1} \cdot$  Class-3 Loss Rate, or
- Class-3 Delay  $\leq 5$  ms.

Here, the first two constraints are relative constraints and the last one is an absolute constraint. The set of constraints can be any mix of relative and absolute constraints. More specifically, JoBS supports the five following types of constraints:

- **Relative delay constraints (RDC)** specify a proportional delay differentiation between classes. As an example, for two classes 1 and 2, the RDC enforces a relationship

$$\frac{\text{Delay of Class 2}}{\text{Delay of Class 1}} \approx \text{constant}.$$

- **Absolute delay constraints (ADC)**: An ADC on class  $i$  requires that the delays of class  $i$  satisfy a worst-case bound  $d_i$ .
- **Relative loss constraints (RLC)** specify a proportional loss differentiation between classes.
- **Absolute loss constraints (ALC)**: An ALC on class  $i$  requires that the loss rate of class  $i$  be bounded by an upper bound  $L_i$ .
- **Absolute rate constraints (ARC)**: An ARC on class  $i$  means that the throughput of class  $i$  is bounded by a lower bound  $\mu_i$ .

JoBS does not rely on admission control or traffic policing, nor does it make any assumption on traffic arrivals. Therefore, a system of constraints may become infeasible, and some constraints may need to be relaxed. QoS constraints are prioritized in the following order.

$$\text{ALC} > \text{ADC}, \text{ARC} > \text{Relative Constraints}.$$

That is, if JoBS is unable to satisfy both absolute and relative constraints, it will give preference to the absolute constraints.

## Mechanisms

JoBS performs scheduling and buffer management in a single pass. JoBS dynamically allocates service rates to classes in order to satisfy the delay constraints. The service rates needed for enforcing absolute delay constraints are allocated upon each packet arrival, while service rates derived from relative delay constraints are computed only every  $N$  packet arrivals. If no feasible service rate allocation exists<sup>1</sup>, or if the packet buffer overflows, packets are dropped according to the loss constraints.

The service rates are translated into packet scheduling decisions by an algorithm resembling Deficit Round Robin. That is, the scheduler tries to achieve the desired service rates by keeping track of the difference between the actual transmission rate for each class and the desired service rate for each class. Scheduling in JoBS is work-conserving.

---

<sup>1</sup>For instance, if the sum of the service rates needed is greater than the output link capacity.

## 7.5.2 Configuration

Running a JoBS simulation requires to create and configure the JoBS “link(s)”, to create and configure the Markers and Demarkers in charge of marking/demarking the traffic, to attach an application-level data source (traffic generator), and to start the traffic generator.

### Initial Setup

```
set ns [new Simulator]                                ; # preamble initialization

Queue/JoBS set drop_front_ false                      ; # use drop-tail
Queue/JoBS set trace_hop_ true                        ; # enable statistic traces
Queue/JoBS set adc_resolution_type_ 0                 ; # see “commands at a glance”
Queue/JoBS set shared_buffer_ 1                      ; # all classes share a common buffer
Queue/JoBS set mean_pkt_size_ 4000                   ; # we expect to receive 500-Byte pkts
Queue/Demarker set demarker_arrvs1_ 0                ; # reset arrivals everywhere
Queue/Demarker set demarker_arrvs2_ 0
Queue/Demarker set demarker_arrvs3_ 0
Queue/Demarker set demarker_arrvs4_ 0
Queue/Marker set marker_arrvs1_ 0
Queue/Marker set marker_arrvs2_ 0
Queue/Marker set marker_arrvs3_ 0
Queue/Marker set marker_arrvs4_ 0

set router(1) [$ns node]                             ; # set first router
set router(2) [$ns node]                             ; # set second router
set source [$ns node]                                ; # set source
set sink [$ns node]                                  ; # set traffic sink
set bw 100000000                                     ; # 10 Mbps
set delay 0.001                                       ; # 1 ms
set buff 500                                          ; # 500 packets
```

### Creating the JoBS links

```
$ns duplex-link $router(1) $router(2) $bw $delay JoBS ; # Creates the JoBS link
$ns_ queue-limit $router(1) $router(2) $buff
set l [$ns_ get-link $router(1) $router(2)]
set q [$l queue]
$q init-rdcs -1 2 2 2 ; # Classes 2, 3 and 4 are bound by proportional delay differentiation with a factor of 2
$q init-rlcs -1 2 2 2 ; # Classes 2, 3 and 4 are bound by proportional loss differentiation with a factor of 2
$q init-alcs 0.01 -1 -1 -1 ; # Class 1 is provided with a loss rate bound of 1%
$q init-adcs 0.005 -1 -1 -1 ; # Class 1 is provided with a delay bound of 5 ms
$q init-arcs -1 -1 -1 500000 ; # Class 4 is provided with a minimum throughput of 500 Kbps
$q link [$l link] ; # The link is attached to the queue (required)
$q trace-file jobstrace ; # Trace per-hop, per-class metrics to the file jobstrace
$q sampling-period 1 ; # Reevaluate rate allocation upon each arrival
$q id 1 ; # Assigns an ID of 1 to the JoBS queue
$q initialize ; # Proceed with the initialization
```



## Marking the traffic

Marking the traffic is handled by Marker objects. Markers are FIFO queues that set the class index of each packet. To ensure accuracy of the simulations, it is best to configure these queues to have a very large buffer, so that no packets are dropped in the Marker. Demarkers are used to gather end-to-end delay statistics.

```
$ns_ simplex-link $source $router(1) $bw $delay Marker ; # set-up marker
$ns_ queue-limit $source $router(1) [expr $buff*10] ; # Select huge buffers for markers
$ns_ queue-limit $router(1) $source [expr $buff*10] ; # to avoid traffic drops
set q [$ns_ get-queue $source $router(1)] ; #in the marker
$q marker_type 2 ; # Statistical marker
$q marker_frc 0.1 0.2 0.3 0.4 ; # 10% Class 1, 20% Class 2, 30% Class 3, 40% Class 4.
$ns_ simplex-link $router(2) $sink $bw $delay Demarker ; # set-up demarker
$ns_ queue-limit $router(2) $sink [expr $buff*10]
$q trace-file e2e ; # trace end-to-end delays to file e2e
```

The remaining steps (attaching agents and traffic generators or applications to the nodes) are explained in Chapters 10 and 37, and are handled as usual. We refer to these chapters and the example scripts provided with your *ns* distribution.

### 7.5.3 Tracing

Tracing in JoBS is handled internally, by the scheduler. Each JoBS queue can generate a trace file containing the following information. Each line of the tracefile consists of 17 columns. The first column is the simulation time, columns 2 to 5 represent the loss rates over the current busy period for classes 1 to 4, columns 6 to 9 represent the delays for each class (average over a 0.5 seconds sliding window), columns 10 to 13 represent the average service rates allocated to each class over the last 0.5 seconds, and columns 14 to 17 represent the instantaneous queue length in packets. Additionally, Demarkers can be used to trace end-to-end delays.

### 7.5.4 Variables

This section summarizes the variables that are used by JoBS, Marker and Demarker objects.

#### JoBS objects

**trace\_hop\_** Can be true or false. If set to true, per-hop, per-class metrics will be traced. (Trace files have then to be specified, using <JoBS object> trace-file <filename>.) Defaults to false.

**drop\_front\_** Can be true or false. If set to true, traffic will be dropped from the front of the queue. Defaults to false (drop-tail).

**adc\_resolution\_type\_** Can be 0 or 1. If set to 0, traffic will be dropped from classes that have an ADC if the ADC cannot be met by adjusting the service rates. If set to 1, traffic will be dropped from all classes. A resolution mode set to 1 is therefore fairer, in the sense that the pain is shared by all classes, but can lead to more deadline violations. Defaults to 0.

**shared\_buffer\_** Can be 0 or 1. If set to 0, all classes use a separate per-class buffer (which is required if only rate guarantees are to be provided). All per-class buffers have the same size. If set to 1, all classes share the same buffer (which is required if loss differentiation is to be provided). Defaults to 1.

**mean\_pkt\_size\_** Used to set the expected mean packet size of packets arriving at a JoBS link. Setting this variable is required to ensure proper delay differentiation.

### Marker objects

**marker\_arrvs1\_** Number of Class-1 packets to have entered a Marker link.

**marker\_arrvs2\_** Number of Class-2 packets to have entered a Marker link.

**marker\_arrvs3\_** Number of Class-3 packets to have entered a Marker link.

**marker\_arrvs4\_** Number of Class-4 packets to have entered a Marker link.

### Demarker objects

**demarker\_arrvs1\_** Number of Class-1 packets to have entered a Demarker link.

**demarker\_arrvs2\_** Number of Class-2 packets to have entered a Demarker link.

**demarker\_arrvs3\_** Number of Class-3 packets to have entered a Demarker link.

**demarker\_arrvs4\_** Number of Class-4 packets to have entered a Demarker link.

## 7.5.5 Commands at a glance

The following is a list of commands used to configure the JoBS, Marker and Demarker objects.

### JoBS objects

```
set q [new Queue/JoBS]
```

This creates an instance of the JoBS queue.

```
$q init-rdcs <k1> <k2> <k3> <k4>
```

This assigns the RDCs for the four JoBS classes. For instance, using a value of 4 for k2 means that Class-3 delays will be roughly equal to four times Class-2 delays. A value of -1 indicates that the class is not concerned by RDCs.

*Important Note:* Since RDCs bound two classes, one would expect only three parameters to be passed (k1, k2, and k3, since k4 theoretically binds Classes 4 and 5, and Class 5 does not exist). However, in this prototype implementation, it is imperative to specify a value different from 0 and -1 to k4 if Class 4 is to be concerned by RDCs.

*Examples:* `$q init-rdcs -1 2 1 -1` specifies that classes 2 and 3 are bound by a delay differentiation factor of 2, `$q init-rdcs 4 4 4 4` specifies that all classes are bound by a delay differentiation factor of 4 and is equivalent to `$q init-rdcs 4 4 4 1`, since the last coefficient is only used to specify that Class 4 is to be bound by proportional differentiation.

```
$q init-rlcs <k'1> <k'2> <k'3> <k'4>
```

This assigns the RLCs for the four JoBS classes. For instance, using a value of 3 for k1 means that Class-2 loss rates will be roughly equal to four times Class-2 loss rates. A value of -1 indicates that the class is not concerned by RLCs. As with RDCs, each RLC binds two classes, thus, one would expect only three parameters to be passed (k'1, k'2, and k'3, since k'4

theoretically bounds Classes 4 and 5, and Class 5 does not exist). As explained above, it is imperative to specify a value different from 0 and -1 to k'4 if Class 4 is to be concerned by RLCs.

```
$q init-alcs <L1> <L2> <L3> <L4>
```

This assigns the absolute loss guarantees (ALCs) to all four classes. L1 to L4 are given in fraction of 1. For instance, setting L1 to 0.05 means that Class-1 loss rate will be guaranteed to be less than 5%. A value of -1 indicates that the corresponding class is not subject to an ALC.

```
$q init-adcs <D1> <D2> <D3> <D4>
```

This assigns the absolute loss guarantees (ADCs) to all four classes. D1 to D4 are given in milliseconds. A value of -1 indicates that the corresponding class is not subject to an ADC.

```
$q trace-file <filename>
```

This specifies the trace file for all per-hop metrics. JoBS uses an internal module to trace loss and delays, service rates, and per-class queue lengths in packets. If filename is set to **null**, no trace will be provided.

```
$q link [<link-object> link]
```

This command is required to bind a link to a JoBS queue. Note that JoBS needs to know the capacity of the link. Thus, this command **has to** be issued before the simulation is started.

```
$q sampling-period <sampling-interval>
```

This command specifies the sampling interval (in packets) at which the service rate adjustments for proportional differentiation will be performed. The default is a sampling interval of 1 packet, meaning that the rate allocation is reevaluated upon each packet arrival. Larger sampling intervals speed up the simulations, but typically result in poorer proportional differentiation.

```
$q id <num_id>
```

This command affects a numerical ID to the JoBS queue.

```
$q initialize
```

This command is required, and should be run after all configuration operations have been performed. This command will perform the final checks and configuration of the JoBS queue.

```
$q copyright-info
```

Displays authors and copyright information.

A simple example script (with nam output), fully annotated and commented can be found in *ns/tcl/ex/jobs-lossdel.tcl*. A more realistic example of a simulation with JoBS queues can be found in *ns/tcl/ex/jobs-cn2002.tcl*. This script is very similar to what was used in a simulation presented in [21]. Associated tracefiles and *gnuplot* scripts for visualization (in case you favor *gnuplot* over *xgraph* can be found in *ns/tcl/ex/jobs-lossdel*, and *ns/tcl/ex/jobs-cn2002*.

## Marker objects

```
$q marker_type <1|2>
```

Selects the type of marker. 1 is DETERMINISTIC, 2 is STATISTICAL.

```
$q marker_class <1|2|3|4>
```

For a deterministic marker, selects which class packets should be marked with.

```
$q marker_frc <f1> <f2> <f3> <f4>
```

For a statistical marker, gives the fraction of packets that should be marked from each class. For instance, using 0.1 for f1 means that 10 percent of the traffic coming to the Marker link will be marked as Class 1.

## Demarker objects

```
$q trace-file <filename>
```

This command specifies the trace file used for the demarker object. filename.1 will contain the end-to-end delays of each Class-1 packet to have reached the Demarker link, filename.2 will contain the end-to-end delays of each Class-2 packet to have reached the Demarker link, and so forth. (There will of course be 4 trace files, one for each class.)

## Chapter 8

# Delays and Links

Delays represent the time required for a packet to traverse a link. A special form of this object (“dynamic link”) also captures the possibility of a link failure. The amount of time required for a packet to traverse a link is defined to be  $s/b + d$  where  $s$  is the packet size (as recorded in its IP header),  $b$  is the speed of the link in bits/sec, and  $d$  is the link delay in seconds. The implementation of link delays is closely associated with the blocking procedures described for Queues in Section 7.1.1.

### 8.1 The LinkDelay Class

The class `LinkDelay` is derived from the base class `Connector`. Its definition is in `~ns/delay.cc`, and is briefly excerpted below:

```
class LinkDelay : public Connector {
public:
    LinkDelay();
    void recv(Packet* p, Handler*);
    void send(Packet* p, Handler*);
    void handle(Event* e);
    double delay();                                /* line latency on this link */
    double bandwidth();                            /* bandwidth on this link */
    inline double txtime(Packet* p) {              /* time to send pkt p on this link */
        hdr_cmh* hdr = (hdr_cmh*) p->access(off_cmh_);
        return (hdr->size() * 8. / bandwidth_);
    }

protected:
    double bandwidth_;                            /* bandwidth of underlying link (bits/sec) */
    double delay_;                                /* line latency */
    int dynamic_;                                 /* indicates whether or not link is ~ */
    Event inTransit_;
    PacketQueue* itq_;                            /* internal packet queue for dynamic links */
    Packet* nextPacket_;                          /* to be delivered for a dynamic link. */
    Event intr_;
};
```

The `recv()` method overrides the base class `Connector` version. It is defined as follows:

```
void LinkDelay::recv(Packet* p, Handler* h)
{
    double txt = txttime(p);
    Scheduler& s = Scheduler::instance();
    if (dynamic_) {
        Event* e = (Event*)p;
        e->time_ = s.clock() + txt + delay_;
        itq_>enqueue(p);
        schedule_next();
    } else {
        s.schedule(target_, p, txt + delay_);
    }
    /* XXX only need one intr_ since upstream object should
     * block until it's handler is called
     *
     * This only holds if the link is not dynamic. If it is, then
     * the link itself will hold the packet, and call the upstream
     * object at the appropriate time. This second interrupt is
     * called inTransit_, and is invoked through schedule_next()
     */
    s.schedule(h, &intr_, txt);
}
```

This object supports one `instproc`-like, `$object dynamic`, to set its variable, `dynamic_`. This variable determines whether the link is dynamic or not (*i.e.*, prone to fail/recover at appropriate times). The internal behavior of the link in each case is different.

For “non-dynamic” links, this method operates by receiving a packet,  $p$ , and scheduling two events. Assume these two events are called  $E_1$  and  $E_2$ , and that event  $E_1$  is scheduled to occur before  $E_2$ .  $E_1$  is scheduled to occur when the upstream node attached to this delay element has completed sending the current packet (which takes time equal to the packet size divided by the link bandwidth).  $E_1$  is usually associated with a `Queue` object, and will cause it to (possibly) become unblocked (see section 7.1.1).  $E_2$  represents the packet arrival event at the downstream neighbor of the delay element. Event  $E_2$  occurs a number of seconds later than  $E_1$  equal to the link delay.

Alternately, when the link is dynamic, and receives  $p$ , then it will schedule  $E_1$  to possibly unblock the queue at the appropriate time. However,  $E_2$  is scheduled only if  $p$  is the only packet currently in transit. Otherwise, there is at least one packet in transit on the link that must be delivered before  $p$  at  $E_2$ . Therefore, packet  $p$  is held in the object’s `inTransit` queue, `itq_`. When the packet just before  $p$  in transit on the link is delivered at the neighbor node, the `DelayLink` object will schedule an event for itself to fire at  $E_2$ . At that appropriate time then, its `handle()` method will directly send  $p$  to its target. The object’s internal `schedule_next()` method will schedule these events for packet  $p$  in transit at the appropriate time.

## 8.2 Commands at a glance

The `LinkDelay` object represents the time required by a packet to transverse the link and is used internally within a `Link`. Hence we donot list any `linkdelay` related commands suitable for simulation scripts here.

## Chapter 9

# Differentiated Services Module in *ns*

**Note:** The Differentiated Services module described in this chapter has been integrated into *ns-2.1b8*.

Differentiated Services, or DiffServ, is an IP QoS architecture based on packet marking that allows packets to be prioritized according to user requirements. During the time of congestion, more low priority packets are discarded than high priority packets. This chapter describes the DiffServ module that was originally implemented by the Advanced IP Networks group in Nortel Networks [28].

### 9.1 Overview

The DiffServ architecture provides QoS by dividing traffic into different categories, marking each packet with a code point that indicates its category, and scheduling packets accordingly. The DiffServ module in *ns* can support four classes of traffic, each of which has three dropping precedences allowing differential treatment of traffic within a single class. Packets in a single class of traffic are enqueued into one corresponding physical RED queue, which contains three virtual queues (one for each drop precedence).

Different RED parameters can be configured for virtual queues, causing packets from one virtual queue to be dropped more frequently than packets from another. A packet with a lower dropping precedence is given better treatment in times of congestion because it is assigned a code point that corresponds to a virtual queue with relatively lenient RED parameters.

The DiffServ module in *ns* has three major components:

**Policy:** Policy is specified by network administrator about the level of service a class of traffic should receive in the network.

**Edge router:** Edge router marks packets with a code point according to the policy specified.

**Core router:** Core router examines packets' code point marking and forwarding them accordingly.

DiffServ attempts to restrict complexity to only the edge routers.

## 9.2 Implementation

The procedures and functions described in this section can be found in `~ns/diffserv/dsred`, `dsredq`, `dsEdge`, `dsCore`, `dsPolicy`.`{cc, h}`.

### 9.2.1 RED queue in DiffServ module

A DiffServ queue (in class `dsREDQueue`) derived from the base class `Queue` is implemented in DiffServ module to provide the basic DiffServ router functionality, see `dsred`.`{h, cc}`. `dsREDQueue` has the following abilities:

- to implement multiple physical RED queues along a single link;
- to implement multiple virtual queues within a physical queue, with individual set of parameters for each virtual queue;
- to determine in which physical and virtual queue a packet is enqueued according to its code point;
- to determine in from which physical and virtual queue a packet is dequeued according to the scheduling scheme chosen.

The class `dsREDQueue` consists of four physical RED queues, each containing three virtual queues. The number of physical and virtual queues are defined in `numPrec` and `numQueues_`. Each combination of physical and virtual queue number is associated with a code point (or a drop preference), which specifies a certain level of service.

The physical queue is defined in class `redQueue`, which enables traffic differentiation by defining virtual queues with independent configuration and state parameters, see `dsredq`.`{h, cc}`. For example, the length of each virtual queue is calculated only on packets mapped to that queue. Thus, packet dropping decisions can be applied based on the state and configuration parameters of that virtual queues. Class `redQueue` is not equivalent to class `REDQueue`, which was already present in *ns*. Instead, it is a modified version of RED implementation with the notion of virtual queues and is only used by class `redQueue` to realize physical queues. All user interaction with class `redQueue` is handled through the command interface of class `dsREDQueue`.

Class `dsREDQueue` contains a data structure known as the Per Hop Behavior (PHB) Table. In DiffServ, edge routers mark packets with code points and core routers simply respond to existing code points; both of them use PHB table to map a code point to a particular physical and virtual queue. The PHB Table is defined as an array with three fields:

```
struct phbParam {
    int codePt_; // corresponding code point
    int queue_; // physical queue
    int prec_; // virtual queue (drop precedence)
};
```

### 9.2.2 Edge and core routers

The DiffServ edge and core routers are defined in class `edgeQueue` and class `coreQueue`, which are derived from class `dsREDQueue`, see `dsEdge`, `dsCore`.`{h, cc}`.

Packet marking is implemented in class `edgeQueue`. A packet is marked with a code point according to the policy specified before it is put into the corresponding physical and virtual queue. Class `edgeQueue` has a reference to an instance of class `PolicyClassifier`, which contains policies for packet marking.



### 9.2.3 Policy

Class `Policy` and its sub-classes (see `dsPolicy.{cc, h}`) define the policies used by edge routers to mark incoming packets. A policy is established between a source and destination node. All flows matching that source-destination pair are treated as a single traffic aggregate. Policy for each different traffic aggregate has an associated policer type, meter type, and initial code point. The meter type specifies the method for measuring the state variables needed by the policer. For example, the TSW Tagger is a meter that measures the average traffic rate, using a specified time window.

When a packet arrives at an edge router, it is examined to determine to which aggregate it belongs. The meter specified by the corresponding policy is invoked to update all state variables. The policer is invoked to determine how to mark the packet depending on the aggregate's state variables: the specified initial code point or a downgraded code point. Then the packet is enqueued accordingly.

Currently, six different policy models are defined:

1. Time Sliding Window with 2 Color Marking (TSW2CMPolicer): uses a CIR and two drop precedences. The lower precedence is used probabilistically when the CIR is exceeded.
2. Time Sliding Window with 3 Color Marking (TSW3CMPolicer): uses a CIR, a PIR, and three drop precedences. The medium drop precedence is used probabilistically when the CIR is exceeded and the lowest drop precedence is used probabilistically when the PIR is exceeded.
3. Token Bucket (tokenBucketPolicer): uses a CIR and a CBS and two drop precedences. An arriving packet is marked with the lower precedence if and only if it is larger than the token bucket.
4. Single Rate Three Color Marker (srTCMPolicer): uses a CIR, CBS, and an EBS to choose from three drop precedences.
5. Two Rate Three Color Marker (trTCMPolicer): uses a CIR, CBS, PIR, and a PBS to choose from three drop precedences.
6. NullPolicer: does not downgrade any packets

The policies above are defined as sub-classes of `dsPolicy`. The specific meter and policer are implemented in functions `applyMeter` and `applyPolicer`, which are defined as virtual functions in class `dsPolicy`. User specified policy can be added in the similar way. Please refer to `NullPolicy` as the simplest example.

All policies are stored in the policy table in class `PolicyClassifier`. This table is an array that includes fields for the source and destination nodes, a policer type, a meter type, an initial code point, and various state information as shown below:

The rates CIR and PIR are specified in bits per second:

CIR: committed information rate

PIR: peak information rate

The buckets CBS, EBS, and PBS are specified in bytes:

CBS: committed burst size

EBS: excess burst size

PBS: peak burst size

C bucket: current size of the committed bucket

E bucket: current size of the excess bucket

P bucket: current size of the peak bucket

Arrival time of last packet

Average sending rate

TSW window length

Class `PolicyClassifier` also contains a `Policer Table` to store the mappings from a policy type and initial code point pair to its associated downgraded code point(s).

## 9.3 Configuration

The number of physical and virtual queues can be configured as:

```
$dsredq set numQueues_ 1
```

```
$dsredq setNumPrec 2
```

Variable `numQueues_` in class `dsREDQueue` specifies the number of physical queues. It has a default value as 4 defined in `~ns/tcl/lib/ns-default.tcl` and can be changed as shown in the example above. Variable `setNumPrec` sets the number of virtual queues within one physical queue.

RED parameters can be configured for each virtual queue as follows:

```
$dsredq configQ 0 1 10 20 0.10
```

The mean packet size (in bytes) is also needed for the average RED queue length calculation.

```
$dsredq meanPktSize 1500
```

The variant of MRED used to calculate queue sizes can be configured.

```
$dsredq setMREDMode RIO-C 0
```

The above command sets the MRED mode of physical queue 0 to RIO-C. If the second argument was not included, all queues would be set to RIO-C which is the default.

The various MRED modes supported in DiffServ module are:

**RIO-C (RIO Coupled):** The probability of dropping an out-of-profile packet is based on the weighted average lengths of all virtual queues; while the probability of dropping an in-profile packet is based solely on the weighted average length of its virtual queue.

**RIO-D (RIO De-coupled):** Similar to RIO-C; except the probability of dropping an out-of-profile packet is based on the size of its virtual queue.

**WRED (Weighted RED):** All probabilities are based on a single queue length.

**DROP:** Same as a drop tail queue with queue limit set by RED minimum threshold: when the queue size reaches the minimum threshold, all packets are dropped regardless of marking.

The following command adds an entry to the PHB Table and maps code point 11 to physical queue 0 and virtual queue 1.

```
$dsredq addPHBEntry 11 0 1
```

In *ns*, packets are defaulted to a code point of zero. Therefore, user must add a PHB entry for the zero code point in order to handle best effort traffic.

In addition, commands are available to allow the user to choose the scheduling mode between physical queues. For example:

```
$dsredq setSchedulerMode WRR
```

```
$dsredq addQueueWeights 1 5
```

The above pair of commands sets the scheduling mode to Weighted Round Robin and the weight for queue 1 to 5. Other scheduling modes supported are Weighted Interleaved Round Robin (WIRR), Round Robin (RR), and Priority (PRI). The default scheduling mode is Round Robin.

For Priority scheduling, priority is arranged in sequential order with queue 0 having the highest priority. Also, one can set the a limit on the maximum bandwidth a particular queue can get using as follows:

```
$dsredq setSchedulerMode PRI
```

```
$dsredq addQueueRate 0 5000000
```

These commands specify the maximum bandwidth that queue 0 can consume is 5Mb.

The `addPolicyEntry` command is used to add an entry to the Policy Table. It takes different parameters depending on what policer type is used. The first two parameters after the command name are always the source and destination node IDs, and the next parameter is the policer type. Following the policer type are the parameters needed by that policer as summarized below:

Null	Initial code point				
TSW2CM	Initial code point	CIR			
TSW3CM	Initial code point	CIR	PIR		
TokenBucket	Initial code point	CIR	CBS		
srTCM	Initial code point	CIR	CBS	EBS	
trTCM	Initial code point	CIR	CBS	PIR	PBS

Note that the Null policer requires only the initial code point. Since this policer does not downgrade packets, other information is not necessary. Consider a Tcl script for which `$q` is a variable for an edge queue, and `$s` and `$d` are source and destination nodes. The following command adds a TSW2CM policer for traffic going from the source to the destination:

```
$q addPolicyEntry [$s id] [$d id] TSW2CM 10 2000000
```

Other parameters could be used for different policers in place of "TSW2CM":

Null	10				
TSW3CM	10	2000000	3000000		
TokenBucket	10	2000000	10000		
srTCM	10	2000000	10000	20000	
trTCM	10	2000000	10000	3000000	10000

Note, however, that only one policy can be applied to any source destination pair.

The following command adds an entry to the Policer Table, specifying that the trTCM has initial (green) code point 10, downgraded (yellow) code point 11 and further downgraded (red) code point 12:

```
$dsredq addPolicerEntry trTCM 10 11 12
```

There must be a Policer Table entry in place for every policer type and initial code point pair. It should be noticed that the Null policer is added in the following way:

```
$dsredq addPolicerEntry Null 10
```

Downgrade code points are not specified because the Null policy does not meter traffic characteristics.

Queries supported:

Output entries in Policy Table, one line at a time:

```
$dsredq printPolicyTable
```

Output entries in Policer Table, one line at a time:

```
$dsredq printPolicerTable
```

Output entries in PHB table, one line at a time:

```
$dsredq printPHBTable
```

Packets statistic results:

```
$dsredq printStats
```

Sample output:

Packets Statistics				
CP	TotPkts	TxPkts	ldrops	edrops
All	249126	249090	21	15
10	150305	150300	0	5
20	98821	98790	21	10

CP: code point

TotPkts: packets received

TxPkts: packets sent

ldrops: packets are dropped due to link overflow

edrops: RED early dropping).

Returns the RED weighted average size of the specified physical queue:

```
$dsredq getAverage 0
```

Returns the current size of the C Bucket (in bytes):

```
$dsredq getCBucket
```

## 9.4 Commands at a glance

The following is a list of related commands commonly used in simulation scripts:

```
$ns simplex-link $edge $core 10Mb 5ms dsRED/edge
$ns simplex-link $core $edge 10Mb 5ms dsRED/core
```

These two commands create the queues along the link between an edge router and a core router.

```
set qEC [$ns link $edge $core] queue]

# Set DS RED parameters from Edge to Core:
$qEC meanPktSize $packetSize
$qEC set numQueues_ 1
$qEC setNumPrec 2
$qEC addPolicyEntry [$s1 id] [$dest id] TokenBucket 10 $cir0 $cbs0
$qEC addPolicyEntry [$s2 id] [$dest id] TokenBucket 10 $cir1 $cbs1
$qEC addPolicerEntry TokenBucket 10 11
$qEC addPHBEntry 10 0 0
$qEC addPHBEntry 11 0 1
$qEC configQ 0 0 20 40 0.02
$qEC configQ 0 1 10 20 0.10
```

This block of code obtains handle to the DiffServ queue from an edge router to a core router and configures all of the parameters for it.

The meanPktSize command is required for the RED state variables to be calculated accurately. Setting the number of physical queues and precedence levels is optional, but it aids efficiency. Because neither the scheduling or MRED mode type are set, they default to Round Robin scheduling and RIO-C Active Queue Management.

The addPolicyEntry commands establish two policies at the edge queue: one between nodes S1 and Dest and one between nodes S2 and Dest. Note that the [\$s1 id] command returns the ID value needed by addPolicyEntry. The CIR and CBS values used in the policies are the ones set at the beginning of the script.

The addPolicerEntry line is required because each policer type and initial code point pair requires an entry in the Policer Table. Each of the policies uses the same policer and initial code point, so only one entry is needed.

The addPHBEntry commands map each code point to a combination of physical and virtual queue. Although each code point in this example maps to a unique combination of physical and virtual queue, multiple code points could receive identical treatment.

Finally, the configQ commands set the RED parameters for each virtual queue. It specifies the virtual queue by first two parameters, for example, 0 and 1. The next three parameters are the minimum threshold, maximum threshold, and the maximum dropping probability. Note that as the precedence value increases, the RED parameters become harsher.

```
set qCE [$ns link $core $el] queue]

# Set DS RED parameters from Core to Edge:
$qCE meanPktSize $packetSize
$qCE set numQueues_ 1
$qCE setNumPrec 2
$qCE addPHBEntry 10 0 0
$qCE addPHBEntry 11 0 1
$qCE configQ 0 0 20 40 0.02
```

```
$qCE configQ 0 1 10 20 0.10
```

Note that the configuration of a core queue matches that of an edge queue, except that there is no Policy Table or Policer Table to configure at a core router. A core router's chief requirement is that it has a PHB entry for all code points that it will see.

```
$qE1C printPolicyTable  
$qCE2 printCoreStats
```

These methods output the policy or policer tables on link and different statistics.

For further information, please refer to the example scripts under *~ns/tcl/ex/diffserv*.

%

# Chapter 10

## Agents

Agents represent endpoints where network-layer packets are constructed or consumed, and are used in the implementation of protocols at various layers. The `class Agent` has an implementation partly in OTcl and partly in C++. The C++ implementation is contained in `~ns/agent.cc` and `~ns/agent.h`, and the OTcl support is in `~ns/tcl/lib/ns-agent.tcl`.

### 10.1 Agent state

The C++ `class Agent` includes enough internal state to assign various fields to a simulated packet before it is sent. This state includes the following:

<code>addr_</code>	node address of myself (source address in packets)
<code>dst_</code>	where I am sending packets to
<code>size_</code>	packet size in bytes (placed into the common packet header)
<code>type_</code>	type of packet (in the common header, see <code>packet.h</code> )
<code>fid_</code>	the IP flow identifier (formerly <i>class</i> in ns-1)
<code>prio_</code>	the IP priority field
<code>flags_</code>	packet flags (similar to ns-1)
<code>defttl_</code>	default IP ttl value

These variables may be modified by any class derived from `Agent`, although not all of them may be needed by any particular agent.

### 10.2 Agent methods

The `class Agent` supports packet generation and reception. The following member functions are implemented by the C++ `Agent` class, and are generally *not* over-ridden by derived classes:

<code>Packet* allocpkt()</code>	allocate new packet and assign its fields
<code>Packet* allocpkt(int)</code>	allocate new packet with a data payload of <code>n</code> bytes and assign its fields



The following member functions are also defined by the class `Agent`, but *are* intended to be over-ridden by classes deriving from `Agent`:

```
void timeout(timeout number)  subclass-specific time out method
void recv(Packet*, Handler*)  receiving agent main receive path
```

The `allocpkt()` method is used by derived classes to create packets to send. The function fills in the following fields in the common packet header (Section 12): `uid`, `ptype`, `size`, and the following fields in the IP header: `src`, `dst`, `flowid`, `prio`, `ttl`. It also zero-fills in the following fields of the Flags header: `ecn`, `pri`, `usr1`, `usr2`. Any packet header information not included in these lists must be handled in the classes derived from `Agent`.

The `recv()` method is the main entry point for an `Agent` which receives packets, and is invoked by upstream nodes when sending a packet. In most cases, Agents make no use of the second argument (the handler defined by upstream nodes).

## 10.3 Protocol Agents

There are several agents supported in the simulator. These are their names in OTcl:

TCP	a "Tahoe" TCP sender (cwnd = 1 on any loss)
TCP/Reno	a "Reno" TCP sender (with fast recovery)
TCP/Newreno	a modified Reno TCP sender (changes fast recovery)
TCP/Sack1	a SACK TCP sender
TCP/Fack	a "forward" SACK sender TCP
TCP/FullTcp	a more full-functioned TCP with 2-way traffic
TCP/Vegas	a "Vegas" TCP sender
TCP/Vegas/RBP	a Vegas TCP with "rate based pacing"
TCP/Vegas/RBP	a Reno TCP with "rate based pacing"
TCP/Asym	an experimental Tahoe TCP for asymmetric links
TCP/Reno/Asym	an experimental Reno TCP for asymmetric links
TCP/Newreno/Asym	an experimental Newreno TCP for asymmetric links
TCPSink	a Reno or Tahoe TCP receiver (not used for FullTcp)
TCPSink/DelAck	a TCP delayed-ACK receiver
TCPSink/Asym	an experimental TCP sink for asymmetric links
TCPSink/Sack1	a SACK TCP receiver
TCPSink/Sack1/DelAck	a delayed-ACK SACK TCP receiver
UDP	a basic UDP agent
RTP	an RTP sender and receiver
RTCP	an RTCP sender and receiver
LossMonitor	a packet sink which checks for losses
IVS/Source	an IVS source
IVS/Receiver	an IVS receiver

CtrMcast/Encap	a “centralised multicast” encapsulator
CtrMcast/Decap	a “centralised multicast” de-encapsulator
Message	a protocol to carry textual messages
Message/Prune	processes multicast routing prune messages
SRM	an SRM agent with non-adaptive timers
SRM/Adaptive	an SRM agent with adaptive timers
Tap	interfaces the simulator to a live network
Null	a degenerate agent which discards packets
rtProto/DV	distance-vector routing protocol agent

Agents are used in the implementation of protocols at various layers. Thus, for some transport protocols (e.g. UDP) the distribution of packet sizes and/or inter-departure times may be dictated by some separate object representing the demands of an application. To this end, agents expose an application programming interface (API) to the application. For agents used in the implementation of lower-layer protocols (e.g. routing agents), size and departure timing is generally dictated by the agent’s own processing of protocol messages.

## 10.4 OTcl Linkage

Agents may be created within OTcl and an agent’s internal state can be modified by use of Tcl’s `set` function and any Tcl functions an Agent (or its base classes) implements. Note that some of an Agent’s internal state may exist only within OTcl, and is thus is not directly accessible from C++.

### 10.4.1 Creating and Manipulating Agents

The following example illustrates the creation and modification of an Agent in OTcl:

```
set newtcp [new Agent/TCP]           ; # create new object (and C++ shadow object)
$newtcp set window_ 20               ; # sets the tcp agent's window to 20
$newtcp target $dest                 ; # target is implemented in Connector class
$newtcp set portID_ 1                ; # exists only in OTcl, not in C++
```

### 10.4.2 Default Values

Default values for member variables, those visible in OTcl only and those linked between OTcl and C++ with `bind` are initialized in the `~ns/tcl/lib/ns-default.tcl` file. For example, `Agent` is initialized as follows:

```
Agent set fid_ 0
Agent set prio_ 0
Agent set addr_ 0
```

```

Agent set dst_ 0
Agent set flags_ 0

```

Generally these initializations are placed in the OTcl namespace before any objects of these types are created. Thus, when an Agent object is created, the calls to bind in the objects' constructors will cause the corresponding member variables to be set to these specified defaults.

### 10.4.3 OTcl Methods

The instance procedures defined for the OTcl Agent class are currently found in `~ns/tcl/lib/ns-agent.tcl`. They are as follows:

```

port      the agent's port identifier
dst-port  the destination's port identifier
attach-source <stype>  create and attach a Source object to an agent

```

## 10.5 Examples: Tcp, TCP Sink Agents

The class `TCP` represents a simplified TCP sender. It sends data to a `TCPSink` agent and processes its acknowledgments. It has a separate object associated with it which represents an application's demand. By looking at the class `TCPAgent` and class `TCPSinkAgent`, we may see how relatively complex agents are constructed. An example from the Tahoe TCP agent `TCPAgent` is also given to illustrate the use of timers.

### 10.5.1 Creating the Agent

The following OTcl code fragment creates a TCP agent and sets it up:

```

set tcp [new Agent/TCP]                ; # create sender agent
$tcp set fid_ 2                        ; # set IP-layer flow ID
set sink [new Agent/TCPSink]          ; # create receiver agent
$ns attach-agent $n0 $tcp             ; # put sender on node $n0
$ns attach-agent $n3 $sink            ; # put receiver on node $n3
$ns connect $tcp $sink                ; # establish TCP connection
set ftp [new Application/FTP]         ; # create an FTP source "application"
$ftp attach-agent $tcp                ; # associate FTP with the TCP sender
$ns at 1.2 "$ftp start"               ; # arrange for FTP to start at time 1.2 sec

```

The OTcl instruction `new Agent/TCP` results in the creation of a C++ `TcpAgent` class object. Its constructor first invokes the constructor of the Agent base class and then performs its own bindings. These two constructors appear as follows:

The `TcpSimpleAgent` constructor (`~ns/tcp.cc`):

```

TcpAgent::TcpAgent() : Agent(PT_TCP), rtt_active_(0), rtt_seq_(-1),
                    rtx_timer_(this), delsnd_timer_(this)
{
    bind("window_", &wnd_);
    bind("windowInit_", &wnd_init_);
}

```

```

        bind("windowOption_", &wnd_option_);
        bind("windowConstant_", &wnd_const_);
        ...
        bind("off_ip_", &off_ip_);
        bind("off_tcp_", &off_tcp_);
        ...
    }

```

The Agent constructor (*~ns/agent.cc*):

```

Agent::Agent(int pkttype) :
    addr_(-1), dst_(-1), size_(0), type_(pkttype), fid_(-1),
    prio_(-1), flags_(0)
{
    memset(pending_, 0, sizeof(pending_));           /* timers */
    //  this is really an IP agent, so set up
    //  for generating the appropriate IP fields...
    bind("addr_", (int*)&addr_);
    bind("dst_", (int*)&dst_);
    bind("fid_", (int*)&fid_);
    bind("prio_", (int*)&prio_);
    bind("flags_", (int*)&flags_);
    ...
}

```

These code fragments illustrate the common case where an agent's constructor passes a packet type identifier to the Agent constructor. The values for the various packet types are used by the packet tracing facility (Section 25.5) and are defined in *~ns/trace.h*. The variables which are bound in the TcpAgent constructor are ordinary instance/member variables for the class with the exception of the special integer values *off\_tcp\_* and *off\_ip\_*. These are needed in order to access a TCP header and IP header, respectively. Additional details are in the section on packet headers (Section 12.1).

Note that the TcpAgent constructor contains initializations for two timers, *rtx\_timer\_* and *delsnd\_timer\_*.

TimerHandler objects are initialized by providing a pointer (the *this* pointer) to the relevant agent.

## 10.5.2 Starting the Agent

The TcpAgent agent is started in the example when its FTP source receives the *start* directive at time 1.2. The *start* operation is an instance procedure defined on the class Application/FTP (Section 37.4). It is defined in *~ns/tcl/lib/ns-source.tcl* as follows:

```

Application/FTP instproc start {} {
    [$self agent] send -1
}

```

In this case, *agent* refers to our simple TCP agent and *send -1* is analogous to sending an arbitrarily large file.

The call to *send* eventually results in the simple TCP sender generating packets. The following function output performs this:

```

void TcpAgent::output(int seqno, int reason)

```

```

{
    Packet* p = allocpkt();
    hdr_tcp *tcph = (hdr_tcp*)p->access(off_tcp_);
    double now = Scheduler::instance().clock();
    tcph->seqno() = seqno;
    tcph->ts() = now;
    tcph->reason() = reason;
    Connector::send(p, 0);
    ...
    if (!(rtx_timer_.status() == TIMER_PENDING))
        /* No timer pending. Schedule one. */
        set_rtx_timer();
}

```

Here we see an illustration of the use of the `Agent::allocpkt()` method. This output routine first allocates a new packet (with its common and IP headers already filled in), but then must fill in the appropriate TCP-layer header fields. To find the TCP header in a packet (assuming it has been enabled (Section 12.2.4)) the `off_tcp_` must be properly initialized, as illustrated in the constructor. The packet `access()` method returns a pointer to the TCP header, its sequence number and time stamp fields are filled in, and the `send()` method of the class `Connector` is called to send the packet downstream one hop. Note that the C++ `::` scoping operator is used here to avoid calling `TcpSimpleAgent::send()` (which is also defined). The check for a pending timer uses the timer method `status()` which is defined in the base class `TimerHandler`. It is used here to set a retransmission timer if one is not already set (a TCP sender only sets one timer per window of packets on each connection).

### 10.5.3 Processing Input at Receiver

Many of the TCP agents can be used with the class `TCPSink` as the peer. This class defines the `recv()` and `ack()` methods as follows:

```

void TcpSink::recv(Packet* pkt, Handler*)
{
    hdr_tcp *th = (hdr_tcp*)pkt->access(off_tcp_);
    acker_->update(th->seqno());
    ack(pkt);
    Packet::free(pkt);
}

void TcpSink::ack(Packet* opkt)
{
    Packet* npkt = allocpkt();

    hdr_tcp *otcp = (hdr_tcp*)opkt->access(off_tcp_);
    hdr_tcp *ntcp = (hdr_tcp*)npkt->access(off_tcp_);
    ntcp->seqno() = acker_->Seqno();
    ntcp->ts() = otcp->ts();

    hdr_ip* oip = (hdr_ip*)opkt->access(off_ip_);
    hdr_ip* nip = (hdr_ip*)npkt->access(off_ip_);
    nip->flowid() = oip->flowid();

    hdr_flags* of = (hdr_flags*)opkt->access(off_flags_);

```

```

        hdr_flags* nf = (hdr_flags*)npkt->access(off_flags_);
        nf->ecn_ = of->ecn_;

        acker_->append_ack((hdr_cmh*)npkt->access(off_cmh_),
                           ntcp, otcp->seqno());
        send(npkt, 0);
    }

```

The `recv()` method overrides the `Agent::recv()` method (which merely discards the received packet). It updates some internal state with the sequence number of the received packet (and therefore requires the `off_tcp_` variable to be properly initialized). It then generates an acknowledgment for the received packet. The `ack()` method makes liberal use of access to packet header fields including separate accesses to the TCP header, IP header, Flags header, and common header. The call to `send()` invokes the `Connector::send()` method.

### 10.5.4 Processing Responses at the Sender

Once the simple TCP's peer receives data and generates an ACK, the sender must (usually) process the ACK. In the `TcpAgent` agent, this is done as follows:

```

/*
 * main reception path - should only see acks, otherwise the
 * network connections are misconfigured
 */
void TcpAgent::recv(Packet *pkt, Handler*)
{
    hdr_tcp *tcph = (hdr_tcp*)pkt->access(off_tcp_);
    hdr_ip* iph = (hdr_ip*)pkt->access(off_ip_);
    ...
    if ((hdr_flags*)pkt->access(off_flags_->ecn_)
        quench(1);
    if (tcph->seqno() > last_ack_) {
        newack(pkt);
        openwnd();
    } else if (tcph->seqno() == last_ack_) {
        if (++dupacks_ == NUMDUPACKS) {
            ...
        }
    }
    Packet::free(pkt);
    send(0, 0, maxburst_);
}

```

This routine is invoked when an ACK arrives at the sender. In this case, once the information in the ACK is processed (by `newack`) the packet is no longer needed and is returned to the packet memory allocator. In addition, the receipt of the ACK indicates the possibility of sending additional data, so the `TcpSimpleAgent::send()` method is invoked which attempts to send more data if the TCP window allows.

## 10.5.5 Implementing Timers

As described in the following chapter (Chapter 11), specific timer classes must be derived from an abstract base class `TimerHandler` defined in `~ns/timer-handler.h`. Instances of these subclasses can then be used as various agent timers. An agent may wish to override the `Agent::timeout()` method (which does nothing). In the case of the Tahoe TCP agent, two timers are used: a delayed send timer `delsnd_timer_` and a retransmission timer `rtx_timer_`. We describe the retransmission timer in TCP (Section 11.1.2) as an example of timer usage.

## 10.6 Creating a New Agent

To create a new agent, one has to do the following:

1. decide its inheritance structure (Section 10.6.1), and create the appropriate class definitions,
2. define the `recv()` and `timeout()` methods (Section 10.6.2),
3. define any necessary timer classes,
4. define OTcl linkage functions (Section 10.6.3),
5. write the necessary OTcl code to access the agent (Section 10.6.4).

The action required to create an agent can be illustrated by means of a very simple example. Suppose we wish to construct an agent which performs the ICMP ECHO REQUEST/REPLY (or “ping”) operations.

### 10.6.1 Example: A “ping” requestor (Inheritance Structure)

Deciding on the inheritance structure is a matter of personal choice, but is likely to be related to the layer at which the agent will operate and its assumptions on lower layer functionality. The simplest type of Agent, connectionless datagram-oriented transport, is the `Agent/UDP` base class. Traffic generators can easily be connected to UDP Agents. For protocols wishing to use a connection-oriented stream transport (like TCP), the various TCP Agents could be used. Finally, if a new transport or “sub-transport” protocol is to be developed, using `Agent` as the base class would likely be the best choice. In our example, we’ll use `Agent` as the base class, given that we are constructing an agent logically belonging to the IP layer (or just above it).

We may use the following class definitions:

```
class ECHO_Timer;

class ECHO_Agent : public Agent {
public:
    ECHO_Agent();
    int command(int argc, const char*const* argv);
protected:
    void timeout(int);
    void sendit();
    double interval_;
    ECHO_Timer echo_timer_;
};

class ECHO_Timer : public TimerHandler {
```

```

public:
    ECHO_Timer(ECHO_Agent *a) : TimerHandler() { a_ = a; }
protected:
    virtual void expire(Event *e);
    ECHO_Agent *a_;
};

```

## 10.6.2 The `recv()` and `timeout()` Methods

The `recv()` method is not defined here, as this agent represents a request function and will generally not be receiving events or packets<sup>1</sup>. By not defining the `recv()` method, the base class version of `recv()` (*i.e.*, `Connector::recv()`) is used. The `timeout()` method is used to periodically send request packets. The following `timeout()` method is used, along with a helper method, `sendit()`:

```

void ECHO_Agent::timeout(int)
{
    sendit();
    echo_timer_.resched(interval_);
}

void ECHO_Agent::sendit()
{
    Packet* p = allocpkt();
    ECHOHeader *eh = ECHOHeader::access(p->bits());
    eh->timestamp() = Scheduler::instance().clock();
    send(p, 0);      // Connector::send()
}

void ECHO_Timer::expire(Event *e)
{
    a_->timeout(0);
}

```

The `timeout()` method simply arranges for `sendit()` to be executed every `interval_` seconds. The `sendit()` method creates a new packet with most of its header fields already set up by `allocpkt()`. The packet is only lacking the current timestamp. The call to `access()` provides for a structured interface to the packet header fields, and is used to set the timestamp field. Note that this agent uses its own special header (“ECHOHeader”). The creation and use of packet headers is described in later chapter (Chapter 12); to send the packet to the next downstream node, `Connector::send()` is invoked without a handler.

## 10.6.3 Linking the “ping” Agent with OTcl

We have the methods and mechanisms for establishing OTcl Linkage earlier (Chapter 3). This section is a brief review of the essential features of that earlier chapter, and describes the minimum functionality required to create the ping agent.

There are three items we must handle to properly link our agent with OTcl. First we need to establish a mapping between the OTcl name for our class and the actual object created when an instantiation of the class is requested in OTcl. This is done as follows:

---

<sup>1</sup>This is perhaps unrealistically simple. An ICMP ECHO REQUEST agent would likely wish to process ECHO REPLY messages.



```

static class ECHOClass : public TclClass {
public:
    ECHOClass() : TclClass("Agent/ECHO") {}
    TclObject* create(int argc, const char*const* argv) {
        return (new ECHO_Agent());
    }
} class_echo;

```

Here, a *static* object “class\_echo” is created. Its constructor (executed immediately when the simulator is executed) places the class name “Agent/ECHO” into the OTcl name space. The mixing of case is by convention; recall from Section 3.5 in the earlier chapters that the “/” character is a hierarchy delimiter for the interpreted hierarchy. The definition of the `create()` method specifies how a C++ shadow object should be created when the OTcl interpreter is instructed to create an object of class “Agent/ECHO”. In this case, a dynamically-allocated object is returned. This is the normal way new C++ shadow objects are created.

Once we have the object creation set up, we will want to link C++ member variables with corresponding variables in the OTcl namespace, so that accesses to OTcl variables are actually backed by member variables in C++. Assume we would like OTcl to be able to adjust the sending interval and the packet size. This is accomplished in the class’s constructor:

```

ECHO_Agent::ECHO_Agent() : Agent(PT_ECHO)
{
    bind_time("interval_", &interval_);
    bind("packetSize_", &size_);
}

```

Here, the C++ variables `interval_` and `size_` are linked to the OTcl instance variables `interval_` and `packetSize_`, respectively. Any read or modify operation to the OTcl variables will result in a corresponding access to the underlying C++ variables. The details of the `bind()` methods are described elsewhere (Section 3.4.2). The defined constant `PT_ECHO` is passed to the `Agent()` constructor so that the `Agent::allocpkt()` method may set the packet type field used by the trace support (Section 25.5). In this case, `PT_ECHO` represents a new packet type and must be defined in `~ns/trace.h` (Section 25.4).

Once object creation and variable binding is set up, we may want to create methods implemented in C++ but which can be invoked from OTcl (Section 3.4.4). These are often control functions that initiate, terminate or modify behavior. In our present example, we may wish to be able to start the ping query agent from OTcl using a “start” directive. This may be implemented as follows:

```

int ECHO_Agent::command(int argc, const char*const* argv)
{
    if (argc == 2) {
        if (strcmp(argv[1], "start") == 0) {
            timeout(0);
            return (TCL_OK);
        }
    }
    return (Agent::command(argc, argv));
}

```

Here, the `start()` method available to OTcl simply calls the C++ member function `timeout()` which initiates the first packet generation and schedules the next. Note this class is so simple it does not even include a way to be stopped.

## 10.6.4 Using the agent through OTcl

The agent we have created will have to be instantiated and attached to a node. Note that a node and simulator object is assumed to have already been created. The following OTcl code performs these functions:

```
set echoagent [new Agent/ECHO]
$simulator attach-agent $node $echoagent
```

To set the interval and packet size, and start packet generation, the following OTcl code is executed:

```
$echoagent set dst_ $dest
$echoagent set fid_ 0
$echoagent set prio_ 0
$echoagent set flags_ 0
$echoagent set interval_ 1.5
$echoagent set packetSize_ 1024
$echoagent start
```

This will cause our agent to generate one 1024-byte packet destined for node \$dest every 1.5 seconds.

## 10.7 The Agent API

Simulated applications may be implemented on top of protocol agents. Chapter 37 describes the API used by applications to access the services provided by the protocol agent.

## 10.8 Different agent objects

Class Agent forms the base class from which different types of objects like Nullobject, TCP etc are derived. The methods for Agent class are described in the next section. Configuration parameters for:

**fid\_** Flowid.

**prio\_** Priority.

**agent\_addr\_** Address of this agent.

**agent\_port\_** Port address of this agent.

**dst\_addr\_** Destination address for the agent.

**dst\_port\_** Destination port address for the agent.

**flags\_**

**ttl\_** TTL defaults to 32.

There are no state variables specific to the generic agent class. Other objects derived from Agent are given below:

**Null Objects** Null objects are a subclass of agent objects that implement a traffic sink. They inherit all of the generic agent object functionality. There are no methods specific to this object. The state variables are:

- sport\_
- dport\_

**LossMonitor Objects** LossMonitor objects are a subclass of agent objects that implement a traffic sink which also maintains some statistics about the received data e.g., number of bytes received, number of packets lost etc. They inherit all of the generic agent object functionality.

`$lossmonitor clear`

Resets the expected sequence number to -1.

State Variables are:

**nlost\_** Number of packets lost.

**npkts\_** Number of packets received.

**bytes\_** Number of bytes received.

**lastPktTime\_** Time at which the last packet was received.

**expected\_** The expected sequence number of the next packet.

**TCP objects** TCP objects are a subclass of agent objects that implement the BSD Tahoe TCP transport protocol as described in paper: "Fall, K., and Floyd, S. Comparisons of Tahoe, Reno, and Sack TCP. December 1995." URL <ftp://ftp.ee.lbl.gov/papers/sacks.ps.Z>. They inherit all of the generic agent functionality. Configuration Parameters are:

**window\_** The upper bound on the advertised window for the TCP connection.

**maxcwnd\_** The upper bound on the congestion window for the TCP connection. Set to zero to ignore. (This is the default.)

**windowInit\_** The initial size of the congestion window on slow-start.

**windowOption\_** The algorithm to use for managing the congestion window.

**windowThresh\_** Gain constant to exponential averaging filter used to compute awnd (see below). For investigations of different window-increase algorithms.

**overhead\_** The range of a uniform random variable used to delay each output packet. The idea is to insert random delays at the source in order to avoid phase effects, when desired [see Floyd, S., and Jacobson, V. On Traffic Phase Effects in Packet-Switched Gateways. Internetworking: Research and Experience, V.3 N.3, September 1992. pp. 115-156 ]. This has only been implemented for the Tahoe ("tcp") version of tcp, not for tcp-reno. This is not intended to be a realistic model of CPU processing overhead.

**ecn\_** Set to true to use explicit congestion notification in addition to packet drops to signal congestion. This allows a Fast Retransmit after a quench() due to an ECN (explicit congestion notification) bit.

**packetSize\_** The size in bytes to use for all packets from this source.

**tcpTick\_** The TCP clock granularity for measuring roundtrip times. Note that it is set by default to the non-standard value of 100ms.

**bugFix\_** Set to true to remove a bug when multiple fast retransmits are allowed for packets dropped in a single window of data.

**maxburst\_** Set to zero to ignore. Otherwise, the maximum number of packets that the source can send in response to a single incoming ACK.

**slow\_start\_restart\_** Set to 1 to slow-start after the connection goes idle. On by default.

Defined Constants are:

**MWS** The Maximum Window Size in packets for a TCP connection. MWS determines the size of an array in `tcp-sink.cc`. The default for MWS is 1024 packets. For Tahoe TCP, the "window" parameter, representing the receiver's advertised window, should be less than MWS-1. For Reno TCP, the "window" parameter should be less than (MWS-1)/2.

State Variables are:

**dupacks\_** Number of duplicate acks seen since any new data was acknowledged.

**seqno\_** Highest sequence number for data from data source to TCP.

**t\_seqno\_** Current transmit sequence number.

**ack\_** Highest acknowledgment seen from receiver. **cwnd\_** Current value of the congestion window.

**awnd\_** Current value of a low-pass filtered version of the congestion window. For investigations of different window-increase algorithms.

**ssthresh\_** Current value of the slow-start threshold.

**rtt\_** Round-trip time estimate.

**srtt\_** Smoothed round-trip time estimate.

**rttvar\_** Round-trip time mean deviation estimate.

**backoff\_** Round-trip time exponential backoff constant.

**TCP/Reno Objects** TCP/Reno objects are a subclass of TCP objects that implement the Reno TCP transport protocol described in paper: "Fall, K., and Floyd, S. Comparisons of Tahoe, Reno, and Sack TCP. December 1995." URL <ftp://ftp.ee.lbl.gov/papers/sacks.ps.Z>. There are no methods, configuration parameters or state variables specific to this object.

**TCP/Newreno Objects** TCP/Newreno objects are a subclass of TCP objects that implement a modified version of the BSD Reno TCP transport protocol. There are no methods or state variables specific to this object.

Configuration Parameters are:

**newreno\_changes\_** Set to zero for the default New Reno described in "Fall, K., and Floyd, S. Comparisons of Tahoe, Reno, and Sack TCP. December 1995". Set to 1 for additional New Reno algorithms [see Hoe, J., Improving the Start-up Behavior of a Congestion Control Scheme for TCP. in SIGCOMM 96, August 1996, pp. 270-280. URL <http://www.acm.org/sigcomm/sigcomm96/papers/hoe.html>.]; this includes the estimation of the ssthresh parameter during slow-start.

**TCP/Vegas Objects** There are no methods or configuration parameters specific to this object. State variables are:

- **v\_alpha\_**
- **v\_beta\_**
- **v\_gamma\_**
- **v\_rtt\_**

**TCP/Sack1 Objects** TCP/Sack1 objects are a subclass of TCP objects that implement the BSD Reno TCP transport protocol with Selective Acknowledgement Extensions described in "Fall, K., and Floyd, S. Comparisons of Tahoe, Reno, and Sack TCP. December 1995". URL <ftp://ftp.ee.lbl.gov/papers/sacks.ps.Z>. They inherit all of the TCP object functionality. There are no methods, configuration parameters or state variables specific to this object.

**TCP/FACK Objects** TCP/Fack objects are a subclass of TCP objects that implement the BSD Reno TCP transport protocol with Forward Acknowledgement congestion control. They inherit all of the TCP object functionality. There are no methods or state variables specific to this object.

Configuration Parameters are:

**ss-div4** Overdamping algorithm. Divides ssthresh by 4 (instead of 2) if congestion is detected within 1/2 RTT of slow-start. (1=Enable, 0=Disable)

**rampdown** Rampdown data smoothing algorithm. Slowly reduces congestion window rather than instantly halving it. (1=Enable, 0=Disable)

**TCP/FULLTCP Objects** This section has not yet been added here. The implementation and the configuration parameters are described in paper: "Fall, K., Floyd, S., and Henderson, T., Ns Simulator Tests for Reno FullTCP. July, 1997." URL <ftp://ftp.ee.lbl.gov/papers/fulltcp.ps>.

**TCPSINK Objects** TCPSink objects are a subclass of agent objects that implement a receiver for TCP packets. The simulator only implements "one-way" TCP connections, where the TCP source sends data packets and the TCP sink sends ACK packets. TCPSink objects inherit all of the generic agent functionality. There are no methods or state variables specific to the TCPSink object. Configuration Parameters are

**packetSize\_** The size in bytes to use for all acknowledgment packets.

**maxSackBlocks\_** The maximum number of blocks of data that can be acknowledged in a SACK option. For a receiver that is also using the time stamp option [RFC 1323], the SACK option specified in RFC 2018 has room to include three SACK blocks. This is only used by the TCPSink/Sack1 subclass. This value may not be increased within any particular TCPSink object after that object has been allocated. (Once a TCPSink object has been allocated, the value of this parameter may be decreased but not increased).

**TCPSINK/DELACK Objects** DelAck objects are a subclass of TCPSink that implement a delayed-ACK receiver for TCP packets. They inherit all of the TCPSink object functionality. There are no methods or state variables specific to the DelAck object. Configuration Parameters are:

**interval\_** The amount of time to delay before generating an acknowledgment for a single packet. If another packet arrives before this time expires, generate an acknowledgment immediately.

**TCPSINK/SACK1 Objects** TCPSink/Sack1 objects are a subclass of TCPSink that implement a SACK receiver for TCP packets. They inherit all of the TCPSink object functionality. There are no methods, configuration parameters or state variables specific to this object.

**TCPSINK/SACK1/DELACK Objects** TCPSink/Sack1/DelAck objects are a subclass of TCPSink/Sack1 that implement a delayed-SACK receiver for TCP packets. They inherit all of the TCPSink/Sack1 object functionality. There are no methods or state variables specific to this object. Configuration Parameters are:

**interval\_** The amount of time to delay before generating an acknowledgment for a single packet. If another packet arrives before this time expires, generate an acknowledgment immediately.

## 10.9 Commands at a glance

Following are the agent related commands used in simulation scripts:

```
ns_ attach-agent <node> <agent>
```

This command attaches the <agent> to the <node>. We assume here that the <agent> has already been created. An agent is typically created by `set agent [new Agent/AgentType]` where Agent/AgentType defines the class definition of the specified agent type.

```
$agent port
```

This returns the port number to which the agent is attached.

```
$agent dst-port
```

This returns the port number of the destination. When any connection is setup between 2 nodes, each agent stores the destination port in its instance variable called `dst_port_`.

```
$agent attach-app <s_type>
```

This commands attaches an application of type <s\_type> to the agent. A handle to the application object is returned. Also note that the application type must be defined as a packet type in packet.h.

```
$agent attach-source <s_type>
```

This used to be the procedure to attach source of type <s\_type> to the agent. But this is obsolete now. Use attach-app (described above) instead.

`$agent attach-tbf <tbf>`

Attaches a token bucket filter (tbf) to the agent.

`$ns_ connect <src> <dst>`

Sets up a connection between the src and dst agents.

`$ns_ create-connection <srctype> <src> <dsttype> <dst> <pktclass>`

This sets up a complete connection between two agents. First creates a source of type <srctype> and binds it to <src>. Then creates a destination of type <dsttype> and binds it to <dst>. Finally connects the src and dst agents and returns a handle to the source agent.

`$ns_ create-connection-list <srctype> <src> <dsttype> <dst> <pktclass>`

This command is exactly similar to create-connection described above. But instead of returning only the source-agent, this returns a list of source and destination agents.

Internal procedures:

`$ns_ simplex-connect <src> <dst>`

This is an internal method that actually sets up an unidirectional connection between the <src> agent and <dst> agent. It simply sets the destination address and destination port of the <src> as <dst>'s agent-address and agent-port. The "connect" described above calls this method twice to set up a bi-directional connection between the src and dst.

`$agent set <args>`

This is an internal procedure used to inform users of the backward compatibility issues resulting from the upgrade to 32-bit addressing space currently used in *ns*.

`$agent attach-trace <file>`

This attaches the <file> to the agent to allow nam-tracing of the agent events.

In addition to the agent related procedures described here, there are additional methods that support different type of agents like Agent/Null, Agent/TCP, Agent/CBR, Agent/TORA, Agent/mcast etc. These additional methods along with the procedures described here can be found in *ns/tcl/lib/(ns-agent.tcl, ns-lib.tcl, ns-mip.tcl, ns-mobilenode.tcl, ns-namsupp.tcl, ns-queue.tcl, ns-route.tcl, ns-sat.tcl, ns-source.tcl)*. They are also described in the previous section.

# Chapter 11

## Timers

Timers may be implemented in C++ or OTcl. In C++, timers are based on an abstract base class defined in `~ns/timer-handler.h`. They are most often used in agents, but the framework is general enough to be used by other objects. The discussion below is oriented towards the use of timers in agents.

The procedures and functions described in this chapter can be found in `~ns/tcl/ex/timer.tcl`, and `~ns/timer-handler.{cc,h}`.

In OTcl, a simple timer class is defined in `~ns/tcl/ex/timer.tcl`. Subclasses can be derived to provide a simple mechanism for scheduling events at the OTcl level.

### 11.1 C++ abstract base class `TimerHandler`

The abstract base class `TimerHandler` contains the following public member functions:

<code>void sched(double delay)</code>	schedule a timer to expire <code>delay</code> seconds in the future
<code>void resched(double delay)</code>	reschedule a timer (similar to <code>sched()</code> , but timer may be pending)
<code>void cancel()</code>	cancel a pending timer
<code>int status()</code>	returns timer status (either <code>TIMER_IDLE</code> , <code>TIMER_PENDING</code> , or <code>TIMER_HANDLING</code> )

The abstract base class `TimerHandler` contains the following protected members:

<code>virtual void expire(Event*e)</code>	<code>=0</code>	this method must be filled in by the timer client
<code>virtual void handle(Event*e)</code>		consumes an event; invokes <code>expire()</code> and sets <code>status_</code> of the timer appropriately
<code>int status_</code>		keeps track of the current timer status
<code>Event event_</code>		event to be consumed upon timer expiration

The pure virtual function `expire()` must be defined by the timer classes deriving from this abstract base class.

Finally, two private inline functions are defined:

```
inline void _sched(double delay) {
```

```

        (void)Scheduler::instance().schedule(this, &event_, delay);
    }
    inline void _cancel() {
        (void)Scheduler::instance().cancel(&event_);
    }
}

```

From this code we can see that timers make use of methods of the Scheduler class.

### 11.1.1 Definition of a new timer

To define a new timer, subclass this function and define handle() if needed (handle() is not always required):

```

class MyTimer : public TimerHandler {
public:
    MyTimer(MyAgentClass *a) : TimerHandler() { a_ = a; }
    virtual double expire(Event *e);
protected:
    MyAgentClass *a_;
};

```

Then define expire:

```

double
MyTimer::expire(Event *e)
{
    // do the work
    // return TIMER_HANDLED;      // => do not reschedule timer
    // return delay;              // => reschedule timer after delay
}

```

Note that expire() can return either the flag TIMER\_HANDLED or a delay value, depending on the requirements for this timer.

Often MyTimer will be a friend of MyAgentClass, or expire() will only call a public function of MyAgentClass.

Timers are not directly accessible from the OTcl level, although users are free to establish method bindings if they so desire.

### 11.1.2 Example: Tcp retransmission timer

TCP is an example of an agent which requires timers. There are three timers defined in the basic Tahoe TCP agent defined in tcp.cc:

```

rtx_timer_;                                /* Retransmission timer */
delsnd_timer_;                             /* Delays sending of packets by a small random amount of time, */
                                           /* to avoid phase effects */
burstsnd_timer_;                           /* Helps TCP to stagger the transmission of a large window */
                                           /* into several smaller bursts */

```



In `ns/tcp.h`, three classes are derived from the base class `TimerHandler`:

```
class RtxTimer : public TimerHandler {
public:
    RtxTimer(TcpAgent *a) : TimerHandler() { a_ = a; }
protected:
    virtual void expire(Event *e);
    TcpAgent *a_;
};

class DelSndTimer : public TimerHandler {
public:
    DelSndTimer(TcpAgent *a) : TimerHandler() { a_ = a; }
protected:
    virtual void expire(Event *e);
    TcpAgent *a_;
};

class BurstSndTimer : public TimerHandler {
public:
    BurstSndTimer(TcpAgent *a) : TimerHandler() { a_ = a; }
protected:
    virtual void expire(Event *e);
    TcpAgent *a_;
};
```

In the constructor for `TcpAgent` in `tcp.cc`, each of these timers is initialized with the `this` pointer, which is assigned to the pointer `a_`.

```
TcpAgent::TcpAgent() : Agent(PT_TCP), rtt_active_(0), rtt_seq_(-1),
    ...
    rtx_timer_(this), delsnd_timer_(this), burstsnd_timer_(this)
{
    ...
}
```

In the following, we will focus only on the retransmission timer. Various helper methods may be defined to schedule timer events; e.g.,

```
/*
 * Set retransmit timer using current rtt estimate. By calling resched()
 * it does not matter whether the timer was already running.
 */
void TcpAgent::set_rtx_timer()
{
    rtx_timer_.resched(rtt_timeout());
}

/*
 * Set new retransmission timer if not all outstanding
 * data has been acked. Otherwise, if a timer is still
```

```

    * outstanding, cancel it.
    */
void TcpAgent::newtimer(Packet* pkt)
{
    hdr_tcp *tcph = (hdr_tcp*)pkt->access(off_tcp_);
    if (t_seqno_ > tcph->seqno())
        set_rtx_timer();
    else if (rtx_timer_.status() == TIMER_PENDING)
        rtx_timer_.cancel();
}

```

In the above code, the `set_rtx_timer()` method reschedules the retransmission timer by calling `rtx_timer_.resched()`. Note that if it is unclear whether or not the timer is already running, calling `resched()` eliminates the need to explicitly cancel the timer. In the second function, examples are given of the use of the `status()` and `cancel(void)` methods.

Finally, the `expire(void)` method for class `RtxTimer` must be defined. In this case, `expire(void)` calls the `timeout(void)` method for `TcpAgent`. This is possible because `timeout()` is a public member function; if it were not, then `RtxTimer` would have had to have been declared a friend class of `TcpAgent`.

```

void TcpAgent::timeout(int tno)
{
    /* retransmit timer */
    if (tno == TCP_TIMER_RTX) {
        if (highest_ack_ == maxseq_ && !slow_start_restart_) {
            /*
             * TCP option:
             * If no outstanding data, then don't do anything.
             */
            return;
        };
        recover_ = maxseq_;
        recover_cause_ = 2;
        closewnd(0);
        reset_rtx_timer(0,1);
        send_much(0, TCP_REASON_TIMEOUT, maxburst_);
    } else {
        /*
         * delayed-send timer, with random overhead
         * to avoid phase effects
         */
        send_much(1, TCP_REASON_TIMEOUT, maxburst_);
    }
}

void RxTimer::expire(Event *e) {
    a_->timeout(TCP_TIMER_RTX);
}

```

The various TCP agents contain additional examples of timers.

## 11.2 OTcl Timer class

A simple timer class is defined in `~ns/tcl/mcast/timer.tcl`. Subclasses of `Timer` can be defined as needed. Unlike the C++ timer API, where a `sched()` aborts if the timer is already set, `sched()` and `resched()` are the same; i.e., no state is kept for the OTcl timers. The following methods are defined in the `Timer` base class:

```
$self sched $delay           ; # causes "$self timeout" to be called $delay seconds in the future
$self resched $delay         ; # same as "$self sched $delay"
$self cancel                 ; # cancels any pending scheduled callback
$self destroy               ; # same as "$self cancel"
$self expire                 ; # calls "$self timeout" immediately
```

## 11.3 Commands at a glance

Following is a list of methods for the class `Timer`. Note that many different types of timers have been derived from this base class (viz. `LogTimer`, `Timer/Iface`, `Timer/Iface/Prune`, `CacheTimer`, `Timer/Scuba` etc).

```
$timer sched <delay>
```

This command cancels any other event that may have been scheduled and re-schedules another event after time `<delay>`.

```
$timer resched <delay>
```

Similar to "sched" described above. Added to have similar APIs as that of the C++ timers.

```
$timer cancel
```

This cancels any scheduled event.

```
$timer destroy
```

This is similar to `cancel`. Cancels any scheduled event.

```
$timer expire
```

This command calls for a time-out. However the time-out procedure needs to be defined in the sub-classes.

All these procedures can be found in `ns/tcl/mcast/timer.tcl`.

## Chapter 12

# Packet Headers and Formats

The procedures and functions described in this chapter can be found in `~ns/tcl/lib/ns-lib.tcl`, `~ns/tcl/lib/ns-packet.tcl`, and `~ns/packet.{cc, h}`.

Objects in the `class Packet` are the fundamental unit of exchange between objects in the simulation. The class `Packet` provides enough information to link a packet on to a list (*i.e.*, in a `PacketQueue` or on a free list of packets), refer to a buffer containing packet headers that are defined on a per-protocol basis, and to refer to a buffer of packet data. New protocols may define their own packet headers or may extend existing headers with additional fields.

New packet headers are introduced into the simulator by defining a C++ structure with the needed fields, defining a static class to provide OTcl linkage, and then modifying some of the simulator initialization code to assign a byte offset in each packet where the new header is to be located relative to others.

When the simulator is initialized through OTcl, a user may choose to enable only a subset of the compiled-in packet formats, resulting in a modest savings of memory during the execution of the simulation. Presently, most configured-in packet formats are enabled. The management of which packet formats are currently enabled in a simulation is handled by a special packet header manager object described below. This object supports an OTcl method used to specify which packet headers will be used in a simulation. If an object in the simulator makes use of a field in a header which has not been enabled, a run-time fatal program abort occurs.

### 12.1 A Protocol-Specific Packet Header

Protocol developers will often wish to provide a specific header type to be used in packets. Doing so allows a new protocol implementation to avoid overloading already-existing header fields. We consider a simplified version of RTP as an example. The RTP header will require a sequence number fields and a source identifier field. The following classes create the needed header (see `~ns/rtp.h` and `~ns/rtp.cc`):

From `rtp.h`:

```
/* rtp packet. For now, just have srcid + seqno. */
struct hdr_rtp {
    u_int32_t srcid_;
    int seqno_;
    /* per-field member functions */
    u_int32_t& srcid() { return (srcid_); }
```

```

int& seqno() { return (seqno_); }

/* Packet header access functions */
static int offset_;
inline static int& offset() { return offset_; }
inline static hdr_rtp* access(const Packet* p) {
    return (hdr_rtp*) p->access(offset_);
}
};

```

From rtp.cc:

```

class RTPHeaderClass : public PacketHeaderClass {
public:
    RTPHeaderClass() : PacketHeaderClass("PacketHeader/RTP",
                                          sizeof(hdr_rtp)) {
        bind_offset(&hdr_rtp::offset_);
    }
} class_rtphdr;

void RTPAgent::sendpkt()
{
    Packet* p = allocpkt();
    hdr_rtp *rh = hdr_rtp::access(p);
    lastpkttime_ = Scheduler::instance().clock();

    /* Fill in srcid_ and seqno_ */
    rh->seqno() = seqno_++;
    rh->srcid() = session_->srcid();
    target_->recv(p, 0);
}

RTPAgent::RTPAgent()
    : session_(0), lastpkttime_(-1e6)
{
    type_ = PT_RTP;
    bind("seqno_", &seqno_);
}

```

The first structure, `hdr_rtp`, defines the layout of the RTP packet header (in terms of words and their placement): which fields are needed and how big they are. This structure definition is only used by the compiler to compute byte offsets of fields; no objects of this structure type are ever directly allocated. The structure also provides member functions which in turn provide a layer of data hiding for objects wishing to read or modify header fields of packets. Note that the static class variable `offset_` is used to find the byte offset at which the rtp header is located in an arbitrary *ns* packet. Two methods are provided to utilize this variable to access this header in any packet: `offset()` and `access()`. The latter is what most users should choose to access this particular header in a packet; the former is used by the packet header management class and should seldom be used. For example, to access the RTP packet header in a packet pointed by `p`, one simply says `hdr_rtp::access(p)`. The actual binding of `offset_` to the position of this header in a packet is done by routines inside `~ns/tcl/lib/ns-packet.tcl` and `~ns/packet.cc`. The `const` in `access()`'s argument provides (presumably) read-only access to a `const Packet`, lthough read-only is enforced since the return pointer is not `const`. One correct way to do this is to provide two methods, one for write access, the other for read-only access. However, this is not currently implemented.

**IMPORTANT:** Notice that this is completely different from the *original* (and obsolete) method to access a packet header, which requires that an integer variable, `off_{hdrname}_`, be defined for any packet header that one needs to access. This

method is now obsolete; its usage is tricky and its misuse can be very difficult to detect.

The static object `class_rtp_hdr` of class `RTPHeaderClass` is used to provide linkage to OTcl when the RTP header is enabled at configuration time. When the simulator executes, this static object calls the `PacketHeaderClass` constructor with arguments "`PacketHeader/RTP`" and `sizeof(hdr_rtp)`. This causes the size of the RTP header to be stored and made available to the packet header manager at configuration time (see below, Section 12.2.4). Notice that `bind_offset()` **MUST** be called in the constructor of this class, so that the packet header manager knows where to store the offset for this particular packet header.

The sample member function `sendpkt()` method of `RTPAgent` creates a new packet to send by calling `allocpkt()`, which handles assignment of all the network-layer packet header fields (in this case, IP). Headers other than IP are handled separately. In this case, the agent uses the `RTPHeader` defined above. The `Packet::access(void)` member function returns the address of the first byte in a buffer used to hold header information (see below). Its return value is cast as a pointer to the header of interest, after which member functions of the `RTPHeader` object are used to access individual fields.

### 12.1.1 Adding a New Packet Header Type

Assuming we wish to create a new header called `newhdr` the following steps are performed:

1. create a new structure defining the raw fields (called `hdr_newhdr`), define `offset_` and access methods.
2. define member functions for needed fields.
3. create a static class to perform OTcl linkage (defines `PacketHeader/Newhdr`), do `bind_offset()` in its constructor.
4. edit `~ns/tcl/lib/ns-packet.tcl` to enable new packet format (see 12.2.2, 12.2.4).

*This is the recommended way to add your packet headers. If you do not follow this method, your simulation may still work, but it may behave in a unpredictable way when more protocols are added into your simulation. The reason is that the BOB (Bag of Bits, Section 12.2.1) in nspacket is a large sparse space, assigning one wrong packet header offset may not trigger failure immediately.*

### 12.1.2 Selectively Including Packet Headers in Your Simulation

By default, ns includes *ALL* packet headers of *ALL* protocols in ns in *EVERY* packet in your simulation. This is a LOT of overhead, and will increase as more protocols are added into ns. For “packet-intensive” simulations, this could be a huge overhead. For instance, as of now (Aug 30, 2000), the size of packet headers of all protocols in ns is about 1.9KB; however, if you turn on only the common header, the IP header and the TCP header, they add up to about 100 bytes. If you are doing large-scale web traffic simulation with many big fat pipes, reducing unused packet headers can lead to major memory saving.

To include only the packet headers that are of interest to you in your specific simulation, follow this pattern (e.g., you want to remove AODV and ARP headers from your simulation):

```
remove-packet-header AODV ARP
.....
set ns [new Simulator]
```

Notice that `remove-packet-header` **MUST** go before the simulator is created. All packet header names are in the forms of `PacketHeader/[hdr]`. You only need to supply the `[hdr]` part, not the prefix. To find the names of packet headers, you may either look them up in `~ns/tcl/lib/ns-packet.tcl`, or run the following simple commands in *ns*:

```

foreach cl [PacketHeader info subclass] {
    puts $cl
}

```

To include only a specific set of headers in your simulation, e.g., IP and TCP, follow this pattern:

```

remove-all-packet-headers
add-packet-header IP TCP
.....
set ns [new Simulator]

```

**IMPORTANT:** You **MUST** never remove common header from your simulation. As you can see in `~ns/tcl/lib/ns-packet.tcl`, this is enforced by these header manipulation procs.

*Notice that by default, all packet headers are included.*

## 12.2 Packet Classes

There are four C++ classes relevant to the handling of packets and packet headers in general: `Packet`, `p_info`, `PacketHeader`, and `PacketHeaderManager`. The class `Packet` defines the type for all packets in the simulation; it is a subclass of `Event` so that packets may be scheduled (e.g. for later arrival at some queue). The class `packet_info` holds all text representations for packet names. The class `PacketHeader` provides a base class for any packet header configured into the simulation. It essentially provides enough internal state to locate any particular packet header in the collection of packet headers present in any given packet. The class `PacketHeaderManager` defines a class used to collect and manage currently-configured headers. It is invoked by a method available to OTcl at simulation configuration time to enable some subset of the compiled-in packet headers.

### 12.2.1 The Packet Class

The class `Packet` defines the structure of a packet and provides member functions to handle a free list for objects of this type. It is illustrated in Figure 12.1 and defined as follows in `packet.h`:

```

class Packet : public Event {
private:
    friend class PacketQueue;
    u_char* bits_;
    u_char* data_;
    u_int datalen_;
protected:
    static Packet* free_;
public:
    Packet* next_;
    static int hdrlen_;
    Packet() : bits_(0), datalen_(0), next_(0) {}
    u_char* const bits() { return (bits_); }
    Packet* copy() const;
    static Packet* alloc();
}

```

*/\* variable size buffer for 'data' \*/*  
*/\* length of variable size buffer \*/*  
*/\* for queues and the free list \*/*

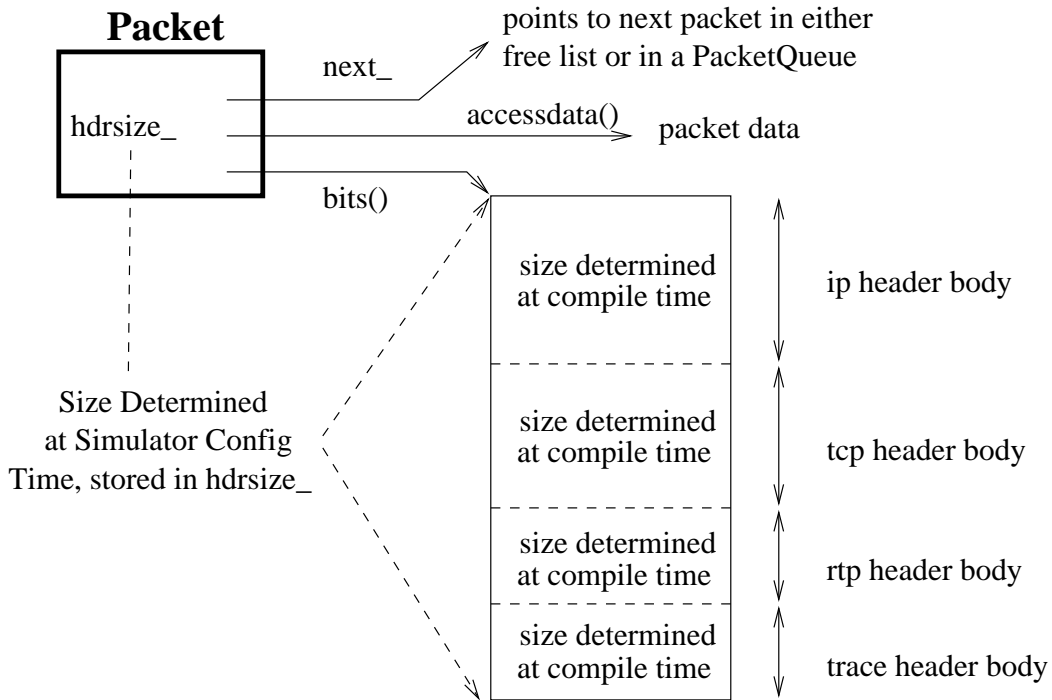


Figure 12.1: A Packet Object

```
static Packet* alloc(int);
inline void allocdata(int);
static void free(Packet*);
inline u_char* access(int off) {
    if (off < 0)
        abort();
    return (&bits_[off]);
}
inline u_char* accessdata() { return data_; }
};
```

This class holds a pointer to a generic array of unsigned characters (commonly called the “bag of bits” or BOB for short) where packet header fields are stored. It also holds a pointer to packet “data” (which is often not used in simulations). The `bits_` variable contains the address of the first byte of the BOB. Effectively BOB is (currently implemented as) a concatenation of all the structures defined for each packet header (by convention, the structures with names beginning `hdr_(something)`) that have been configured in. BOB generally remains a fixed size throughout a simulation, and the size is recorded in the `Packet::hdrlen_` member variable. This size is updated during simulator configuration by OTcl<sup>1</sup>.

The other methods of the class `Packet` are for creating new packets and storing old (unused) ones on a private free list. Such allocation and deallocation is performed by the following code (in `~ns/packet.h`):

```
inline Packet* Packet::alloc()
{
    Packet* p = free_;
    if (p != 0)
        free_ = p->next_;
```

<sup>1</sup>It is not intended to be updated after configuration time. Doing so *should* be possible, but is currently untested.



```

        else {
            p = new Packet;
            p->bits_ = new u_char[hdrsize_];
            if (p == 0 || p->bits_ == 0)
                abort();
        }
        return (p);
    }

    /* allocate a packet with an n byte data buffer */
    inline Packet* Packet::alloc(int n)
    {
        Packet* p = alloc();
        if (n > 0)
            p->alloccdata(n);
        return (p);
    }

    /* allocate an n byte data buffer to an existing packet */
    inline void Packet::alloccdata(int n)
    {
        datalen_ = n;
        data_ = new u_char[n];
        if (data_ == 0)
            abort();
    }

    inline void Packet::free(Packet* p)
    {
        p->next_ = free_;
        free_ = p;
        if (p->datalen_) {
            delete p->data_;
            p->datalen_ = 0;
        }
    }

    inline Packet* Packet::copy() const
    {
        Packet* p = alloc();
        memcpy(p->bits(), bits_, hdrsize_);
        if (datalen_) {
            p->datalen_ = datalen_;
            p->data_ = new u_char[datalen_];
            memcpy(p->data_, data_, datalen_);
        }
        return (p);
    }
}

```

The `alloc()` method is a support function commonly used to create new packets. It is called by `Agent::allocpkt()` method on behalf of agents and is thus not normally invoked directly by most objects. It first attempts to locate an old packet on the free list and if this fails allocates a new one using the C++ new operator. Note that `Packet` class objects and BOBs are allocated separately. The `free()` method frees a packet by returning it to the free list. Note that *packets are never returned to*

the system's memory allocator. Instead, they are stored on a free list when `Packet::free()` is called. The `copy()` member creates a new, identical copy of a packet with the exception of the `uid_` field, which is unique. This function is used by Replicator objects to support multicast distribution and LANs.

### 12.2.2 p\_info Class

This class is used as a “glue” to bind numeric packet type values with their symbolic names. When a new packet type is defined, its numeric code should be added to the enumeration `packet_t` (see `~ns/packet.h`)<sup>2</sup> and its symbolic name should be added to the constructor of `p_info`:

```
enum packet_t {
    PT_TCP,
    ...
    PT_NTTYPE // This MUST be the LAST one
};

class p_info {
public:
    p_info() {
        name_[PT_TCP]= "tcp";
        ...
    }
}
```

### 12.2.3 The hdr\_cmh Class

Every packet in the simulator has a “common” header which is defined in `~ns/packet.h` as follows:

```
struct hdr_cmh {
    double    ts_;                /* timestamp: for q-delay measurement */
    packet_t  ptype_;            /* packet type (see above) */
    int       uid_;              /* unique id */
    int       size_;             /* simulated packet size */
    int       iface_;            /* receiving interface (label) */

    /* Packet header access functions */
    static int offset_;
    inline static int& offset() { return offset_; }
    inline static hdr_cmh* access(Packet* p) {
        return (hdr_cmh*) p->access(offset_);
    }

    /* per-field member functions */
    int& ptype() { return (ptype_); }
    int& uid() { return (uid_); }
    int& size() { return (size_); }
    int& iface() { return (iface_); }
```

---

<sup>2</sup>Note: `PT_NTTYPE` should remain the last element of this enumeration.

```

        double& timestamp() { return (ts_); }
};

```

This structure primarily defines fields used for tracing the flow of packets or measuring other quantities. The `time_stamp` field is used to measure queuing delay at switch nodes. The `ptype_` field is used to identify the type of packets, which makes reading traces simpler. The `uid_` field is used by the scheduler in scheduling packet arrivals. The `size_` field is of general use and gives the simulated packet's size in bytes. Note that the actual number of bytes consumed in the simulation may not relate to the value of this field (i.e., `size_` has *no* relationship to `sizeof(struct_hdr_cmh)` or other ns structures). Rather, it is used most often in computing the time required for a packet to be delivered along a network link. As such it should be set to the sum of the application data size and IP-, transport-, and application-level headers for the simulated packet. The `iface_` field is used by the simulator when performing multicast distribution tree computations. It is a label indicating (typically) on which link a packet was received.

## 12.2.4 The PacketHeaderManager Class

An object of the class `PacketHeaderManager` is used to manage the set of currently-active packet header types and assign each of them unique offsets in the BOB. It is defined in both the C++ and OTcl code:

From `tcl/lib/ns-packet.tcl`:

```

PacketHeaderManager set hdrlen_ 0
.....
foreach prot {
    AODV
    ARP
    aSRM
    Common
    CtrMcast
    Diffusion
    .....
    TORA
    UMP
} {
    add-packet-header $prot
}
Simulator instproc create_packetformat {} {
    PacketHeaderManager instvar tab_
    set pm [new PacketHeaderManager]
    foreach cl [PacketHeader info subclass] {
        if [info exists tab_($cl)] {
            set off [$pm allochdr $cl]
            $cl offset $off
        }
    }
    $self set packetManager_ $pm
}
PacketHeaderManager instproc allochdr cl {
    set size [$cl set hdrlen_]
    $self instvar hdrlen_
    set NS_ALIGN 8 ;# round up to nearest NS_ALIGN bytes, (needed on sparc/solaris)
    set incr [expr ($size + (NS_ALIGN-1)) & ~(NS_ALIGN-1)]
    set base $hdrlen_

```

```

        incr hdrlen_ $incr
        return $base
    }

```

From packet.cc:

```

/* manages active packet header types */
class PacketHeaderManager : public TclObject {
public:
    PacketHeaderManager() {
        bind("hdrlen_", &Packet::hdrlen_);
    }
};

```

The code in `~ns/tcl/lib/ns-packet.tcl` is executed when the simulator initializes. Thus, the `foreach` statement is executed before the simulation begins, and initializes the OTcl class array `tab_` to contain the mapping between class the name and the names of the currently active packet header classes. As discussed above (12.1), packet headers should be accessed using `hdr_<hdrname>::access()`.

The `create_packetformat{}` instance procedure is part of the basic Simulator class and is called one time during simulator configuration. It first creates a single `PacketHeaderManager` object. The C++ constructor links the OTcl instance variable `hdrlen_` (of class `PacketHeaderManager`) to the C++ variable `Packet::hdrlen_` (a static member of the `Packet` class). This has the effect of setting `Packet::hdrlen_` to zero. Note that binding across class types in this fashion is unusual.

After creating the packet manager, the `foreach` loop enables each of the packet headers of interest. This loop iterates through the list of defined packet headers of the form  $(h_i, o_i)$  where  $h_i$  is the name of the  $i$ th header and  $o_i$  is the name of the variable containing the location of the  $h_i$  header in BOB. The placement of headers is performed by the `allochdr` instproc of the `PacketHeaderManager` OTcl class. The procedure keeps a running variable `hdrlen_` with the current length of BOB as new packet headers are enabled. It also arranges for 8-byte alignment for any newly-enabled packet header. This is needed to ensure that when double-world length quantities are used in packet headers on machines where double-word alignment is required, access faults are not produced.<sup>3</sup>

## 12.3 Commands at a glance

Following is a list of packet-header related procedures:

`Simulator::create_packetformat`

This is an internal simulator procedure and is called once during the simulator configuration to setup a `PacketHeaderManager` object.

`PacketHeaderManager::allochdr`

This is another internal procedure of Class `PacketHeaderManager` that keeps track of a variable called `hdrlen_` as new packet-headers are enabled. It also allows 8-byte alignment for any newly-enabled pkt header.

`add-packet-header` takes a list of arguments, each of which is a packet header name (without `PacketHeader/` prefix). This global proc will tell simulator to include the specified packet header(s) in your simulation.

---

<sup>3</sup>In some processor architectures, including the Sparc and HP-PA, double-word access must be performed on a double-word boundary (i.e. addresses ending in 0 mod 8). Attempting to perform unaligned accesses result in an abnormal program termination.

`remove-packet-header` operates in the same syntax, but it removes the specified headers from your simulation; notice that it does not remove the common header even it is instructed to do so.

`remove-all-packet-headers` is a global Tcl proc. It takes no argument and removes all packet headers, except the common header, from your simulation. `add-all-packet-headers` is its counterpart.

# Chapter 13

## Error Model

This chapter describes the implementation and configuration of error models, which introduces packet losses into a simulation.

In addition to the basic class `ErrorModel` described in details below, there are several other types of error modules not being completely documented yet, which include:

- `SRModelErrorModel`, `PGModelErrorModel`: error model for SRM and PGM.
- `ErrorModel/Trace`: error model that reads a loss trace (instead of a math/computed model)
- `MrouteErrorModel`: error model for multicast routing, now inherits from trace.
- `ErrorModel/Periodic`: models periodic packet drops (drop every *nth* packet we see). This model can be conveniently combined with a flow-based classifier to achieve drops in particular flows
- `SelectErrorModel`: for Selective packet drop.
- `ErrorModel/TwoState`: Two-State: error-free and error
- `ErrorModel/TwoStateMarkov`, `ErrorModel/Expo`, `ErrorModel/Empirical`: inherit from `ErrorModel/TwoState`.
- `ErrorModel/List`: specify a list of packets/bytes to drop, which could be in any order.

Their definitions can be found in `~ns/queue/errmodel.{cc,h}` and `~ns/tcl/lib/ns-errmodel.tcl`, `ns-default.tcl`.

### 13.1 Implementation

The procedures and functions described in this section can be found in `~ns/errmodel.{cc,h}`.

Error model simulates link-level errors or loss by either marking the packet's error flag or dumping the packet to a drop target. In simulations, errors can be generated from a simple model such as the packet error rate, or from more complicated statistical and empirical models. To support a wide variety of models, the unit of error can be specified in term of packet, bits, or time-based.

The `ErrorModel` class is derived from the `Connector` base class. As the result, it inherits some methods for hooking up objects such as `target` and `drop-target`. If the drop target exists, it will received corrupted packets from `ErrorModel`.

Otherwise, `ErrorModel` just marks the `error_` flag of the packet's common header, thereby, allowing agents to handle the loss. The `ErrorModel` also defines additional Tcl method `unit` to specify the unit of error and `ranvar` to specify the random variable for generating errors. If not specified, the unit of error will be in packets, and the random variable will be uniform distributed from 0 to 1. Below is a simple example of creating an error model with the packet error rate of 1 percent (0.01):

```
# create a loss_module and set its packet error rate to 1 percent
set loss_module [new ErrorModel]
$loss_module set rate_ 0.01

# optional: set the unit and random variable
$loss_module unit pkt ;# error unit: packets (the default)
$loss_module ranvar [new RandomVariable/Uniform]

# set target for dropped packets
$loss_module drop-target [new Agent/Null]
```

In C++, the `ErrorModel` contains both the mechanism and policy for dropping packets. The packet dropping mechanism is handled by the `recv` method, and packet corrupting policy is handled by the `corrupt` method.

```
enum ErrorUnit { EU_PKT=0, EU_BIT, EU_TIME };

class ErrorModel : public Connector {
public:
    ErrorModel();
    void recv(Packet*, Handler*);
    virtual int corrupt(Packet*);
    inline double rate() { return rate_; }
protected:
    int command(int argc, const char*const* argv);
    ErrorUnit eu_; /* error unit in pkt, bit, or time */
    RandomVariable* ranvar_;
    double rate_;
};
```

The `ErrorModel` only implements a simple policy based on a single error rate, either in packets of bits. More sophisticated dropping policy can be implemented in C++ by deriving from `ErrorModel` and redefining its `corrupt` method.

## 13.2 Configuration

The previous section talked about error model, in this section we discuss how to use error models in ns over either wired networks or wireless ones.

To use an error model for wired networks, at first it has to be inserted into a `SimpleLink` object. Because a `SimpleLink` is a composite object (Chapter 6), an error model can be inserted to many places. Currently we provide the following methods to insert an error module into three different places.

- Insert an error module in a `SimpleLink` BEFORE the queue module. This is provided by the following two OTcl methods:

`SimpleLink::errormodule` args When an error model is given as a parameter, it inserts the error module into the simple link, right after the queue module, and set the drop-target of the error model to be the drop trace object of the simple link. Note that this requires the following configuration order: `ns namtrace-all` followed by link configurations, followed by error model insertion. When no argument is given, it returns the current error model in the link, if there's any. This method is defined in `ns/tcl/lib/ns-link.tcl`

`Simulator::lossmodel` `<em>` `<src>` `<dst>` Call `SimpleLink::errormodule` to insert the given error module into the simple link (`src`, `dst`). It's simply a wrapper for the above method. This method is defined in `ns/tcl/lib/ns-lib.tcl`.

- Insert an error module in a SimpleLink AFTER the queue but BEFORE the delay link. This is provided by the following two methods:

`SimpleLink::insert-linkloss` args This method's behavior is identical to that of `SimpleLink::errormodule`, except that it inserts an error module immediately after the queue object. It's defined in `ns/tcl/lib/ns-link.tcl`

`Simulator::link-lossmodel` `<em>` `<src>` `<dst>` This is a wrapper for `SimpleLink::insert-linkloss`. It's defined in `ns/tcl/lib/ns-lib.tcl`

The nam traces generated by error models inserted using these two methods do not require special treatment and can be visualized using an older version of nam.

- Insert an error module in a Link AFTER the delay link module. This can be done by `Link::install-error`. Currently this API doesn't produce any trace. It only serves as a placeholder for possible future extensions.

To add an error model over wireless networks, each node can insert a given statistical error model either over outgoing or incoming wireless channels. Precisely, the instantiated error model would be stuck between mac and netif modules depicted in Figure 16.2. For the outgoing link, the error module would be pointed by `downtarget_` of the above mac module while for the incoming link it would be linked by `uptaget_` pointer of the below netif module. And in each case the `target_` of the error module points to either the netif or the mac respectively. The difference of placing over the two different locations is that the outgoing causes all the receivers to receive the packet suffering the same degree of errors since the error is determined before wireless channel module copies the packet. On the other hand, the incoming error module lets each receiver get the packet corrupted with different degree of error since the error is independently computed in each error module.

The insertion into the wireless protocol stack can be done by calling `node-config` command explained in Section ?? with the two options `IncomingErrProc` and `OutgoingErrProc`. We can use these two options at the same time or each one separately. The argument of the two option is the name of the global procedure which creates an error model object, feeds it with necessary initial values appropriate for the new error module, and finally returns the object pointer. The following shows a simple TCL example script to add an error module into the wireless protocol stack.

```
$ns node-config -IncomingErrProc UniformErr -OutgoingErrProc UniformErr

proc UniformErr
set err [new ErrorModel]
$err unit packet
return $err
```

### 13.3 Multi-state error model

Contributed by Jianping Pan (jpan@bbr.uwaterloo.ca).



The multi-state error model implements time-based error state transitions. Transitions to the next error state occur at the end of the duration of the current state. The next error state is then selected using the transition state matrix.

To create a multi-state error model, the following parameters should be supplied (as defined in `ns/tcl/lib/ns-errmodel.tcl`):

- `states`: an array of states (error models).
- `periods`: an array of state durations.
- `trans`: the transition state model matrix.
- `transunit`: one of `[pkt|byte|time]`.
- `sttype`: type of state transitions to use: either `time` or `pkt`.
- `nstates`: number of states.
- `start`: the start state.

Here is a simple example script to create a multi-state error model:

```
set tmp [new ErrorModel/Uniform 0 pkt]
set tmp1 [new ErrorModel/Uniform .9 pkt]
set tmp2 [new ErrorModel/Uniform .5 pkt]

# Array of states (error models)
set m_states [list $tmp $tmp1 $tmp2]
# Durations for each of the states, tmp, tmp1 and tmp2, respectively
set m_periods [list 0 .0075 .00375]
# Transition state model matrix
set m_transmx { {0.95 0.05 0}
                {0 0 1}
                {1 0 0} }
set m_trunit pkt
# Use time-based transition
set m_sttype time
set m_nstates 3
set m_nstart [lindex $m_states 0]

set em [new ErrorModel/MultiState $m_states $m_periods $m_transmx
      $m_trunit $m_sttype $m_nstates $m_nstart]
```

## 13.4 Commands at a glance

The following is a list of error-model related commands commonly used in simulation scripts:

```
set em [new ErrorModel]
$em unit pkt
$em set rate_ 0.02
$em ranvar [new RandomVariable/Uniform]
$em drop-target [new Agent/Null]
```

This is a simple example of how to create and configure an error model. The commands to place the error-model in a simple link will be shown next.

```
$simplelink errormodule <args>
```

This command inserts the error-model before the queue object in simple link. However in this case the error-model's drop-target points to the link's drophead\_ element.

```
$ns_ lossmodel <em> <src> <dst>
```

This command places the error-model before the queue in a simplelink defined by the <src> and <dst> nodes. This is basically a wrapper for the above method.

```
$simplelink insert-linkloss <args>
```

This inserts a loss-module after the queue, but right before the delay link\_ element in the simple link. This is because nam can visualize a packet drop only if the packet is on the link or in the queue. The error-module's drop-target points to the link's drophead\_ element.

```
$ns_ link-lossmodel <em> <src> <dst>
```

This too is a wrapper method for insert-linkloss method described above. That is this inserts the error-module right after the queue element in a simple link (src-dst).

## Chapter 14

# Local Area Networks

The characteristics of the wireless and local area networks (LAN) are inherently different from those of point-to-point links. A network consisting of multiple point-to-point links cannot capture the sharing and contention properties of a LAN. To simulate these properties, we created a new type of a Node, called `LanNode`. The OTcl configurations and interfaces for `LanNode` reside in the following two files in the main *ns* directory:

```
tcl/lan/vlan.tcl
tcl/lan/ns-ll.tcl
tcl/lan/ns-mac.tcl
```

### 14.1 Tcl configuration

The interface for creating and configuring a LAN slightly differs from those of point-to-point link. At the top level, the OTcl class `Simulator` exports a new method called `make-lan`. The parameters to this method are similar to the method `duplex-link`, except that `make-lan` only accepts a list of nodes as a single parameter instead of 2 parameters as in `duplex-link`:

```
Simulator instproc make-lan {nodes bw delay lltype ifqtype mactype chantype}
```

The optional parameters to `make-lan` specify the type of objects to be created for the link layer (LL), the interface queue, the MAC layer (Mac), and the physical layer (Channel). Below is an example of how a new CSMA/CD (Ethernet) LAN is created.

Example:

```
$ns make-lan "$n1 $n2" $bw $delay LL Queue/DropTail Mac/Csma/Cd
```

creates a LAN with basic link-layer, drop-tail queue, and CSMA/CD MAC.

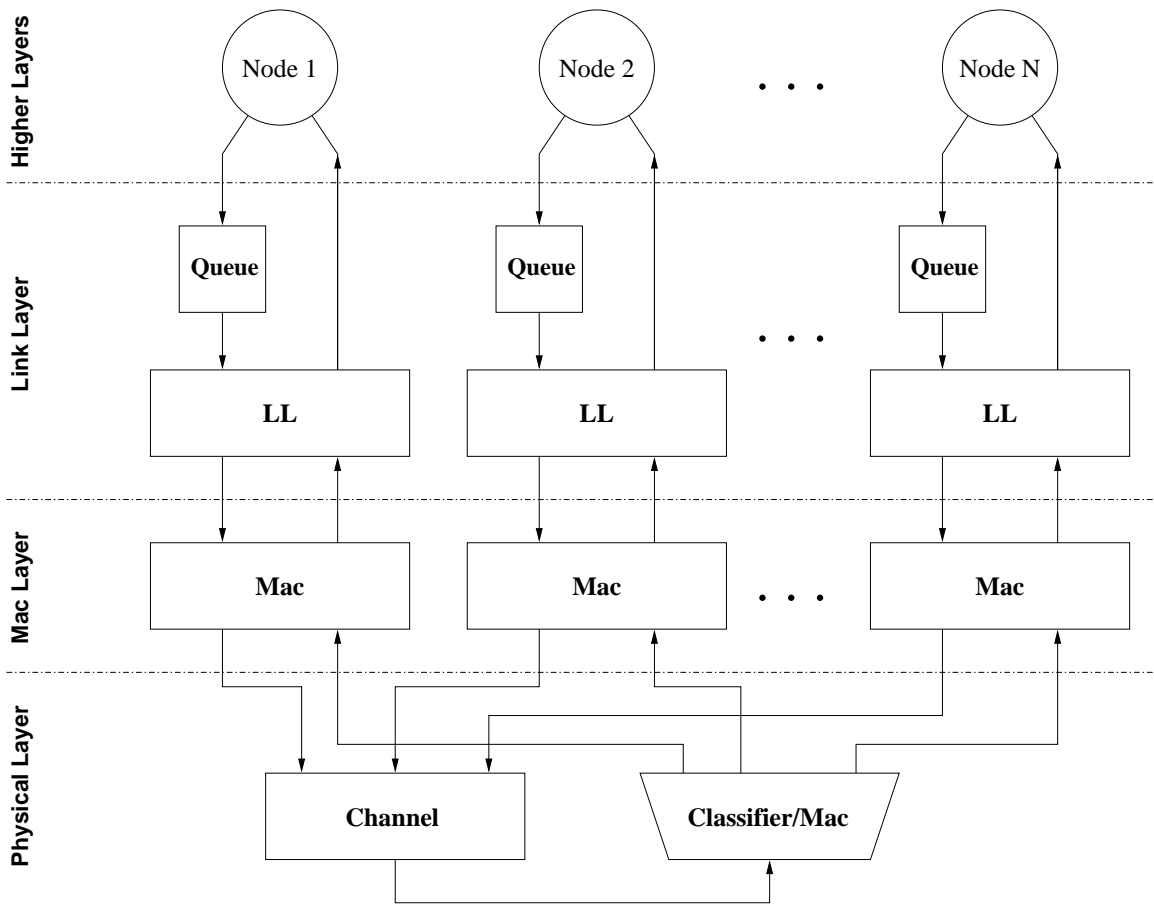


Figure 14.1: Connectivity within a LAN

## 14.2 Components of a LAN

LanLink captures the functionality of the three lowest layers in the network stack:

1. Link Layer (LL)
2. Medium Access Control (MAC) Layer
3. Physical (PHY) Layer

Figure 14.1 illustrates the extended network stack that makes simulations of local area network possible in *ns*. A packet sent down the stack flows through the link layer (Queue and LL), the MAC layer (Mac), and the physical layer (Channel to Classifier/Mac). The packet then makes its way up the stack through the Mac, and the LL.

At the bottom of the stack, the physical layer is composed of two simulation objects: the Channel and Classifier/Mac. The Channel object simulates the shared medium and supports the medium access mechanisms of the MAC objects on the sending side of the transmission. On the receiving side, the Classifier/Mac is responsible for delivering and optionally replicating packets to the receiving MAC objects.

Depending on the type of physical layer, the MAC layer must contain a certain set of functionalities such as: carrier sense, collision detection, collision avoidance, etc. Since these functionalities affect both the sending and receiving sides, they are

implemented in a single `Mac` object. For sending, the `Mac` object must follow a certain medium access protocol before transmitting the packet on the channel. For receiving, the `MAC` layer is responsible for delivering the packet to the link layer.

Above the `MAC` layer, the link layer can potentially have many functionalities such as queuing and link-level retransmission. The need of having a wide variety of link-level schemes leads to the division of functionality into two components: `Queue` and `LL` (link-layer). The `Queue` object, simulating the interface queue, belongs to the same `Queue` class that is described in Chapter 7. The `LL` object implements a particular data link protocol, such as ARQ. By combining both the sending and receiving functionalities into one module, the `LL` object can also support other mechanisms such as piggybacking.

## 14.3 Channel Class

The `Channel` class simulates the actual transmission of the packet at the physical layer. The basic `Channel` implements a shared medium with support for contention mechanisms. It allows the `MAC` to carry out carrier sense, contention, and collision detection. If more than one transmissions overlaps in time, a channel raises the collision flag. By checking this flag, the `MAC` object can implement collision detection and handling.

Since the transmission time is a function of the number of bits in the packet and the modulation speed of each individual interface (`MAC`), the `Channel` object only sets its busy signal for the duration requested by the `MAC` object. It also schedules the packets to be delivered to the destination `MAC` objects after the transmission time plus the propagation delay.

### 14.3.1 Channel State

The C++ class `Channel` includes enough internal state to schedule packet delivery and detect collisions. It exports the following OTcl configuration parameter:

```
delay_    propagation delay on the channel
```

### 14.3.2 Example: Channel and classifier of the physical layer

```
set channel_ [new Channel]
$channel_ set delay_ 4us          # propagation delay

set mcl_ [new Classifier/Mac]
$channel_ target $mcl_
$mcl_ install $mac_DA $recv_iface
. . .
```

### 14.3.3 Channel Class in C++

In C++, the class `Channel` extends the `Connector` object with several new methods to support a variety of `MAC` protocols. The class is defined as follows in `~ns/channel.h`:

```
class Channel : public Connector {
public:
    Channel();
```

```

        void recv(Packet* p, Handler*);
        virtual int send(Packet* p, double txttime);
        virtual void contention(Packet*, Handler*);
        int hold(double txttime);
        virtual int collision() { return numtx_ > 1; }
        virtual double txstop() { return txstop_; }
        . . .
};

```

The important methods of the class Channel are:

- txstop() method returns the time when the channel will become idle, which can be used by the MAC to implement carrier sense.
- contention() method allows the MAC to contend for the channel before sending a packet. The channel then use this packet to signal the corresponding Mac object at the end of each contention period.
- collision() method indicates whether a collision occurs during the contention period. When the Channel signal the end of the contention period, the MAC can use the collision() method to detect collision.
- send() method allows the MAC object to transmit a packet on the channel for a specified duration of time.
- hold() method allows the MAC object to hold the channel for a specified duration of time without actually transmitting any packets. This is useful in simulating the jamming mechanism of some MAC protocols.

## 14.4 MacClassifier Class

The MacClassifier class extends the Classifier class to implement a simple broadcasting mechanism. It modifies the recv() method in the following way: since the replication of a packet is expensive, normally a unicast packet will be classified by the MAC destination address macDA\_ and delivered directly to the MAC object with such an address. However, if the destination object cannot be found or if the MAC destination address is explicitly set to the broadcast address BCAST\_ADDR, the packet will be replicated and sent to all MACs on the lan excluding the one that is the source of the packet. Finally, by setting the bound variable MacClassifier::bcast\_ to a non-zero value, will cause MacClassifier always to replicate packets.

```

class MacClassifier : public Classifier {
public:
    void recv(Packet*, Handler*);
};

void MacClassifier::recv(Packet* p, Handler*)
{
    Mac* mac;
    hdr_mac* mh = hdr_mac::access(p);

    if (bcast_ || mh->macDA() == BCAST_ADDR || (mac = (Mac *)find(p)) == 0) {
        // Replicate packets to all slots (broadcast)
        . . .
        return;
    }
    mac->recv(p);
}

```

## 14.5 MAC Class

The `Mac` object simulates the medium access protocols that are necessary in the shared medium environment such as the wireless and local area networks. Since the sending and receiving mechanisms are tightly coupled in most types of MAC layers, it is essential for the `Mac` object to be duplex.

On the sending side, the `Mac` object is responsible for adding the MAC header and transmitting the packet onto the channel. On the receiving side, the `Mac` object asynchronously receives packets from the classifier of the physical layer. After MAC protocol processing, it passes the data packet to the link layer.

### 14.5.1 Mac State

The C++ class `Mac` contains enough internal state to simulate the particular MAC protocol. It also exports the following OTcl configuration parameter:

<code>bandwidth_</code>	modulation rate of the MAC
<code>hlen_</code>	additional bytes added to packet for MAC header
<code>label_</code>	MAC address

### 14.5.2 Mac Methods

The class `Mac` added several Tcl methods for configuration, in particular, linking with other simulation objects:

<code>channel</code>	specify the channel for transmission
<code>classifier</code>	the classifier that deliver packets to receiving MAC
<code>maclist</code>	a link list of MAC interfaces on the same node

### 14.5.3 Mac Class in C++

In C++, the `Mac` class derives from `Connector`. When the `recv()` method gets a packet, it identifies the direction of the packet based on the presence of a callback handler. If there is a callback handler, the packet is outgoing, otherwise, it is incoming.

```
class Mac : public Connector {
public:
    Mac();
    virtual void recv(Packet* p, Handler* h);
    virtual void send(Packet* p);
    virtual void resume(Packet* p = 0);
    . . .
};
```

When a `Mac` object receives a packet via its `recv()` method, it checks whether the packet is outgoing or incoming. For an outgoing packet, it assumes that the link-layer of the sender has obtained the destination MAC address and filled in the `macDA_` field of the MAC header, `hdr_mac`. The `Mac` object fills in the rest of the MAC header with the source MAC address

and the frame type. It then passes the packet to its `send()` method, which carries out the medium access protocol. For the basic `Mac` object, the `send()` method calls `txtime()` to compute the transmission time, then invokes `Channel::send` to transmit the packet. Finally, it schedules itself to resume after the transmission time has elapsed.

For an incoming packet, the `MAC` object does its protocol processing and passes the packet to the link-layer.

### 14.5.4 CSMA-based MAC

The class `CsmaMac` extends the `Mac` class with new methods that implements carrier sense and backoff mechanisms. The `CsmaMac::send()` method detects when the channel becomes idle using `Channel::txtime()`. If the channel is busy, the `MAC` schedules the next carrier sense at the moment the channel turns idle. Once the channel is idle, the `CsmaMac` object initiates the contention period with `Channel::contention()`. At the end of the contention period, the `endofContention()` method is invoked. At this time, the basic `CsmaMac` just transmits the packet using `Channel::send`.

```
class CsmaMac : public Mac {
public:
    CsmaMac();
    void send(Packet* p);
    void resume(Packet* p = 0);
    virtual void endofContention(Packet* p);
    virtual void backoff(Handler* h, Packet* p, double delay=0);
    . . .
};

class CsmaCdMac : public CsmaMac {
public:
    CsmaCdMac();
    void endofContention(Packet*);
};

class CsmaCaMac : public CsmaMac {
public:
    CsmaCaMac();
    virtual void send(Packet*);
};
```

The `CsmaCdMac` extends `CsmaMac` to carry out collision detection procedure of the CSMA/CD (Ethernet) protocol. When the channel signals the end of contention period, the `endofContention` method checks for collision using the `Channel::collision` method. If there is a collision, the `MAC` invokes its `backoff` method to schedule the next carrier sense to retransmit the packet.

The `CsmaCaMac` extends the `send` method of `CsmaMac` to carry out the collision avoidance (CSMA/CA) procedure. Instead of transmitting immediately when the channel is idle, the `CsmaCaMac` object backs off a random number of slots, then transmits if the channel remains idle until the end of the backoff period.



## 14.6 LL (link-layer) Class

The link-layer object is responsible for simulating the data link protocols. Many protocols can be implemented within this layer such as packet fragmentation and reassembly, and reliable link protocol.

Another important function of the link layer is setting the MAC destination address in the MAC header of the packet. In the current implementation this task involves two separate issues: finding the next-hop-node's IP address (routing) and resolving this IP address into the correct MAC address (ARP). For simplicity, the default mapping between MAC and IP addresses is one-to-one, which means that IP addresses are re-used at the MAC layer.

### 14.6.1 LL Class in C++

The C++ class LL derives from the LinkDelay class. Since it is a duplex object, it keeps a separate pointer for the send target, sendtarget, and the receive target, recvtarget. It also defines the methods recvfrom() and sendto() to handle the incoming and outgoing packets respectively.

```
class LL : public LinkDelay {
public:
    LL();
    virtual void recv(Packet* p, Handler* h);
    virtual Packet* sendto(Packet* p, Handler* h = 0);
    virtual Packet* recvfrom(Packet* p);

    inline int seqno() return seqno_;
    inline int ackno() return ackno_;
    inline int macDA() return macDA_;
    inline Queue* ifq() return ifq_;
    inline NsObject* sendtarget() return sendtarget_;
    inline NsObject* recvtarget() return recvtarget_;

protected:
    int command(int argc, const char*const* argv);
    void handle(Event* e) recv((Packet*)e, 0);
    inline virtual int arp(int ip_addr) return ip_addr;
    int seqno_; // link-layer sequence number
    int ackno_; // ACK received so far
    int macDA_; // destination MAC address
    Queue* ifq_; // interface queue
    NsObject* sendtarget_; // for outgoing packet
    NsObject* recvtarget_; // for incoming packet

    LanRouter* lanrouter_; // for lookups of the next hop
};
```

### 14.6.2 Example: Link Layer configuration

```
set ll_ [new LL]
set ifq_ [new Queue/DropTail]
$ll_ lanrouter [new LanRouter $ns $lan] # LanRouter is one object
```

```

                                # per LAN
$ll_ set delay_ $delay          # link-level overhead
$ll_ set bandwidth_ $bw         # bandwidth
$ll_ sendtarget $mac            # interface queue at the sender side
$ll_ recvtarget $iif            # input interface of the receiver
. . .

```

## 14.7 LanRouter class

By default, there is just one `LanRouter` object per LAN, which is created when a new `LanNode` is initialized. For every node on the LAN, the link layer object (LL) has a pointer to the `LanRouter`, so it is able to find the next hop for the packet that is sent on the LAN:

```

Packet* LL::sendto(Packet* p, Handler* h)
{
    int nh = (lanrouter_) ? lanrouter_>next_hop(p) : -1;
    . . .
}

```

`LanRouter` is able to find the next hop by querying the current `RouteLogic`:

```

int LanRouter::next_hop(Packet *p) {
    int next_hopIP;
    if (enableHrouting_) {
        routelogic_>lookup_hier(lanaddr_, adst, next_hopIP);
    } else {
        routelogic_>lookup_flat(lanaddr_, adst, next_hopIP);
    }
}

```

One limitation of this is that `RouteLogic` may not be aware of dynamic changes to the routing. But it is always possible to derive a new class from `LanRouter` so that to re-define its `next_hop` method to handle dynamic changes appropriately.

## 14.8 Other Components

In addition to the C++ components described above, simulating local area networks also requires a number of existing components in *ns* such as `Classifier`, `Queue`, and `Trace`, `networkinterface`, etc. Configuring these objects requires knowledge of what the user wants to simulate. The default configuration is implemented in the two Tcl files mentioned at the beginning of this chapter. To obtain more realistic simulations of wireless networks, one can use the `ErrorModel` described in Chapter 13.

## 14.9 LANs and *ns* routing

When a LAN is created using either `make-lan` or `newLan`, a “virtual LAN node” `LanNode` is created. `LanNode` keeps together all shared objects on the LAN: `Channel`, `Classifier/Mac`, and `LanRouter`. Then for each node on the LAN,

a `LanIface` object is created. `LanIface` contains all other objects that are needed on the per-node basis: a `Queue`, a link layer (LL), `Mac`, etc. It should be emphasized that `LanNode` is a node only for routing algorithms: `Node` and `LanNode` have very little in common. One of few things that they share is an identifier taken from the `Node ID`-space. If *hierarchical routing* is used, `LanNode` *has to be assigned a hierarchical address* just like any other node. From the point of view of *ns* (static) routing, `LanNode` is just another node connected to every node on the LAN. Links connecting the `LanNode` with

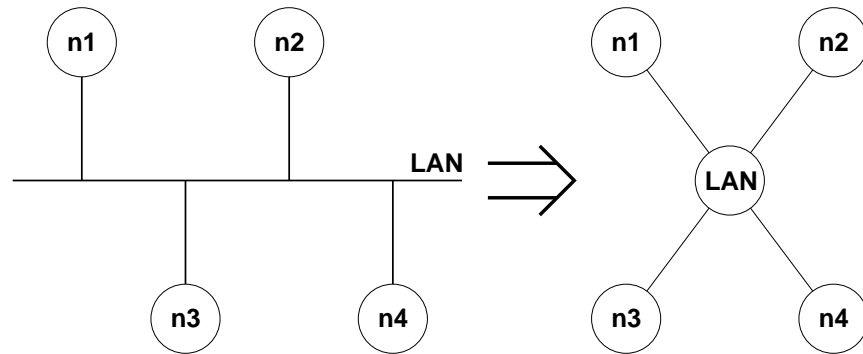


Figure 14.2: Actual LAN configuration (left) and as seen by *ns* routing (right)

the nodes on the LAN are also “virtual” (`Vlink`). The default routing cost of such a link is  $1/2$ , so the cost of traversing two `Vlinks` (e.g. `n1`  $\rightarrow$  `LAN`  $\rightarrow$  `n2`) is counted as just one hop.

Most important method of `Vlink` is the one that gives the head of the link:

```
Vlink instproc head {} {
    $self instvar lan_ dst_ src_
    if {$src_ == [$lan_ set id_]} {
        # if this is a link FROM the lan vnode,
        # it doesn't matter what we return, because
        # it's only used by $lan add-route (empty)
        return ""
    } else {
        # if this is a link TO the lan vnode,
        # return the entry to the lanIface object
        set src_lif [$lan_ set lanIface_($src_)]
        return [$src_lif entry]
    }
}
```

This method is used by static (default) routing to install correct routes at a node (see `Simulator` methods `compute-flat-routes` and `compute-hier-routes` in `tcl/lib/ns-route.tcl`, as well as `Node` methods `add-route` and `add-hroute` in `tcl/lib/ns-node.tcl`).

From the code fragment above it can be seen that it returns `LAN` interface of the node as a head of the link to be installed in the appropriate classifier.

Thus, `Vlink` *does not impose any delay on the packet* and serves the only purpose to install `LAN` interfaces instead of normal links at nodes' classifiers.

Note, that this design allows to have nodes connected by parallel `LANs`, while in the current implementation it is impossible to have nodes connected by parallel simple links and use them both (the array `Simulator instvar link_` holds the link object for each connected pair of source and destination, and it can be only one object per source/destination pair).

## 14.10 Commands at a glance

The following is a list of lan related commands commonly used in simulation scripts:

```
$ns_ make-lan <nodelist> <bw> <delay> <LL> <ifq> <MAC> <channel> <phy>
```

Creates a lan from a set of nodes given by <nodelist>. Bandwidth, delay characteristics along with the link-layer, Interface queue, Mac layer and channel type for the lan also needs to be defined. Default values used are as follows:

```
<LL> .. LL
```

```
<ifq>.. Queue/DropTail
```

```
<MAC>.. Mac
```

```
<channel>.. Channel and
```

```
<phy>.. Phy/WiredPhy
```

```
$ns_ newLan <nodelist> <BW> <delay> <args>
```

This command creates a lan similar to make-lan described above. But this command can be used for finer control whereas make-lan is a more convenient and easier command. For example newLan may be used to create a lan with hierarchical addresses. See *ns/tcl/ex/vlantest-hier.tcl*, *vlantest-mcst.tcl*, *lantest.tcl*, *mac-test.tcl* for usage of newLan. The possible argument types that can be passed are LL, ifq, MAC, channel, phy and address.

```
$lannode cost <c>
```

This assigns a cost of  $c/2$  to each of the (uni-directional) links in the lan.

```
$lannode cost?
```

Returns the cost of (bi-directional) links in the lan, i.e  $c$ .

Internal procedures :

```
$lannode addNode <nodes> <bw> <delay> <LL> <ifq> <MAC> <phy>
```

Lan is implemented as a virtual node. The LanNode mimics a real node and uses an address (id) from node's address space. This command adds a list of <nodes> to the lan represented by lannode. The bandwidth, delay and network characteristics of nodes are given by the above arguments. This is an internal command used by make-lan and newLan.

```
$lannode id
```

Returns the virtual node's id.

```
$lannode node-addr
```

Returns virtual nodes's address.

```
$lannode dump-namconfig
```

This command creates a given lan layout in nam. This function may be changed to redefine the lan layout in a different way.

```
$lannode is-lan?
```

This command always returns 1, since the node here is a virtual node representing a lan. The corresponding command for base class Node `$node is-lan?` always returns a 0.

## Chapter 15

# The (Revised) Addressing Structure in NS

This chapter describes the internals of the revised addressing format implemented in *ns*. The chapter consists of five sections. We describe the APIs that can be used for allocating bits to the *ns* addressing structure. The address space as described in chapter 3, can be thought of a contiguous field of  $n$  bits, where  $n$  may vary as per the address requirement of the simulation. The default value of  $n$  is 16 (as defined by *MAXADDRSIZE\_*). The maximum value of  $n$  is set to 32 (defined as *MAXADDRSIZE\_*). These default and maximum address sizes are defined in *~ns/tcl/lib/ns-default.tcl*.

The address space consists of 2 parts, the node-id and the port-id. The higher bits are assigned as the node's address or id\_ and remaining lower bits are assigned to form port-id or the identification of the agent attached to the node. Of the higher bits, 1 bit is assigned for multicast. The address space consists of 32 bits and port id space consists of 32 bits as well. The higher 32 bits for node-id, the MSB for multicast and the lower 32 bits for port-id. Additionally, the address space may also be set in hierarchical format, consisting of multiple levels of addressing hierarchy. We shall be describing the APIs for setting address structure in different formats as described above as well as expanding the address space. The procedures and functions described in this chapter can be found in *~ns/tcl/lib/ns-address.tcl*, *address.cc* and *address.h*.

### 15.1 The Default Address Format

The default settings allocates 32 lower bits for port-id, 1 higher bit for mcast and the rest 32 higher bits for node-id. The procedure to set the address format in default mode is called during initialisation of the simulator as:

```
# The preamble  
set ns [new Simulator] ;# initialise the simulation
```

It can also be called explicitly set as:

```
$ns set-address-format def
```

### 15.2 The Hierarchical Address Format

There are two options for setting an address to hierarchical format, the default and the specified.

### 15.2.1 Default Hierarchical Setting

The default hierarchical node-id consists of 3 levels with (10 11 11) bits in the three levels. The hierarchical configuration may be invoked as follows:

```
$ns set-address-format hierarchical
```

This sets :

- \* 32 bits for port-id, \* 32 bits for node-id assigned in - 3 levels of hierarchy - (10 11 11) bits for the three levels.
- or (9 11 11) if multicast is enabled.

### 15.2.2 Specific Hierarchical Setting

The second option allows a hierarchical address to be set with specified number of levels with number of bits assigned for each level. The API would be as the following:

```
$ns set-address-format hierarchical <#n hierarchy levels> <#bits for level1> <#bits for level 2> ....<#bits for nth level>
```

An example configuration would be:

```
$ns set-address-format hierarchical 2 8 15
```

where 2 levels of hierarchy is specified, assigning 8 bits for the 1st level and 15 bits for the second.

## 15.3 The Expanded Node-Address Format

NOTE: Please note that this command is now obsolete given that node address and port address spaces are 32 bits wide.

On the event of requirement of more bits to the address space, the expanded address API may be used as:

```
$ns set-address-format expanded
```

This expands the address space to 30 bits, allocating 22 higher bits to node-id and lower 8 bits to port-id.

## 15.4 Expanding port-id field

NOTE: Please note that this command is now obsolete given that node address and port address spaces are 32 bits wide.

This primitive may be used in case of need to expand portid in the event of requirement to attach a large number of agents to the nodes. This may be used in conjunction with set-address-format command (with different options) explained above. Synopsis for this command shall be:

```
expand-port-field-bits <#bits for portid>
```

expand-port-field-bits checks and raises error in the following if the requested portsize cannot be accommodated (i.e if sufficient num.of free bits are not available) or if requested portsize is less than or equal to the existing portsize.

## 15.5 Errors in setting address format

Errors are returned for both *set-address-format* and *expand-port-field-bits* primitives in the following cases:

- \* if number of bits specified is less than 0.
- \* if bit positions clash (contiguous number of requested free bits not found).
- \* if total number of bits exceed MAXADDRSIZE\_.
- \* if expand-port-field-bits is attempted with portbits less than or equal to the existing portsize.
- \* if number of hierarchy levels donot match with number of bits specified (for each level).

## 15.6 Commands at a glance

The following is a list of address-format related commands used in simulation scripts:

```
$ns_ set-address-format def
```

This command is used internally to set the address format to its default value of 32 lower bits for port-id, 1 higher bit for mcast and the rest 31 higher bits for node-id. However this API has been replaced by the new node API

```
$ns_ node-config -addressType flat.
```

```
$ns_ set-address-format hierarchical
```

This command is used to set the address format to the hierarchical configuration that consists of 3 levels with 8bits assigned to each level and 32 lower bits for port-id. However this API has been replaced by the new node API

```
$ns_ node-config -addressType hierarchical.
```

```
$ns_ set-address-format hierarchical <levels> <args>
```

This command is used to set the address format to a specific hierarchical setting. The <levels> indicate the number of levels of hierarchy in the addressing structure, while the args define number of bits for each level. An example would be `$ns_ set-address-format hierarchical 3 4 4 16`, where 4, 4 and 16 defines the number of bits to be used for the address space in level 1, 2 and 3 respectively.

```
$ns_ set-address-format expanded
```

THIS COMMAND IS NOW OBSOLETE This command was used to expand the address space to 30 bits, allocating 22 higher bits for node-id and lower 8 bits for port-id. However this command is obsoleted now by 32 bit addressing, i.e node-id field is 32 bit wide.

```
expand-port-field-bits <bits-for-portid>
```

THIS COMMAND IS NOW OBSOLETE Similar to the command above, this was used to expand the address space for the port-id field to <bits-for-portid> number of bits. However this command is obsolete now that the ports are 32 bit wide.

## Chapter 16

# Mobile Networking in ns

This chapter describes the wireless model that was originally ported as CMU's Monarch group's mobility extension to *ns*. This chapter consists of two sections and several subsections. The first section covers the original mobility model ported from CMU/Monarch group. In this section, we cover the internals of a mobilenode, routing mechanisms and network components that are used to construct the network stack for a mobilenode. The components that are covered briefly are Channel, Network-interface, Radio propagation model, MAC protocols, Interface Queue, Link layer and Address resolution protocol model (ARP). CMU trace support and Generation of node movement and traffic scenario files are also covered in this section. The original CMU model allows simulation of pure wireless LANs or multihop ad-hoc networks. Further extensions were made to this model to allow combined simulation of wired and wireless networks. MobileIP was also extended to the wireless model. These are discussed in the second section of this chapter.

### 16.1 The basic wireless model in ns

The wireless model essentially consists of the `MobileNode` at the core, with additional supporting features that allows simulations of multi-hop ad-hoc networks, wireless LANs etc. The `MobileNode` object is a split object. The C++ class `MobileNode` is derived from parent class `Node`. Refer to Chapter 5 for details on `Node`. A `MobileNode` thus is the basic `Node` object with added functionalities of a wireless and mobile node like ability to move within a given topology, ability to receive and transmit signals to and from a wireless channel etc. A major difference between them, though, is that a `MobileNode` is not connected by means of `Links` to other nodes or mobilenodes. In this section we shall describe the internals of `MobileNode`, its routing mechanisms, the routing protocols `dsdv`, `aodv`, `tora` and `dsr`, creation of network stack allowing channel access in `MobileNode`, brief description of each stack component, trace support and movement/traffic scenario generation for wireless simulations.

#### 16.1.1 Mobilenode: creating wireless topology

`MobileNode` is the basic *nsNode* object with added functionalities like movement, ability to transmit and receive on a channel that allows it to be used to create mobile, wireless simulation environments. The class `MobileNode` is derived from the base class `Node`. `MobileNode` is a split object. The mobility features including node movement, periodic position updates, maintaining topology boundary etc are implemented in C++ while plumbing of network components within `MobileNode` itself (like classifiers, `dmux`, `LL`, `Mac`, `Channel` etc) have been implemented in `Otcl`. The functions and procedures described in this subsection can be found in `~ns/mobilenode.{cc,h}`, `~ns/tcl/lib/ns-mobilenode.tcl`, `~ns/tcl/mobility/dsdv.tcl`, `~ns/tcl/mobility/dsr.tcl`, `~ns/tcl/mobility/tora.tcl`. Example scripts can be found in `~ns/tcl/ex/wireless-test.tcl` and `~ns/tcl/ex/wireless.tcl`.



While the first example uses a small topology of 3 nodes, the second example runs over a topology of 50 nodes. These scripts can be run simply by typing

```
$ns tcl/ex/wireless.tcl (or /wireless-test.tcl)
```

The four ad-hoc routing protocols that are currently supported are Destination Sequence Distance Vector (DSDV), Dynamic Source Routing (DSR), Temporally ordered Routing Algorithm (TORA) and Adhoc On-demand Distance Vector (AODV). The primitive to create a mobilenode is described below. Please note that the old APIs for creating a mobilenode depended on which routing protocol was used, like

```
set mnode [$opt(rp)-create-mobile-node $id]
```

where

```
$opt(rp)
```

defines "dsv", "aodv", "tora" or "dsr" and id is the index for the mobilenode. But the old API's use is being deprecated and the new API is described as follows:.

```
$ns_ node-config -adhocRouting $opt(adhocRouting)
                 -llType $opt(ll)
                 -macType $opt(mac)
                 -ifqType $opt(ifq)
                 -ifqLen $opt(ifqlen)
                 -antType $opt(ant)
                 -propInstance [new $opt(prop)]
                 -phyType $opt(netif)
                 -channel [new $opt(chan)]
                 -topoInstance $topo
                 -wiredRouting OFF
                 -agentTrace ON
                 -routerTrace OFF
                 -macTrace OFF
```

The above API configures for a mobilenode with all the given values of adhoc-routing protocol, network stack, channel, topography, propagation model, with wired routing turned on or off (required for wired-cum-wireless scenarios) and tracing turned on or off at different levels (router, mac, agent). In case hierarchical addressing is being used, the hierarchical address of the node needs to be passed as well. For more info about this command (part of new node APIs) see chapter titled "Restructuring ns node and new Node APIs" in ns Notes and Documentation.

Next actually create the mobilenodes as follows:

```
for { set j 0 } { $j < $opt(nn) } {incr j} {
    set node_($j) [ $ns_ node ]
    $node_($j) random-motion 0 ;# disable random motion
}
```

The above procedure creates a mobilenode (split) object, creates an adhoc-routing routing agent as specified, creates the network stack consisting of a link layer, interface queue, mac layer, and a network interface with an antenna, uses the defined

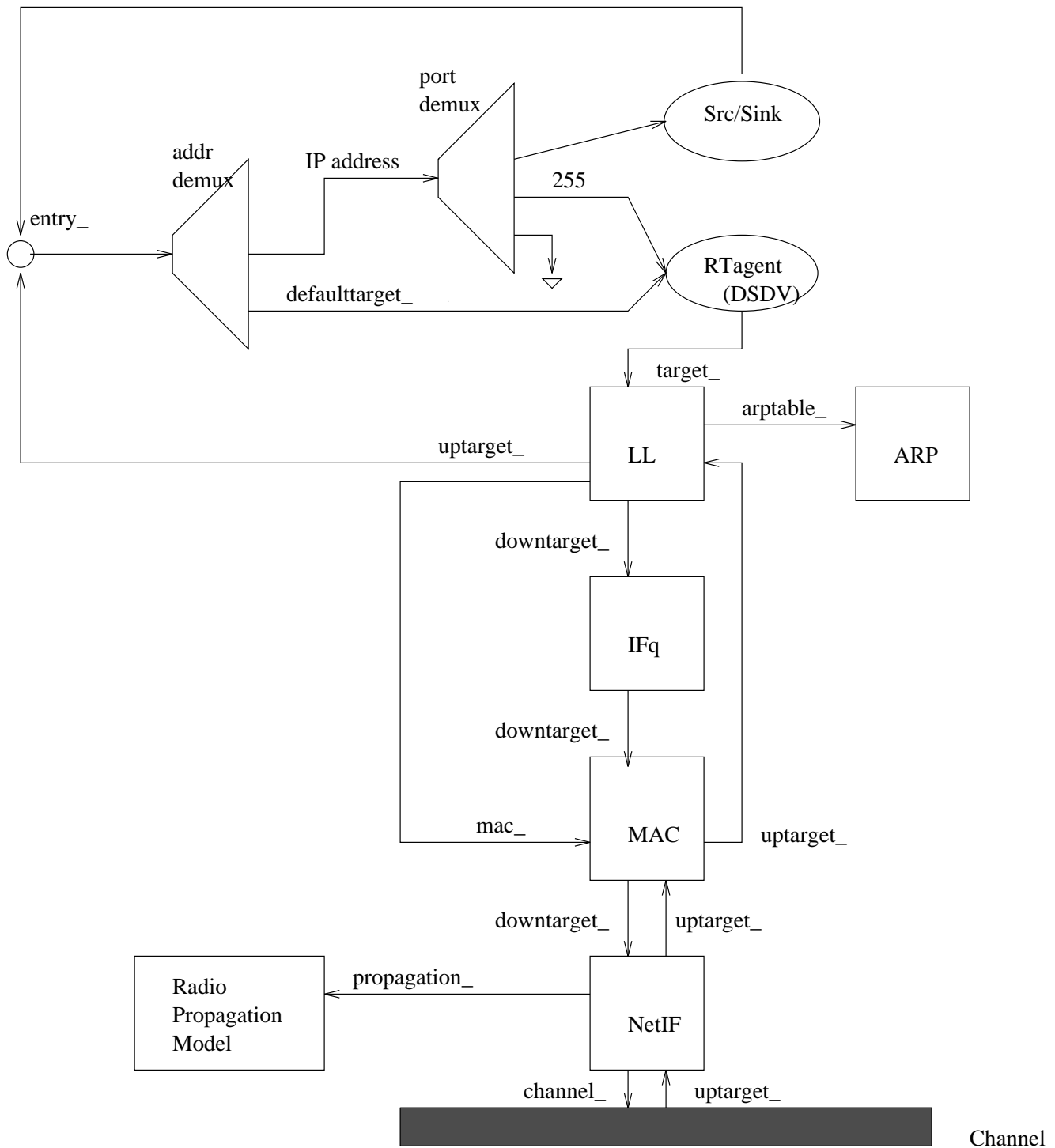


Figure 16.1: Schematic of a mobilenode under the CMU monarch's wireless extensions to *ns*

propagation model, interconnects these components and connects the stack to the channel. The mobilenode now looks like the schematic in Figure 16.1.

The mobilenode structure used for DSR routing is slightly different from the mobilenode described above. The class *SRNode* is derived from class *MobileNode*. *SRNode* does not use address demux or classifiers and all packets received by the node

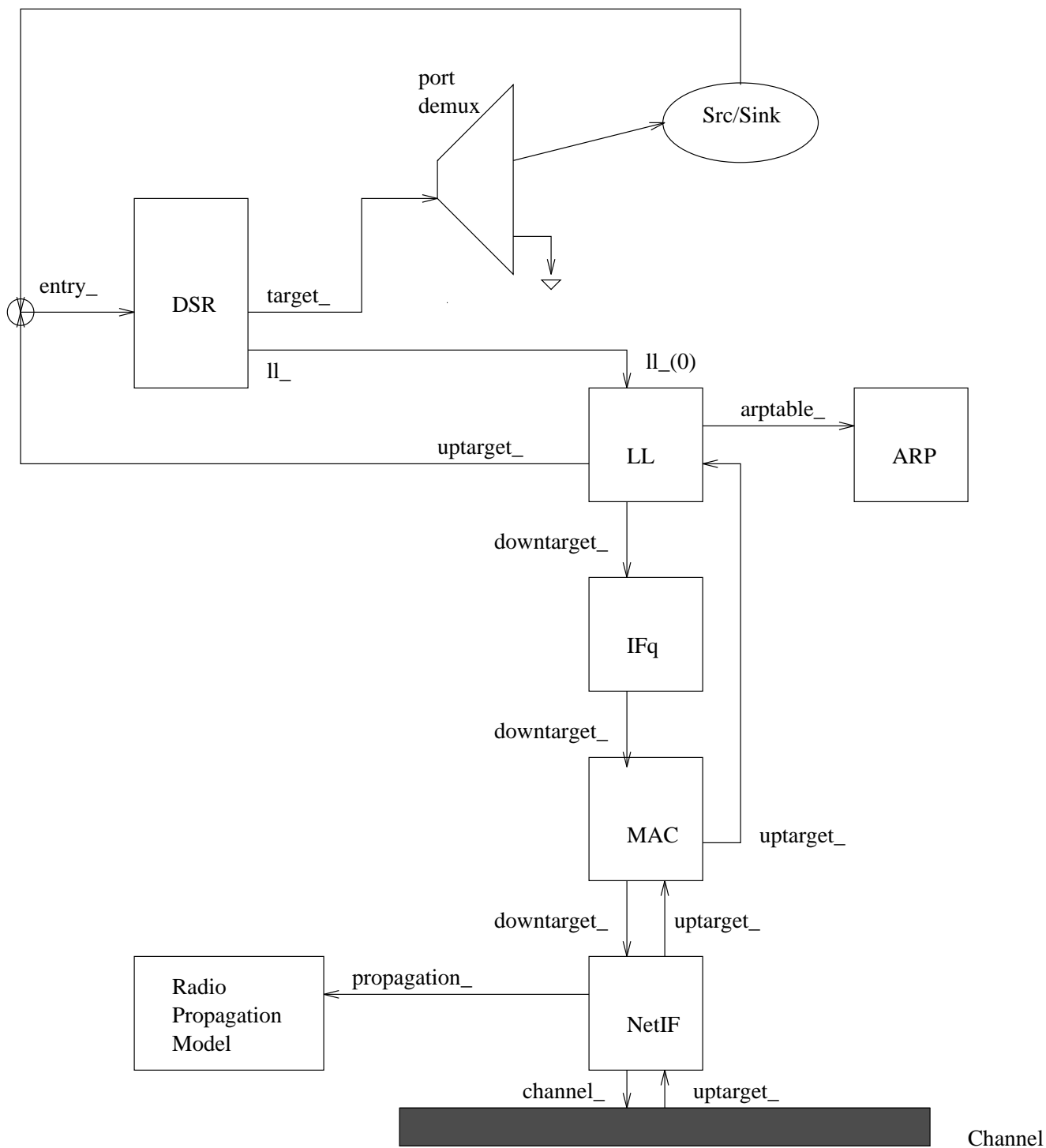


Figure 16.2: Schematic of a SRNode under the CMU monarch's wireless extensions to *ns*

are handed down to the DSR routing agent by default. The DSR routing agent either receives pkts for itself by handing it over to the port demux or forwards pkts as per source routes in the pkt hdr or sends out route requests and route replies for fresh packets. Details on DSR routing agent may be found in section 16.1.5. The schematic model for a SRNode is shown in Figure 16.2.

## 16.1.2 Creating Node movements

The mobilenode is designed to move in a three dimensional topology. However the third dimension (Z) is not used. That is the mobilenode is assumed to move always on a flat terrain with Z always equal to 0. Thus the mobilenode has X, Y, Z(=0) co-ordinates that is continually adjusted as the node moves. There are two mechanisms to induce movement in mobilenodes. In the first method, starting position of the node and its future destinations may be set explicitly. These directives are normally included in a separate movement scenario file.

The start-position and future destinations for a mobilenode may be set by using the following APIs:

```
$node set X_ <x1>
$node set Y_ <y1>
$node set Z_ <z1>

$ns at $time $node setdest <x2> <y2> <speed>
```

At \$time sec, the node would start moving from its initial position of (x1,y1) towards a destination (x2,y2) at the defined speed.

In this method the node-movement-updates are triggered whenever the position of the node at a given time is required to be known. This may be triggered by a query from a neighbouring node seeking to know the distance between them, or the setdest directive described above that changes the direction and speed of the node.

An example of a movement scenario file using the above APIs, can be found in `~ns/tcl/mobility/scen/scen-670x670-50-600-20-0`. Here 670x670 defines the length and width of the topology with 50 nodes moving at a maximum speed of 20m/s with average pause time of 600s. These node movement files may be generated using CMU's scenario generator to be found under `~ns/indep-utils/cmu-scen-gen/setdest`. See subsection 16.1.8 for details on generation of node movement scenarios.

The second method employs random movement of the node. The primitive to be used is:

```
$mobilenode start
```

which starts the mobilenode with a random position and have routined updates to change the direction and speed of the node. The destination and speed values are generated in a random fashion. We have not used the second method and leave it to the user to explore the details. The mobilenode movement is implemented in C++. See methods in `~ns/mobilenode.{cc,h}` for the implementational details.

Irrespective of the methods used to generate node movement, the topography for mobilenodes needs to be defined. It should be defined before creating mobilenodes. Normally flat topology is created by specifying the length and width of the topography using the following primitive:

```
set topo [new Topography]
$topo load_flatgrid $opt(x) $opt(y)
```

where opt(x) and opt(y) are the boundaries used in simulation.

The movement of mobilenodes may be logged by using a procedure like the following:

```
proc log-movement {} {
```

```

global logtimer ns_ ns

set ns $ns_
source ../mobility/timer.tcl
Class LogTimer -superclass Timer
LogTimer instproc timeout {} {
    global opt node_
    for {set i 0} {$i < $opt(nn)} {incr i} {
        $node_($i) log-movement
    }
    $self sched 0.1
}

set logtimer [new LogTimer]
$logtimer sched 0.1
}

```

In this case, mobilenode positions would be logged every 0.1 sec.

### 16.1.3 Network Components in a mobilenode

The network stack for a mobilenode consists of a link layer(LL), an ARP module connected to LL, an interface priority queue(IFq), a mac layer(MAC), a network interface(netIF), all connected to the channel. These network components are created and plumbed together in OTcl. The relevant MobileNode method add-interface() in *~ns/tcl/lib/ns-mobilenode.tcl* is shown below:

```

#
# The following setups up link layer, mac layer, network interface
# and physical layer structures for the mobile node.
#
Node/MobileNode instproc add-interface { channel pmodel
    lltype mactype qtype qlen iftype anttype } {

    $self instvar arptable_ nifs_
    $self instvar netif_ mac_ ifq_ ll_

    global ns_ MacTrace opt

    set t $nifs_
    incr nifs_

    set netif_($t) [new $iftype]           ;# net-interface
    set mac_($t)   [new $mactype]          ;# mac layer
    set ifq_($t)   [new $qtype]            ;# interface queue
    set ll_($t)    [new $lltype]           ;# link layer
    set ant_($t)   [new $anttype]

    #
    # Local Variables
    #
    set nullAgent_ [$ns_ set nullAgent_]
}

```

```

set netif $netif_($t)
set mac $mac_($t)
set ifq $ifq_($t)
set ll $ll_($t)

#
# Initialize ARP table only once.
#
if { $arptable_ == "" } {
    set arptable_ [new ARPTable $self $mac]
    set drpT [cmu-trace Drop "IFQ" $self]
    $arptable_ drop-target $drpT
}

#
# Link Layer
#
$ll arptable $arptable_
$ll mac $mac
$ll up-target [$self entry]
$ll down-target $ifq

#
# Interface Queue
#
$ifq target $mac
$ifq set qlim_ $qlen
set drpT [cmu-trace Drop "IFQ" $self]
$ifq drop-target $drpT

#
# Mac Layer
#
$mac netif $netif
$mac up-target $ll
$mac down-target $netif
$mac nodes $opt(nn)

#
# Network Interface
#
$netif channel $channel
$netif up-target $mac
$netif propagation $pmodel           ;# Propagation Model
$netif node $self                     ;# Bind node <---> interface
$netif antenna $ant_($t)              ;# attach antenna

#
# Physical Channel
#
$channel addif $netif                 ;# add to list of interfaces

# =====
# Setting up trace objects

```

```

if { $MacTrace == "ON" } {
    #
    # Trace RTS/CTS/ACK Packets
    #
    set rcvT [cmu-trace Recv "MAC" $self]
    $mac log-target $rcvT

    #
    # Trace Sent Packets
    #
    set sndT [cmu-trace Send "MAC" $self]
    $sndT target [$mac sendtarget]
    $mac sendtarget $sndT

    #
    # Trace Received Packets
    #
    set rcvT [cmu-trace Recv "MAC" $self]
    $rcvT target [$mac rcvtarget]
    $mac rcvtarget $rcvT

    #
    # Trace Dropped Packets
    #
    set drpT [cmu-trace Drop "MAC" $self]
    $mac drop-target $drpT
} else {
    $mac log-target [$ns_ set nullAgent_]
    $mac drop-target [$ns_ set nullAgent_]
}

# =====

$self addif $netif
}

```

The plumbing in the above method creates the network stack we see in Figure 16.1.

Each component is briefly described here. Hopefully more detailed documentation from CMU shall be available in the future.

**Link Layer** The LL used by mobilenode is same as described in Chapter 14. The only difference being the link layer for mobilenode, has an ARP module connected to it which resolves all IP to hardware (Mac) address conversions. Normally for all outgoing (into the channel) packets, the packets are handed down to the LL by the Routing Agent. The LL hands down packets to the interface queue. For all incoming packets (out of the channel), the mac layer hands up packets to the LL which is then handed off at the `node_entry_` point. The class LL is implemented in `~ns/ll.{cc,h}` and `~ns/tcl/lan/ns-ll.tcl`.

**ARP** The Address Resolution Protocol (implemented in BSD style) module receives queries from Link layer. If ARP has the hardware address for destination, it writes it into the mac header of the packet. Otherwise it broadcasts an ARP query, and caches the packet temporarily. For each unknown destination hardware address, there is a buffer for a single packet. In case additional packets to the same destination is sent to ARP, the earlier buffered packet is dropped. Once

the hardware address of a packet's next hop is known, the packet is inserted into the interface queue. The `class ARPTable` is implemented in `~ns/arp.{cc,h}` and `~ns/tcl/lib/ns-mobilenode.tcl`.

**Interface Queue** The `class PriQueue` is implemented as a priority queue which gives priority to routing protocol packets, inserting them at the head of the queue. It supports running a filter over all packets in the queue and removes those with a specified destination address. See `~ns/priqueue.{cc,h}` for interface queue implementation.

**Mac Layer** The IEEE 802.11 distributed coordination function (DCF) Mac protocol has been implemented by CMU. It uses a RTS/CTS/DATA/ACK pattern for all unicast packets and simply sends out DATA for all broadcast packets. The implementation uses both physical and virtual carrier sense. The `class Mac802_11` is implemented in `~ns/mac-802_11.{cc,h}`.

**Tap Agents** Agents that subclass themselves as `class Tap` defined in `mac.h` can register themselves with the mac object using method `installTap()`. If the particular Mac protocol permits it, the tap will promiscuously be given all packets received by the mac layer, before address filtering is done. See `~ns/mac.{cc,h}` for `class Tap` implementation.

**Network Interfaces** The Network Interphase layer serves as a hardware interface which is used by `mobilenode` to access the channel. The wireless shared media interface is implemented as `class Phy/WirelessPhy`. This interface subject to collisions and the radio propagation model receives packets transmitted by other node interfaces to the channel. The interface stamps each transmitted packet with the meta-data related to the transmitting interface like the transmission power, wavelength etc. This meta-data in `pkt` header is used by the propagation model in receiving network interface to determine if the packet has minimum power to be received and/or captured and/or detected (carrier sense) by the receiving node. The model approximates the DSSS radio interface (Lucent WaveLan direct-sequence spread-spectrum). See `~ns/phy.{cc,h}` and `~ns/wireless-phy.{cc,h}` for network interface implementations.

**Radio Propagation Model** It uses Friis-space attenuation ( $1/r^2$ ) at near distances and an approximation to Two ray Ground ( $1/r^4$ ) at far distances. The approximation assumes specular reflection off a flat ground plane. See `~ns/tworayground.{cc,h}` for implementation.

**Antenna** An omni-directional antenna having unity gain is used by `mobilenodes`. See `~ns/antenna.{cc,h}` for implementation details.

## 16.1.4 Different MAC layer protocols for mobile networking

In *ns*, two MAC layer protocols are implemented for mobile networks, which are 802.11 and TDMA. In this section we briefly discuss each of them.

### 802.11 MAC protocol

See `~ns/mac-802_11.{cc,h}` for implementation details.

### Preamble based TDMA protocol

**Note:** this works is still at a preliminary stage, some practical issues, such as: contention in the preamble phase and time slot reuse in a multi-hop environment are not considered.

Unlike contention based MAC protocol (802.11, for example), a TDMA MAC protocol allocates different time slots for nodes to send and receive packets. The superset of these time slots is called a TDMA frame.

Currently, *ns* supports a single hop, preamble-based TDMA MAC protocol. With this protocol, a TDMA frame contains preamble besides the data transmission slots. Within the preamble, every node has a dedicated subslot and uses it to broadcast



the destination node id of outgoing packet. Other nodes listen in the preamble and record the time slots to receive packets. Like other common TDMA protocols (GSM, for example), each node has a data transmission slot to send packets.

To avoid unnecessary power consumption, each node turns its radio on and off explicitly by invoking node API `set_node_sleep()`. The radio only needs to be on when: in the preamble phase (takes one slot time) and there is a packet to send and receive.

The preamble is implemented as a central data structure `tdma_preamble_`, which is accessible to all the nodes. At the beginning of a frame, each node writes the destination node id into its subslot in preamble if it has a packet to send. Following preamble phase, each node sends packet in its data transmission slot and checks the preamble to determine if there is a packet to receive in other slots.

The following parameters are user configurable: the wireless link bandwidth `bandwidth_`, the slot length `packet_slot_len_`, and the number of nodes `max_node_num_`. See `~ns/mac-tdma.{cc,h}` for implementation details.

## 16.1.5 Different types of Routing Agents in mobile networking

The four different ad-hoc routing protocols currently implemented for mobile networking in *ns* are dsdv, dsr, aodv and tora. In this section we shall briefly discuss each of them.

### DSDV

In this routing protocol routing messages are exchanged between neighbouring mobilenodes (i.e. mobilenodes that are within range of one another). Routing updates may be triggered or routine. Updates are triggered in case a routing information from one of the neighbours forces a change in the routing table. A packet for which the route to its destination is not known is cached while routing queries are sent out. The pkts are cached until route-replies are received from the destination. There is a maximum buffer size for caching the pkts waiting for routing information beyond which pkts are dropped.

All packets destined for the mobilenode are routed directly by the address dmux to its port dmux. The port dmux hands the packets to the respective destination agents. A port number of 255 is used to attach routing agent in mobilenodes. The mobilenodes also use a default-target in their classifier (or address demux). In the event a target is not found for the destination in the classifier (which happens when the destination of the packet is not the mobilenode itself), the pkts are handed to the default-target which is the routing agent. The routing agent assigns the next hop for the packet and sends it down to the link layer.

The routing protocol is mainly implemented in C++. See `~ns/dsdv` directory and `~ns/tcl/mobility/dsdv.tcl` for all procedures related to DSDV protocol implementation.

### DSR

This section briefly describes the functionality of the dynamic source routing protocol. As mentioned earlier the `SRNode` is different from the `MobileNode`. The `SRNode`'s `entry_` points to the DSR routing agent, thus forcing all packets received by the node to be handed down to the routing agent. This model is required for future implementation of piggy-backed routing information on data packets which otherwise would not flow through the routing agent.

The DSR agent checks every data packet for source-route information. It forwards the packet as per the routing information. In case it does not find routing information in the packet, it provides the source route, if route is known, or caches the packet and sends out route queries if route to destination is not known. Routing queries, always triggered by a data packet with no route to its destination, are initially broadcast to all neighbours. Route-replies are sent back either by intermediate nodes or the destination node, to the source, if it can find routing info for the destination in the route-query. It hands over all packets

destined to itself to the port dmux. In SRNode the port number 255 points to a null agent since the packet has already been processed by the routing agent.

See `~ns/dsr` directory and `~ns/tcl/mobility/dsr.tcl` for implementation of DSR protocol.

## TORA

Tora is a distributed routing protocol based on "link reversal" algorithm. At every node a separate copy of TORA is run for every destination. When a node needs a route to a given destination it broadcasts a QUERY message containing the address of the destination for which it requires a route. This packet travels through the network until it reaches the destination or an intermediate node that has a route to the destination node. This recipient node then broadcasts an UPDATE packet listing its height wrt the destination. As this node propagates through the network each node updates its height to a value greater than the height of the neighbour from which it receives the UPDATE. This results in a series of directed links from the node that originated the QUERY to the destination node. If a node discovers a particular destination to be unreachable it sets a local maximum value of height for that destination. In case the node cannot find any neighbour having finite height wrt this destination it attempts to find a new route. In case of network partition, the node broadcasts a CLEAR message that resets all routing states and removes invalid routes from the network.

TORA operates on top of IMEP (Internet MANET Encapsulation Protocol) that provides reliable delivery of route-messages and informs the routing protocol of any changes of the links to its neighbours. IMEP tries to aggregate IMEP and TORA messages into a single packet (called block) in order to reduce overhead. For link-status sensing and maintaining a list of neighbour nodes, IMEP sends out periodic BEACON messages which is answered by each node that hears it by a HELLO reply message. See `ns/tora` directory and `ns/tcl/mobility/tora.tcl` for implementation of tora in `ns`.

## AODV

AODV is a combination of both DSR and DSDV protocols. It has the basic route-discovery and route-maintenance of DSR and uses the hop-by-hop routing, sequence numbers and beacons of DSDV. The node that wants to know a route to a given destination generates a ROUTE REQUEST. The route request is forwarded by intermediate nodes that also creates a reverse route for itself from the destination. When the request reaches a node with route to destination it generates a ROUTE REPLY containing the number of hops requires to reach destination. All nodes that participates in forwarding this reply to the source node creates a forward route to destination. This state created from each node from source to destination is a hop-by-hop state and not the entire route as is done in source routing. See `ns/aodv` and `ns/tcl/lib/ns-lib.tcl` for implementational details of aodv.

### 16.1.6 Trace Support

The trace support for wireless simulations currently use cmu-trace objects. In the future this shall be extended to merge with trace and monitoring support available in ns, which would also include nam support for wireless modules. For now we will explain briefly with cmu-trace objects and how they may be used to trace packets for wireless scenarios.

The cmu-trace objects are of three types - CMUTrace/Drop, CMUTrace/Recv and CMUTrace/Send. These are used for tracing packets that are dropped, received and sent by agents, routers, mac layers or interface queues in `ns`. The methods and procedures used for implementing wireless trace support can be found under `~ns/trace.{cc,h}` and `~ns/tcl/lib/ns-cmutrace.tcl`.

A cmu-trace object may be created by the following API:

```
set sndT [cmu-trace Send "RTR" $self]
```

which creates a trace object, `sndT`, of the type `CMUTrace`/`Send` for tracing all packets that are sent out in a router. The trace objects may be used to trace packets in MAC, agents (routing or others), routers or any other `NsObject`.

The `cmu-trace` object `CMUTrace` is derived from the base class `Trace`. See Chapter 25 for details on class `Trace`. The class `CMUTrace` is defined as the following:

```
class CMUTrace : public Trace {
public:
    CMUTrace(const char *s, char t);
    void      recv(Packet *p, Handler *h);
    void      recv(Packet *p, const char* why);

private:
    int off_arp_;
    int off_mac_;
    int off_sr_;

    char      tracename[MAX_ID_LEN + 1];
    int       tracetype;
    MobileNode *node_;

    int initialized() { return node_ && 1; }

    int      command(int argc, const char*const* argv);
    void     format(Packet *p, const char *why);

    void     format_mac(Packet *p, const char *why, int offset);
    void     format_ip(Packet *p, int offset);

    void     format_arp(Packet *p, int offset);
    void     format_dsr(Packet *p, int offset);
    void     format_msg(Packet *p, int offset);
    void     format_tcp(Packet *p, int offset);
    void     format_rtp(Packet *p, int offset);
};
```

The `type` field (described in `Trace` class definition) is used to differentiate among different types of traces. For `cmu-trace` this can be **s** for sending, **r** for receiving or **D** for dropping a packet. A fourth type **f** is used to denote forwarding of a packet (When the node is not the originator of the packet). Similar to the method `Trace::format()`, the `CMUTrace::format()` defines and dictates the trace file format. The method is shown below:

```
void CMUTrace::format(Packet* p, const char *why)
{
    hdr_cmh *ch = HDR_CMH(p);
    int offset = 0;

    /*
     * Log the MAC Header
     */
    format_mac(p, why, offset);
    offset = strlen(why);
```

```

switch(ch->ptype()) {

case PT_MAC:
    break;

case PT_ARP:
    format_arp(p, offset);
    break;

default:
    format_ip(p, offset);
    offset = strlen(wrk_);

    switch(ch->ptype()) {

        case PT_DSR:
            format_dsr(p, offset);
            break;

        case PT_MESSAGE:
        case PT_UDP:
            format_msg(p, offset);
            break;

        case PT_TCP:
        case PT_ACK:
            format_tcp(p, offset);
            break;

        case PT_CBR:
            format_rtp(p, offset);
            break;
        .....
    }
}
}

```

The above function calls different format functions depending on the type of the packet being traced. All traces are written to the buffer wrk\_. A count of the offset for the buffer is kept and is passed along the different trace functions. The most basic format is defined by format\_mac() and is used to trace all pkt types. The other format functions print additional information as defined by the packet types. The mac format prints the following:

```

#ifdef LOG_POSITION
    double x = 0.0, y = 0.0, z = 0.0;
    node_>getLoc(&x, &y, &z);
#endif

    sprintf(wrk_ + offset,
#ifdef LOG_POSITION
        "%c %.9f %d (%6.2f %6.2f) %3s %4s %d %s %d [%x %x %x %x] ",
#else
        "%c %.9f _%d_ %3s %4s %d %s %d [%x %x %x %x] ",

```

```

#endif
                op,                // s, r, D or f
                Scheduler::instance().clock(), // time stamp
                src_,                // the nodeid for this node
#ifdef LOG_POSITION
                x,                // x co-ord
                y,                // y co-ord
#endif
                tracename,        // name of object type tracing
                why,                // reason, if any

                ch->uid(),          // identifier for this event
                packet_info.name(ch->ptype()), // packet type
                ch->size(),          // size of cmn header
                mh->dh_duration,     // expected time to send data
                ETHER_ADDR(mh->dh_da), // mac_destination address
                ETHER_ADDR(mh->dh_sa), // mac_sender address
                GET_ETHER_TYPE(mh->dh_body)); // type - arp or IP

```

If the LOG\_POSITION is defined the x and y co-ordinates for the mobilenode is also printed. The descriptions for different fields in the mac trace are given in the comments above. For all IP packets additional IP header fields are also added to the above trace. The IP trace is described below:

```

sprintf(wrk_ + offset, "----- [%d:%d %d:%d %d %d] ",
        src,                // IP src address
        ih->sport_,          // src port number
        dst,                // IP dest address
        ih->dport_,          // dest port number
        ih->ttl_,            // TTL value
        (ch->next_hop_ < 0) ? 0 : ch->next_hop_); // next hopaddress, if any.

```

An example of a trace for a tcp packet is as follows:

```

r 160.093884945 _6_ RTR --- 5 tcp 1492 [a2 4 6 800] ----- [655
36:0 16777984:0 31 16777984] [1 0] 2 0

```

Here we see a TCP data packet being received by a node with id of 6. UID of this pkt is 5 with a cmn hdr size of 1492. The mac details shows an IP pkt (ETHERTYPE\_IP is defined as 0x0800, ETHERTYPE\_ARP is 0x0806), mac-id of this receiving node is 4. That of the sending node is 6 and expected time to send this data pkt over the wireless channel is a2 (hex2dec conversion: 160+2 sec). Additionally, IP traces information about IP src and destination addresses. The src translates (using a 3 level hier-address of 8/8/8) to a address string of 0.1.0 with port of 0. The dest address is 1.0.3 with port address of 0. The TTL value is 31 and the destination was a hop away from the src. Additionally TCP format prints information about tcp seqno of 1, ackno of 0. See other formats described in `~ns//cmu-trace.cc` for DSR, UDP/MESSAGE, TCP/ACK and CBR packet types.

Other trace formats are also used by the routing agents (TORA and DSR) to log certain special routing events like "originating" (adding a SR header to a packet) or "ran off the end of a source route" indicating some sort of routing problem with the source route etc. These special event traces begin with "S" for DSR and "T" for Tora and may be found in `~ns//tora/tora.cc` for TORA and `~ns//dsr/dsrgent.cc` for DSR routing agent.

## 16.1.7 Revised format for wireless traces

In an effort to merge wireless trace, using `cmu-trace` objects, with `ns` tracing, a new, improved trace format has been introduced. This revised trace support is backwards compatible with the old trace formatting and can be enabled by the following command:

```
$ns use-newtrace
```

This command should be called before the universal trace command `$ns trace-all <trace-fd>`. Primitive `use-newtrace` sets up new format for wireless tracing by setting a simulator variable called `newTraceFormat`. Currently this new trace support is available for wireless simulations only and shall be extended to rest of `ns` in the near future.

An example of the new trace format is shown below:

```
s -t 0.267662078 -Hs 0 -Hd -1 -Ni 0 -Nx 5.00 -Ny 2.00 -Nz 0.00 -Ne  
-1.000000 -Nl RTR -Nw --- -Ma 0 -Md 0 -Ms 0 -Mt 0 -Is 0.255 -Id -1.255 -It  
message -Il 32 -If 0 -Ii 0 -Iv 32  
s -t 1.511681090 -Hs 1 -Hd -1 -Ni 1 -Nx 390.00 -Ny 385.00 -Nz 0.00 -Ne  
-1.000000 -Nl RTR -Nw --- -Ma 0 -Md 0 -Ms 0 -Mt 0 -Is 1.255 -Id -1.255 -It  
message -Il 32 -If 0 -Ii 1 -Iv 32  
s -t 10.0000000000 -Hs 0 -Hd -2 -Ni 0 -Nx 5.00 -Ny 2.00 -Nz 0.00 -Ne  
-1.000000 -Nl AGT -Nw --- -Ma 0 -Md 0 -Ms 0 -Mt 0 -Is 0.0 -Id 1.0 -It tcp -Il 1000 -If  
2 -Ii 2 -Iv 32 -Pn tcp -Ps 0 -Pa 0 -Pf 0 -Po 0  
r -t 10.0000000000 -Hs 0 -Hd -2 -Ni 0 -Nx 5.00 -Ny 2.00 -Nz 0.00 -Ne  
-1.000000 -Nl RTR -Nw --- -Ma 0 -Md 0 -Ms 0 -Mt 0 -Is 0.0 -Id 1.0 -It tcp -Il 1000 -If  
2 -Ii 2 -Iv 32 -Pn tcp -Ps 0 -Pa 0 -Pf 0 -Po 0  
r -t 100.004776054 -Hs 1 -Hd 1 -Ni 1 -Nx 25.05 -Ny 20.05 -Nz 0.00 -Ne  
-1.000000 -Nl AGT -Nw --- -Ma a2 -Md 1 -Ms 0 -Mt 800 -Is 0.0 -Id 1.0 -It  
tcp -Il 1020 -If 2 -Ii 21 -Iv 32 -Pn tcp -Ps 0 -Pa 0 -Pf 1 -Po 0  
s -t 100.004776054 -Hs 1 -Hd -2 -Ni 1 -Nx 25.05 -Ny 20.05 -Nz 0.00 -Ne  
-1.000000 -Nl AGT -Nw --- -Ma 0 -Md 0 -Ms 0 -Mt 0 -Is 1.0 -Id 0.0 -It ack -Il 40  
-If 2 -Ii 22 -Iv 32 -Pn tcp -Ps 0 -Pa 0 -Pf 0 -Po 0
```

### Explanation of new trace format

The new trace format as seen above can be divided into the following fields :

**Event type** In the traces above, the first field (as in the older trace format) describes the type of event taking place at the node and can be one of the four types:

- s** send
- r** receive
- d** drop
- f** forward

**General tag** The second field starting with "-t" may stand for time or global setting

- t** time
- t \*** (global setting)

**Node property tags** This field denotes the node properties like node-id, the level at which tracing is being done like agent, router or MAC. The tags start with a leading "-N" and are listed as below:

- Ni: node id
- Nx: node's x-coordinate
- Ny: node's y-coordinate
- Nz: node's z-coordinate
- Ne: node energy level
- Nl: trace level, such as AGT, RTR, MAC
- Nw: reason for the event. The different reasons for dropping a packet are given below:
  - "END" DROP\_END\_OF\_SIMULATION
  - "COL" DROP\_MAC\_COLLISION
  - "DUP" DROP\_MAC\_DUPLICATE
  - "ERR" DROP\_MAC\_PACKET\_ERROR
  - "RET" DROP\_MAC\_RETRY\_COUNT\_EXCEEDED
  - "STA" DROP\_MAC\_INVALID\_STATE
  - "BSY" DROP\_MAC\_BUSY
  - "NRTE" DROP\_RTR\_NO\_ROUTE i.e no route is available.
  - "LOOP" DROP\_RTR\_ROUTE\_LOOP i.e there is a routing loop
  - "TTL" DROP\_RTR\_TTL i.e TTL has reached zero.
  - "TOUT" DROP\_RTR\_QTIMEOUT i.e packet has expired.
  - "CBK" DROP\_RTR\_MAC\_CALLBACK
  - "IFQ" DROP\_IFQ\_QFULL i.e no buffer space in IFQ.
  - "ARP" DROP\_IFQ\_ARP\_FULL i.e dropped by ARP
  - "OUT" DROP\_OUTSIDE\_SUBNET i.e dropped by base stations on receiving routing updates from nodes outside its domain.

**Packet information at IP level** The tags for this field start with a leading "-I" and are listed along with their explanations as following:

- Is: source address.source port number
- Id: dest address.dest port number
- It: packet type
- Il: packet size
- If: flow id
- Ii: unique id
- Iv: ttl value

**Next hop info** This field provides next hop info and the tag starts with a leading "-H".

- Hs: id for this node
- Hd: id for next hop towards the destination.

**Packet info at MAC level** This field gives MAC layer information and starts with a leading "-M" as shown below:

- Ma: duration
- Md: dst's ethernet address
- Ms: src's ethernet address
- Mt: ethernet type

**Packet info at "Application level"** The packet information at application level consists of the type of application like ARP, TCP, the type of adhoc routing protocol like DSDV, DSR, AODV etc being traced. This field consists of a leading "-P" and list of tags for different application is listed as below:

**-P arp** Address Resolution Protocol. Details for ARP is given by the following tags:

- Po: ARP Request/Reply
- Pm: src mac address
- Ps: src address
- Pa: dst mac address
- Pd: dst address

**-P dsr** This denotes the adhoc routing protocol called Dynamic source routing. Information on DSR is represented by the following tags:

- Pn: how many nodes traversed
- Pq: routing request flag
- Pi: route request sequence number
- Pp: routing reply flag
- Pl: reply length
- Pe: src of srcrouting->dst of the source routing
- Pw: error report flag ?
- Pm: number of errors
- Pc: report to whom
- Pb: link error from linka->linkb

**-P cbr** Constant bit rate. Information about the CBR application is represented by the following tags:

- Pi: sequence number
- Pf: how many times this pkt was forwarded
- Po: optimal number of forwards

**-P tcp** Information about TCP flow is given by the following subtags:

- Ps: seq number
- Pa: ack number
- Pf: how many times this pkt was forwarded
- Po: optimal number of forwards

This field is still under development and new tags shall be added for other applications as they get included along the way.

## 16.1.8 Generation of node-movement and traffic-connection for wireless scenarios

Normally for large topologies, the node movement and traffic connection patterns are defined in separate files for convenience. These movement and traffic files may be generated using CMU's movement- and connection-generators. In this section we shall describe both separately.

### MobileNode Movement

Some examples of node movement files may be found in `~ns/tcl/mobility/scen/scen-670x670-50-600-20-*`. These files define a topology of 670 by 670m where 50 nodes move with a speed of 20m/s with pause time of 600s. each node is assigned



a starting position. The information regarding number of hops between the nodes is fed to the central object "GOD" (XXX but why/where is this information used??-answer awaited from CMU.) Next each node is a speed and a direction to move to.

The generator for creating node movement files are to be found under `~ns/indep-utils/cmu-scen-gen/setdest/` directory. Compile the files under setdest to create an executable. run setdest with arguments in the following way:

```
./setdest -n <num_of_nodes> -p <pausetime> -s <maxspeed> -t <simtime>  
          -x <maxx> -y <maxy> > <outdir>/<scenario-file>
```

Note that the index used for nodes now start from 0 instead of 1 as was in the original CMU version, to match with *ns*'s tradition of assigning node indices from 0.

### Generating traffic pattern files

The examples for traffic patterns may be found in `~ns/tcl/mobility/scenecbr-50-{10-4-512,20-4-512}`.

The traffic generator is located under `~ns/indep-utils/cmu-scen-gen/` and are called `cbrgen.tcl` and `tcpgen.tcl`. They may be used for generating CBR and TCP connections respectively.

To create CBR connections, run

```
ns cbrgen.tcl [-type cbr|tcp] [-nn nodes] [-seed seed]  
              [-mc connections] [-rate rate]
```

To create TCP connections, run

```
ns tcpgen.tcl [-nn nodes] [-seed seed]
```

You will need to pipe the outputs from above to a `cbr-*` or a `tcp-*` file.

## 16.2 Extensions made to CMU's wireless model

As mentioned earlier, the original CMU wireless model allows simulation of wireless LANs and ad-hoc networks. However in order to use the wireless model for simulations using both wired and wireless nodes we had to add certain extensions to cmu model. We call this wired-cum-wireless feature. Also SUN's MobileIP (implemented for wired nodes) was integrated into the wireless model allowing mobileIP to run over wireless mobilenodes. The following two subsections describe these two extensions to the wireless model in *ns*.

### 16.2.1 wired-cum-wireless scenarios

The mobilenodes described so far mainly supports simulation of multi-hop ad-hoc networks or wireless LANs. But what if we need to simulate a topology of multiple wireless LANs connected through wired nodes, or may need to run mobileIP on top of these wireless nodes? The extensions made to the CMU wireless model allows us to do that.

The main problem facing the wired-cum-wireless scenario was the issue of routing. In ns, routing information is generated based on the connectivity of the topology, i.e how nodes are connected to one another through `Link`s. Mobilenodes on the other hand have no concept of links. They route packets among themselves, within the wireless topology, using their routing protocol. so how would packets be exchanged between these two types of nodes?

So a node called `BaseStationNode` is created which plays the role of a gateway for the wired and wireless domains. The `BaseStationNode` is essentially a hybrid between a Hierarchical node<sup>1</sup> (`HierNode`) and a `MobileNode`. The basestation node is responsible for delivering packets into and out of the wireless domain. In order to achieve this we need Hierarchical routing.

Each wireless domain along with its base-station would have an unique domain address assigned to them. All packets destined to a wireless node would reach the base-station attached to the domain of that wireless node, who would eventually hand the packet over to the destination (mobilenode). And mobilenodes route packets, destined to outside their (wireless) domain, to their base-station node. The base-station knows how to forward these packets towards the (wired) destination. The schematic of a `BaseStationNode` is shown in Figure 16.3.

The mobilenodes in wired-cum-wireless scenario are required to support hierarchical addressing/routing. Thus the `MobileNode` looks exactly like the `BaseStationNode`. The `SRNode`, however, simply needs to have its own hier-address since it does not require any address demuxes and thus is not required to support hier routing<sup>2</sup>.

The DSDV agent on having to forward a packet checks to see if the destination is outside its (wireless) subnet. If so, it tries to forward the packet to its base-station node. In case no route to base-station is found the packet is dropped. Otherwise the packet is forwarded to the `next_hop` towards the base-station. Which is then routed towards the wired network by base-station's classifiers.

The DSR agent, on receiving a pkt destined outside its subnet, sends out a route-query for its base-station in case the route to base-station is not known. The data pkt is temporarily cached while it waits to hear route replies from base-station. On getting a reply the packet is provided with routing information in its header and send away towards the base-station. The base-station address demuxes routes it correctly toward the wired network.

The example script for a wired-cum-wireless simulation can be found at `~ns/tcl/ex/wired-cum-wireless-sim.tcl`. The methods for wired-cum-wireless implementations are defined in `~ns/tcl/lib/ns-bsnode.tcl`, `~ns/tcl/mobility/{com.tcl,dsr.tcl,dsdv.tcl}`, `~ns/dsdv/dsdv.{cc,h}` and `~ns/dsr/dsragent.{cc,h}`.

## 16.2.2 MobileIP

The wired-cum-wireless extensions for the wireless model paved the path for supporting wireless MobileIP in ns. Sun Microsystem's (Charlie Perkins *et al*) MobileIP model was based on ns's wired model (consisting of `Node`'s and `Link`'s) and thus didnot use CMU's mobility model.

Here we briefly describe the wireless MobileIP implementation. We hope that Sun would provide the detailed version of the documentation in the future.

The mobileIP scenario consists of Home-Agents(HA) and Foreign-Agents(FA) and have Mobile-Hosts(MH) moving between their HA and FAs. The HA and FA are essentially base-station nodes we have described earlier. While MHs are basically the mobileNodes described in section 16.1.1. The methods and procedures for MobileIP extensions are described in `~ns/mip.{cc,h}`, `~ns/mip-reg.cc`, `~ns/tcl/lib/ns-mip.tcl` and `~ns/tcl/lib/ns-wireless-mip.tcl`.

The HA and FA nodes are defined as `MobileNode/MIPBS` having a registering agent (`regagent_`) that sends beacon out to

---

<sup>1</sup>Refer to Chapter 31 for details on hierarchical routing and internals of `HierNode`.

<sup>2</sup>In order to do away with all these different variations of the definition of a node, we are planning to revise the node architecture that would allow a more flexible and modularised construction of a node without the necessity of having to define and be limited to certain Class definitions only.

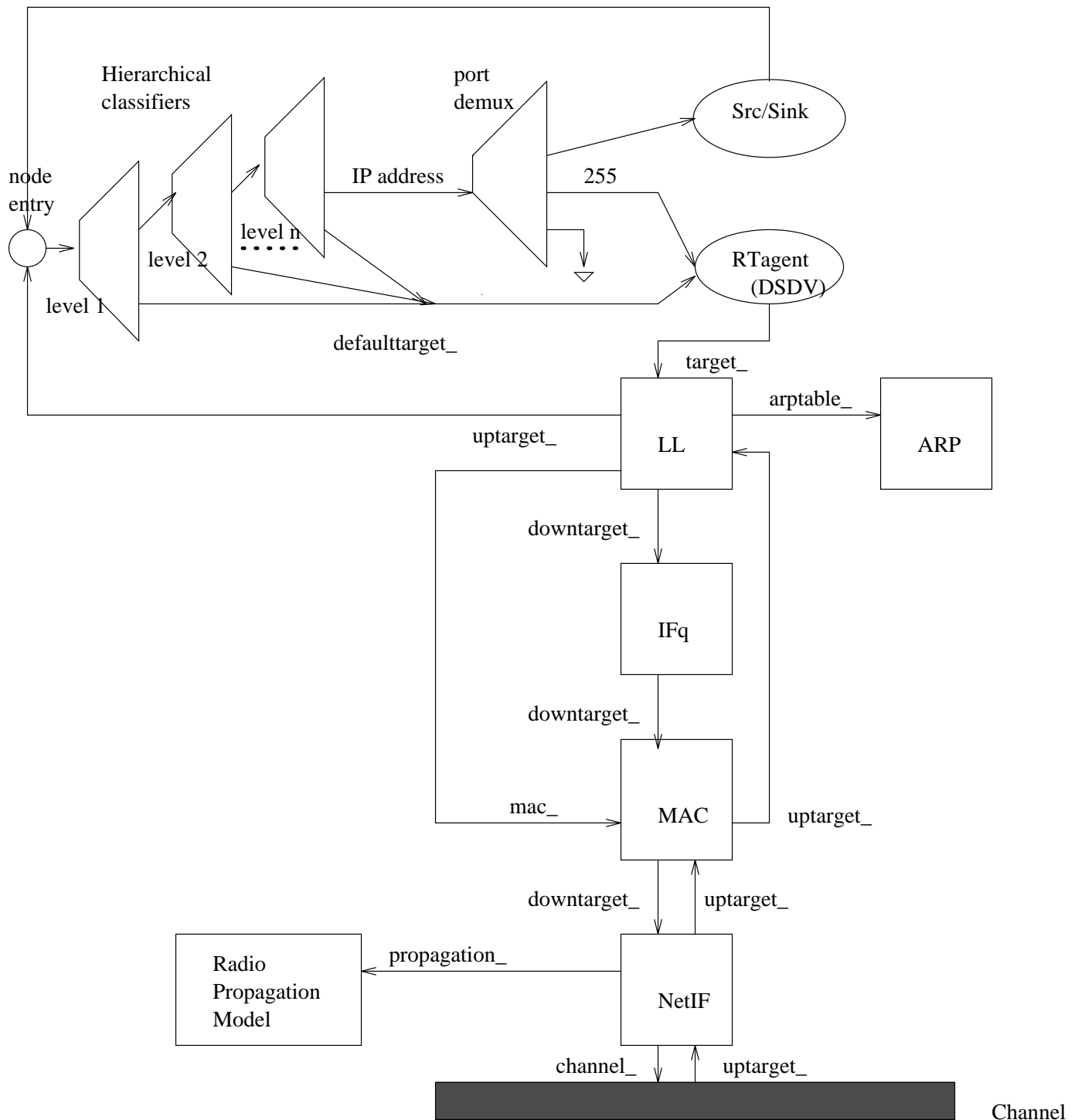


Figure 16.3: Schematic of a `baseStationNode`

the mobilenodes, sets up encapsulator and decapsulator, as required and replies to solicitations from MHs. The MH nodes are defined as `MobileNode/MIPMH` which too have a `regagent_` that receives and responds to beacons and sends out solicitations to HA or FAs. Figure 16.4 illustrates the schematic of a `MobileNode/MIPBS` node. The `MobileNode/MIPMH` node is very similar to this except for the fact that it doesnot have any encapsulator or decapsulator. As for the `SRNode` version of a MH, it doesnot have the hierarchical classifiers and the RA agent forms the entry point of the node. See Figure 16.2 for

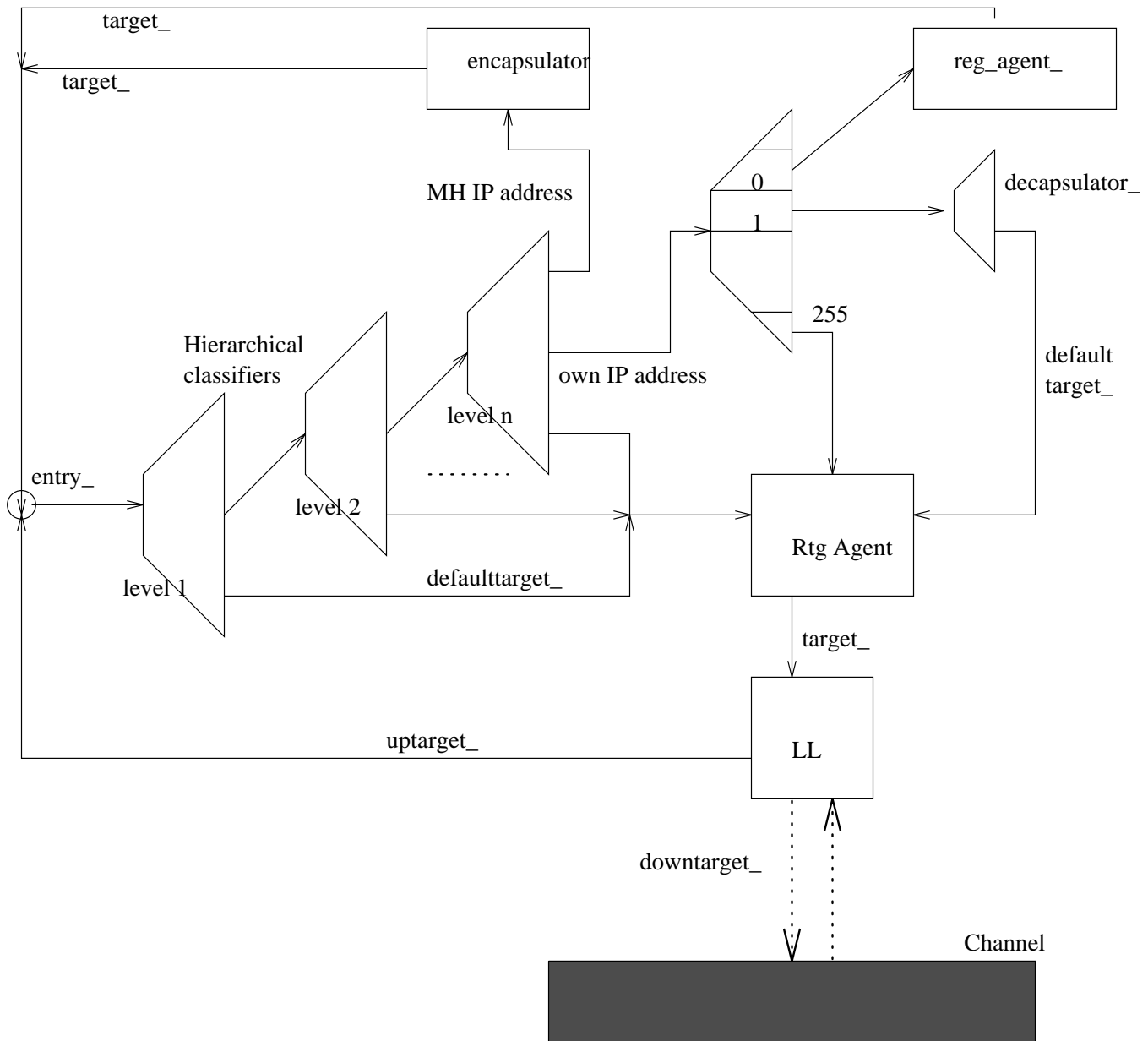


Figure 16.4: Schematic of a Wireless MobileIP BaseStation Node

model of a SRNode.

The MobileNode/MIPBS node routinely broadcasts beacon or advertisement messages out to MHs. A solicitation from a mobilenode generates an ad that is send directly to the requesting MH. The address of the base-station sending out beacon is heard by MH and is used as the COA (care-of-address) of the MH. Thus as the MH moves from its native to foreign domains, its COA changes. Upon receiving reg\_request (as reply to ads) from a mobilehost the base-station checks to see if it is the HA for the MH. If not, it sets up its decapsulator and forwards the reg\_request towards the HA of the MH.

In case the base-station *is* the HA for the requesting MH but the COA doesnt match its own, it sets up an encapsulator and sends reg-request-reply back to the COA (address of the FA) who has forwarded the reg\_request to it. so now all packets destined to the MH reaching the HA would be tunneled through the encapsulator which encapsulates the IP pkthdr with a IPinIP hdr, now destined to the COA instead of MH. The FA's decapsulator recives this packet, removes the encapsulation

and sends it to the MH.

If the COA matches that of the HA, it just removes the encapsulator it might have set up (when its mobilehost was roaming into foreign networks) and sends the reply directly back to the MH, as the MH have now returned to its native domain.

The mobilehost sends out solicitations if it doesnot hear any ads from the base-stations. Upon receiving ads, it changes its COA to the address of the HA/FA it has heard the ad from, and replies back to the COA with a request for registration (`reg-request`). Initially the MH maybe in the range of the HA and receives all pkts directly from its COA which is HA in this case. Eventually as the MH moves out of range of its HA and into the a foreign domain of a FA, the MH's COA changes from its HA to that of the FA. The HA now sets up an encapsulator and tunnels all pkts destined for MH towards the FA. The FA decapsulates the pkts and hands them over to the MH. The data from MH destined for the wired world is always routed towards its current COA. An example script for wireless mobileIP can be found at `~ns/tcl/ex/wireless-mip-test.tcl`. The simulation consists of a MH moving between its HA and a FA. The HA and FA are each connected to a wired domain on one side and to their wireless domains on the other. TCP flows are set up between the MH and a wired node.

## 16.3 Lists of changes for merging code developed in older version of ns (2.1b5 or later) into the current version (2.1b8)

The CMU-wireless model developed by David Johnhson's Monarch project was merged into ns around 1998-99 in what was then the ns-2.1b5 version. Since then the ns versions used by Monarch and by us here at ISI have forked quite a bit. Recently we ported a newer version of DSR developed by the Monarch group back into ns and in the process have created a list of changes that were required to be made for the merge. Hopefully this list will be helpful for those who have been working on older versions of ns from around that time or or later, to have their stuff merged in to the current version of ns-2.1b8.

The following lists of changes are required for merging the cmu version of ns (2.1b5) in to current version of 2.1b8. Each change is followed by a brief explanation for why the change was made.

Methods for accessing pkt hdrs have changed from

```
(hdr_sr *)p->access(off_sr)
```

to a static access method defi ned for each hdr, as

```
hdr_sr::access(p)
```

where for class `hdr_sr` a static method `access()` is defi ned as

```
inline static hdr_sr* access(const Packet* p)
return (hdr_sr*)p->access(offset_);
```

why: This change avoids using casts everywhere.

As the method for accessing hdrs have changed, there is no need to explicitly bind the hdr offset values. This is now done while establishing tcl linkage for the individual hdr classes. so lines like

```
bind("off_SR_", &off_sr_);
bind("off_ll_", &off_ll_);
bind("off_mac_", &off_mac_);
bind("off_ip_", &off_ip_);
should be removed.
```

AF\_ enumerations replaced by NS\_AF\_ as in

```
enum ns_af_enum NS_AF_NONE, NS_AF_ILINK, NS_AF_INET ;
```

why: This avoids header clashes between ns and the OS.

The ip\_hdr (dst/src) address fields that used to be integers are now defined as structures called ns\_addr\_t. ns\_addr\_t has 2 members address\_ and port\_ that are both defined as int. Hence lines like

```
iph->src() should change to  
iph->saddr() & iph->sport();
```

Also lines like

```
dst_ = (IP_BROADCAST << 8) | RT_PORT  
should be replaced by
```

```
dst_.addr_ = IP_BROADCAST;  
dst_.port_ = RT_PORT;
```

Why: This extension supports 32bit addressing.

The addr\_ member for hsr\_sr class has a separate function for returning its value. Thus need to call hsr.addr\_() instead of hsr.addr.

why: addr\_ is now a private variable which is accessed by public function addr\_().

All includes that had absolute paths by using <> were replaced by ". ". Thus

```
<cmu/dsr/dsragent.h>  
was changed to  
"cmu/dsr/dsragent.h"
```

The tcl command "ip-addr" was changed to "addr".

Other new tcl commands like "node", "port-dmux" and "trace-target" were added.

why: Part of support for mobileIP and wired-cum-wireless simulations.

Need to convert address in string format into int format; so use

```
Address::instance().str2addr(argv[2])  
instead of
```

```
atoi(argv[2])
```

why: This is required for supporting hier-addressing/routing.

The array packet\_names[] has changed to packet\_info.name()

why: In order to remove a bunch of #defines for pkt types, an enumeration called packet\_t now describes all packet types in ns. class p\_info was created that now describes an array name\_ that has replaced packet\_names array used previously.

Have to explicitly set direction of new pkts to DOWN before sending them down to the LL.

why: A variable direction\_ in hsr\_cmh is now used. This is used in the lower layers like LL, mac, phy etc to determine the direction of the pkt flow. All incoming pkts are marked as UP by channel, which should be remarked as DOWN by agents before sending them out into the network again.

Instead of logtarget->buffer, should now call logtarget->pt->buffer.

why: This change reflects support for eventtracing. Tracing has evolved into two types, packet tracing and event tracing. Class Trace essentially supports packet tracing. However in addition to the basic tracing properties that it derives from a BaseTrace class, pkt-tracing also requires to inherit some of the Connector class properties as well. Hence pt\_, a basetrace object represents the pure tracing functionalities required for a trace object.

The parameter used to describe the reason a pkt was dropped used to be an integer. This was changed to char\*. Hence needed to define different pkt-drop reasons in string formats.

Why: Allows greater expandability and flexibility.

linkHead changed to dsrLinkHead.

why: name clashed with linkHead used elsewhere in ns.

The older cmu model used an incoming\_ flag added in all pkts to figure out direction of pkt flow in the lower layers like ll, mac etc. Later this was replaced by a variable called direction\_ added in cmn\_hdr. direction value can be set to UP, DOWN or NONE. all pkts created with a DOWN dir by default.

why: Both these flags were being used which is not really reqd. so incoming\_ flag has been replaced with direction\_.

## 16.4 Commands at a glance

Following is a list of commands used in wireless simulations:

```
$ns_ node-config -addressingType <usually flat or hierarchical used for
                                wireless topologies>
                    -adhocRouting <adhoc routing protocol like DSDV, DSR,
                                TORA, AODV etc>
                    -llType <LinkLayer>
                    -macType <MAC type like Mac/802_11>
                    -propType <Propagation model like
                                Propagation/TwoRayGround>
                    -ifqType <interface queue type like
                                Queue/DropTail/PriQueue>
                    -ifqLen <interface queue length like 50>
                    -phyType <network interface type like
                                Phy/WirelessPhy>
                    -antType <antenna type like Antenna/OmniAntenna>
                    -channelType <Channel type like Channel/WirelessChannel>
                    -topoInstance <the topography instance>
                    -wiredRouting <turning wired routing ON or OFF>
                    -mobileIP <setting the flag for mobileIP ON or OFF>
                    -energyModel <EnergyModel type>
                    -initialEnergy <specified in Joules>
                    -rxPower <specified in W>
                    -txPower <specified in W>
                    -agentTrace <tracing at agent level turned ON or OFF>
                    -routerTrace <tracing at router level turned ON or OFF>
                    -macTrace <tracing at mac level turned ON or OFF>
                    -movementTrace <mobilenode movement logging turned
                                ON or OFF>
```

This command is used typically to configure for a mobilenode. For more info about this command (part of new node APIs) see chapter titled "Restructuring ns node and new Node APIs" in ns Notes and Documentation.

```
$ns_ node <optional:hier address>
```

This command is used to create a mobilenode after node configuration is done as shown in the node-config command. In case hierarchical addressing is being used, the hier address of the node needs to be passed as well.

```
$node log-movement
```

This command previously used to enable logging of mobilenode's movement has now been replaced by \$ns\_ node-config -movementTrace <ON or OFF>.

```
create-god <num_nodes>
```

This command is used to create a God instance. The number of mobilenodes is passed as argument which is used by God to create a matrix to store connectivity information of the topology.

`$topo load_flatgrid <X> <Y> <optional:res>`

This initializes the grid for the topography object. <X> and <Y> are the x-y co-ordinates for the topology and are used for sizing the grid. The grid resolution may be passed as <res>. A default value of 1 is normally used.

`$topo load_demfile <file-descriptor>`

For loading DEMFile objects into topography. See *ns/dem.cc,h* for details on DEMFiles.

`$ns_ namtrace-all-wireless <namtrace> <X> <Y>`

This command is used to initialize a namtrace file for logging node movements to be viewed in nam. The namtrace file descriptor, the X and Y co-ordinates of the wireless topology is passed as parameters with this command.

`$ns_ nam-end-wireless <stop-time>`

This command is used to tell nam the simulation stop time given by <stop-time>.

`$ns_ initial_node_pos <node> <size>`

This command defines the node initial position in nam. <size> denotes the size of node in nam. This function must be called after mobility model has been defined.

`$mobilenode random-motion <0 or 1>`

Random-motion is used to turn on random movements for the mobilenode, in which case random destinations are assigned to the node. 0 disables and 1 enables random-motion.

`$mobilenode setdest <X> <Y> <s>`

This command is used to setup a destination for the mobilenode. The mobile node starts moving towards destination given by <X> and <Y> at a speed of <s> m/s.

`$mobilenode reset`

This command is used to reset all the objects in the nodes (network components like LL, MAC, phy etc).

#### Internal procedures

Following is a list of internal procedures used in wireless networking:

`$mobilenode base-station <BSnode-hier-addr>`

This is used for wired-cum-wireless scenarios. Here the mobilenode is provided with the base-station node info for its domain. The address is hierarchical since wired-cum-wireless scenarios typically use hierarchical addressing.

`$mobilenode log-target <target-object>`

The <target-object>, which is normally a trace object, is used to log mobilenode movements and their energy usage, if energy model is provided.

`$mobilenode topography <topoinstance>`

This command is used to provide the node with a handle to the topography object.

`$mobilenode addif`

A mobilenode may have more than one network interface. This command is used to pass handle for a network interface to the node.

`$mobilenode namattach <namtracefd>`

This command is used to attach the namtrace file descriptor <namtracefd> to the mobilenode. All nam traces for the node are then written into this namtrace file.



`$mobilenode radius <r>`

The radius <r> denotes the node's range. All mobilenodes that fall within the circle of radius <r> with the node at its center are considered as neighbours. This info is typically used by the gridkeeper.

`$mobilenode start`

This command is used to start off the movement of the mobilenode.

# Chapter 17

## Satellite Networking in *ns*

This chapter describes extensions that enable the simulation of satellite networks in *ns*. In particular, these extensions enable *ns* to model the following: i) traditional geostationary “bent-pipe” satellites with multiple users per uplink/downlink and asymmetric links, ii) geostationary satellites with processing payloads (either regenerative payloads or full packet switching), and iii) polar orbiting LEO constellations such as Iridium and Teledesic. These satellite models are principally aimed at using *ns* to study networking aspects of satellite systems; in particular, MAC, link layer, routing, and transport protocols.

### 17.1 Overview of satellite models

Exact simulation of satellite networks requires a detailed modelling of radio frequency characteristics (interference, fading), protocol interactions (e.g., interactions of residual burst errors on the link with error checking codes), and second-order orbital effects (precession, gravitational anomalies, etc.). However, in order to study fundamental characteristics of satellite networks from a *networking* perspective, certain features may be abstracted out. For example, the performance of TCP over satellite links is impacted little by using an approximate rather than detailed channel model—performance can be characterized to first order by the overall packet loss probability. This is the approach taken in this simulation model—to create a framework for studying transport, routing, and MAC protocols in a satellite environment consisting of geostationary satellites or constellations of polar-orbiting low-earth-orbit (LEO) satellites. Of course, users may extend these models to provide more detail at a given layer.

#### 17.1.1 Geostationary satellites

Geostationary satellites orbit the Earth at an altitude of 22,300 miles above the equator. The position of the satellites is specified in terms of the longitude of the nadir point (subsatellite point on the Earth’s surface). In practice, geostationary satellites can drift from their designated location due to gravitational perturbations—these effects are not modelled in *ns*.

Two kinds of geostationary satellites can be modelled. Traditional “bent-pipe” geostationary satellites are merely repeaters in orbit— all packets received by such satellites on an uplink channel are piped through at RF frequencies to a corresponding downlink, and the satellite node is not visible to routing protocols. Newer satellites will increasingly use baseband processing, both to regenerate the digital signal and to perform fast packet switching on-board the spacecraft. In the simulations, these satellites can be modelled more like traditional *ns* nodes with classifiers and routing agents.

Previously, users could simulate geostationary satellite links by simply simulating a long delay link using traditional *ns* links and nodes. The key enhancement of these satellite extensions with respect to geostationary satellites is the capability to

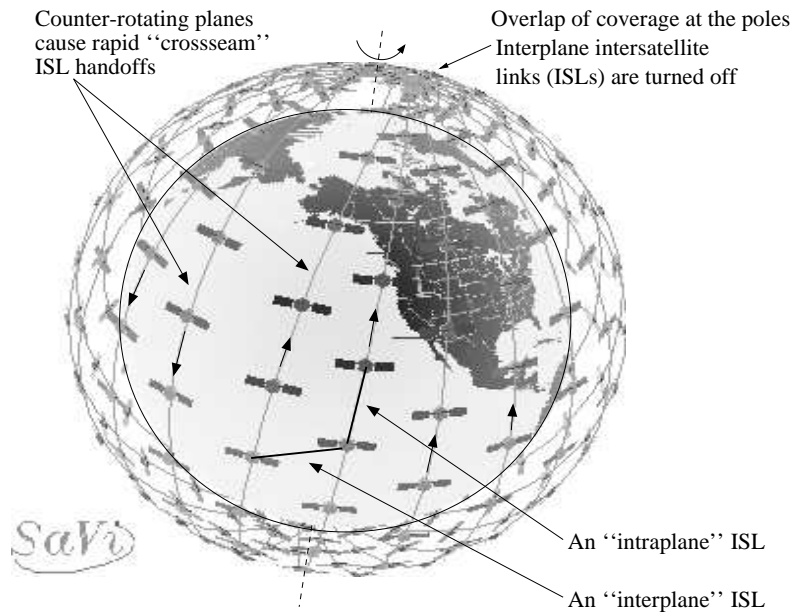


Figure 17.1: Example of a polar-orbiting LEO constellation. This figure was generated using the SaVi software package from the geometry center at the University of Minnesota.

simulate MAC protocols. Users can now define many terminals at different locations on the Earth’s surface and connect them to the same satellite uplink and downlink channels, and the propagation delays in the system (which are slightly different for each user) are accurately modelled. In addition, the uplink and downlink channels can be defined differently (perhaps with different bandwidths or error models).

### 17.1.2 Low-earth-orbiting satellites

Polar orbiting satellite systems, such as Iridium and the proposed Teledesic system, can be modelled in *ns*. In particular, the simulator supports the specification of satellites that orbit in purely circular planes, for which the neighboring planes are co-rotating. There are other non-geostationary constellation configurations possible (e.g., Walker constellations)—the interested user may develop new constellation classes to simulate these other constellation types. In particular, this would mainly require defining new intersatellite link handoff procedures.

The following are the parameters of satellite constellations that can currently be simulated:

- **Basic constellation definition** Includes satellite altitude, number of satellites, number of planes, number of satellites per plane.
- **Orbits** Orbit inclination can range continuously from 0 to 180 degrees (inclination greater than 90 degrees corresponds to retrograde orbits). Orbit eccentricity is not modeled. Nodal precession is not modeled. Intersatellite spacing within a given plane is fixed. Relative phasing between planes is fixed (although some systems may not control phasing between planes).
- **Intersatellite (ISL) links** For polar orbiting constellations, intraplane, interplane, and crossseam ISLs can be defined. Intraplane ISLs exist between satellites in the same plane and are never deactivated or handed off. Interplane ISLs exist between satellites of neighboring co-rotating planes. These links are deactivated near the poles (above the “ISL latitude threshold” in the table) because the antenna pointing mechanism cannot track these links in the polar regions. Like intraplane ISLs, interplane ISLs are never handed off. Crossseam ISLs may exist in a constellation between satellites

in counter-rotating planes (where the planes form a so-called “seam” in the topology). GEO ISLs can also be defined for constellations of geostationary satellites.

- **Ground to satellite (GSL) links** Multiple terminals can be connected to a single GSL satellite channel. GSL links for GEO satellites are static, while GSL links for LEO channels are periodically handed off as described below.
- **Elevation mask** The elevation angle above which a GSL link can be operational. Currently, if the (LEO) satellite serving a terminal drops below the elevation mask, the terminal searches for a new satellite above the elevation mask. Satellite terminals check for handoff opportunities according to a timeout interval specified by the user. Each terminal initiates handoffs asynchronously; it would be possible also to define a system in which each handoff occurs synchronously in the system.

The following table lists parameters used for example simulation scripts of the Iridium<sup>1</sup> and Teledesic<sup>2</sup> systems.

	<b>Iridium</b>	<b>Teledesic</b>
<b>Altitude</b>	780 km	1375 km
<b>Planes</b>	6	12
<b>Satellites per plane</b>	11	24
<b>Inclination (deg)</b>	86.4	84.7
<b>Interplane separation (deg)</b>	31.6	15
<b>Seam separation (deg)</b>	22	15
<b>Elevation mask (deg)</b>	8.2	40
<b>Intraplane phasing</b>	yes	yes
<b>Interplane phasing</b>	yes	no
<b>ISLs per satellite</b>	4	8
<b>ISL bandwidth</b>	25 Mb/s	155 Mb/s
<b>Up/downlink bandwidth</b>	1.5 Mb/s	1.5 Mb/s
<b>Cross-seam ISLs</b>	no	yes
<b>ISL latitude threshold (deg)</b>	60	60

Table 17.1: Simulation parameters used for modeling a broadband version of the Iridium system and the proposed 288-satellite Teledesic system. Both systems are examples of polar orbiting constellations.

<sup>1</sup>Aside from the link bandwidths (Iridium is a narrowband system only), these parameters are very close to what a broadband version of the Iridium system might look like.

<sup>2</sup>These Teledesic constellation parameters are subject to change; thanks to Marie-Jose Montpetit of Teledesic for providing tentative parameters as of January 1999. The link bandwidths are not necessarily accurate.

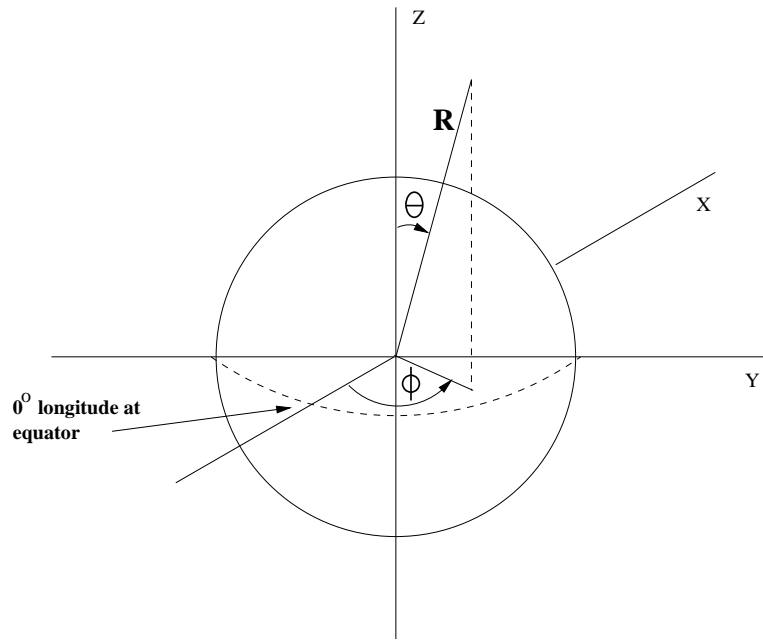


Figure 17.2: Spherical coordinate system used by satellite nodes

## 17.2 Using the satellite extensions

### 17.2.1 Nodes and node positions

There are two basic kinds of satellite nodes: *geostationary* and *non-geostationary* satellite nodes. In addition, *terminal* nodes can be placed on the Earth's surface. As is explained later in Section 17.3, each of these three different types of nodes is actually implemented with the same class `SatNode` object, but with different position, handoff manager, and link objects attached. The position object keeps track of the satellite node's location in the coordinate system as a function of the elapsed simulation time. This position information is used to determine link propagation delays and appropriate times for link handoffs. Section 5.3 introduced the "node-config" utility used to prime the node generator for different types of satellite nodes.

Figure 17.2 illustrates the spherical coordinate system, and the corresponding Cartesian coordinate system. The coordinate system is centered at the Earth's center, and the  $z$  axis coincides with the Earth's axis of rotation.  $(R, \theta, \phi) = (6378km, 90^\circ, 0^\circ)$  corresponds to  $0^\circ$  longitude (prime meridian) on the equator.

Specifically, there is one class of satellite node `Class Node/SatNode`, to which one of three types of `Position` objects may be attached. Each `SatNode` and `Position` object is a split OTcl/C++ object, but most of the code resides in C++. The following types of position objects exist:

- `Position/Sat/Term` A terminal is specified by its latitude and longitude. Latitude ranges from  $[-90, 90]$  and longitude ranges from  $[-180, 180]$ , with negative values corresponding to south and west, respectively. As simulation time evolves, the terminals move along with the Earth's surface. The node generator can be used to create a terminal with an attached position object as follows:

```
$ns node-config -satNodeType terminal \
(other node config commands go here...)
set nl [$ns node]
$nl set-position $lat $lon; # in decimal degrees
```

- **Position/Sat/Geo** A geostationary satellite is specified by its longitude above the equator. As simulation time evolves, the geostationary satellite moves through the coordinate system with the same orbital period as that of the Earth's rotation. The longitude ranges from  $[-180, 180]$  degrees. As we describe further below, two flavors of geostationary nodes exist: “geo” (for processing satellites) and “geo-repeater” (for bent-pipe satellites). The node generator can be used to create a geostationary satellite with an attached position object as follows:

```
$ns node-config -satNodeType geo (or ``geo-repeater``) \  
(other node config commands go here...)\  
set nl [$ns node]  
$nl set-position $lon; # in decimal degrees
```

- **Position/Sat/Polar** A polar orbiting satellite has a purely circular orbit along a fixed plane in the coordinate system; the Earth rotates underneath this orbital plane, so there is both an east-west and a north-south component to the track of a polar satellite's footprint on the Earth's surface. Strictly speaking, the polar position object can be used to model the movement of any circular orbit in a fixed plane; we use the term “polar” here because we later use such satellites to model polar-orbiting constellations.

Satellite orbits are usually specified by six parameters: *altitude*, *semi-major axis*, *eccentricity*, *right ascension of ascending node*, *inclination*, and *time of perigee passage*. The polar orbiting satellites in *ns* have purely circular orbits, so we simplify the specification of the orbits to include only three parameters: *altitude*, *inclination*, and *longitude*, with a fourth parameter *alpha* specifying initial position of the satellite in the orbit, as described below. **Altitude** is specified in kilometers above the Earth's surface, and **inclination** can range from  $[0, 180]$  degrees, with 90 corresponding to pure polar orbits and angles greater than 90 degrees corresponding to “retrograde” orbits. The *ascending node* refers to the point where the footprint of the satellite orbital track crosses the equator moving from south to north. In this simulation model, the parameter **longitude of ascending node** specifies the earth-centric longitude at which the satellite's nadir point crosses the equator moving south to north.<sup>3</sup> *Longitude of ascending node* can range from  $[-180, 180]$  degrees. The fourth parameter, **alpha**, specifies the initial position of the satellite along this orbit, starting from the ascending node. For example, an *alpha* of 180 degrees indicates that the satellite is initially above the equator moving from north to south. *Alpha* can range from  $[0, 360]$  degrees. Finally, a fifth parameter, **plane**, is specified when creating polar satellite nodes— all satellites in the same plane are given the same plane index. The node generator used to create a polar satellite with an attached position object as follows:

```
$ns node-config -satNodeType polar \  
(other node config commands go here...)\  
set nl [$ns node]  
$nl set-position $alt $inc $lon $alpha $plane
```

## 17.2.2 Satellite links

Satellite links resemble wireless links, which are described in Chapter 16. Each satellite node has one or more satellite network interface stacks, to which channels are connected to the physical layer object in the stack. Figure 17.3 illustrates the major components. Satellite links differ from *ns* wireless links in two major respects: i) the transmit and receive interfaces must be connected to different channels, and ii) there is no ARP implementation. Currently, the *Radio Propagation Model* is a placeholder for users to add more detailed error models if so desired; the current code does not use a propagation model.

Network interfaces can be added with the following instproc of `Class Node/SatNode`:

```
$node add-interface $type $ll $qtype $qlim $mac $mac_bw $phy
```

<sup>3</sup>Traditionally, the “right ascension” of the ascending node is specified for satellite orbits— the right ascension corresponds to the *celestial* longitude. In our case, we do not care about the orientation in a celestial coordinate system, so we specify the earth-centric longitude instead.

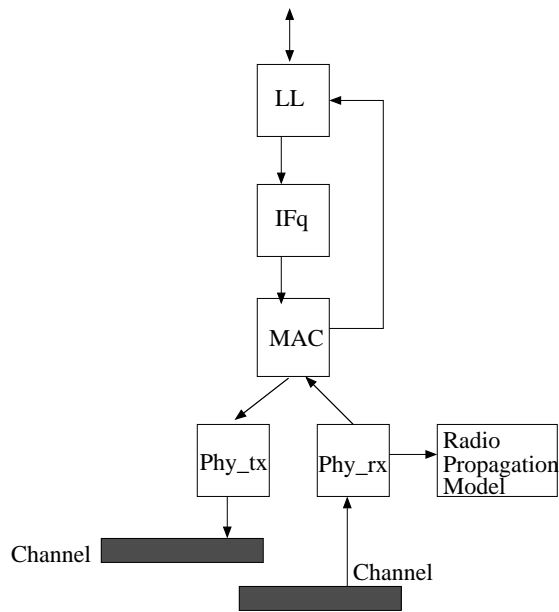


Figure 17.3: Main components of a satellite network interface

The `add-interface` instproc returns an index value that can be used to access the network interface stack later in the simulation. By convention, the first interface created on a node is attached to the uplink and downlink channels of a satellite or terminal. The following parameters must be provided:

- **type:** The following link types can be indicated: `geo` or `polar` for links from a terminal to a geo or polar satellite, respectively, `gsl` and `gsl-repeater` for links from a satellite to a terminal, and `intraplane`, `interplane`, and `crossseam` ISLs. The `type` field is used internally in the simulator to identify the different types of links, but structurally they are all very similar.
- **ll:** The link layer type (`class LL/Sat` is currently the only one defined).
- **qtype:** The queue type (e.g., `class Queue/DropTail`). Any queue type may be used— however, if additional parameters beyond the length of the queue are needed, then this instproc may need to be modified to include more arguments.
- **qlim:** The length of the interface queue, in packets.
- **mac:** The MAC type. Currently, two types are defined: `class Mac/Sat`— a basic MAC for links with only one receiver (i.e., it does not do collision detection), and `Class Mac/Sat/UnslottedAloha`— an implementation of unslotted Aloha.
- **mac\_bw:** The bandwidth of the link is set by this parameter, which controls the transmission time how fast the MAC sends. The packet size used to calculate the transmission time is the sum of the values `size()` in the common packet header and `LINK_HDRSIZE`, which is the size of any link layer headers. The default value for `LINK_HDRSIZE` is 16 bytes (settable in `satlink.h`). The transmission time is encoded in the packet header for use at the receive MAC (to simulate waiting for a whole packet to arrive).
- **phy:** The physical layer— currently two `Phys` (`Class Phy/Sat` and `Class Phy/Repeater`) are defined. The `class Phy/Sat` just pass the information up and down the stack— as in the wireless code described in Chapter 16, a radio propagation model could be attached at this point. The `class Phy/Repeater` pipes any packets received on a receive interface straight through to a transmit interface.

An ISL can be added between two nodes using the following instproc:

```
$ns add-isl $ltype $node1 $node2 $bw $qtype $qlim
```

This creates two channels (of type Channel/Sat), and appropriate network interfaces on both nodes, and attaches the channels to the network interfaces. The bandwidth of the link is set to `bw`. The linktype (`ltype`) must be specified as either `intraplane`, `interplane`, or `crossseam`.

A GSL involves adding network interfaces and a channel on board the satellite (this is typically done using the wrapper methods described in the next paragraph), and then defining the correct interfaces on the terrestrial node and attaching them to the satellite link, as follows:

```
$node add-gsl $type $ll $qtype $qlim $mac $bw_up $phy \
    [$node_satellite set downlink_] [$node_satellite set uplink_]
```

Here, the `type` must be either `geo` or `polar`, and we make use of the `downlink_` and `uplink_` instvars of the satellite; therefore, the satellite's uplink and downlink must be created before this instproc is called.

By default, the node generator for satellite nodes (described in Section 5.3) will create nodes of a given type, give them an uplink and downlink interface, and create and attach an (initial) uplink and downlink channel, based on the interface options specified.

### 17.2.3 Handoffs

Satellite handoff modelling is a key component of LEO satellite network simulations. It is difficult to predict exactly how handoffs will occur in future LEO systems because the subject is not well treated in the literature. In these satellite extensions, we establish certain criteria for handoffs, and allow nodes to independently monitor for situations that require a handoff. An alternative would be to have all handoff events synchronized across the entire simulation—it would not be difficult to change the simulator to work in such a manner.

There are no link handoffs involving geostationary satellites, but there are two types of links to polar orbiting satellites that must be handed off: GSLs to polar satellites, and crossseam ISLs. A third type of link, interplane ISLs, are not handed off but are deactivated at high latitudes as we describe below.

Each terminal connected to a polar orbiting satellite runs a timer that, upon expiry, causes the `HandoffManager` to check whether the current satellite has fallen below the elevation mask of the terminal. If so, the handoff manager detaches the terminal from that satellite's up and down links, and searches through the linked list of satellite nodes for another possible satellite. First, the "next" satellite in the current orbital plane is checked—a pointer to this satellite is stored in the `Position` object of each polar satellite node and is set during simulation configuration using the `Node/SatNode` instproc "`$node set_next $next_node`." If the next satellite is not suitable, the handoff manager searches through the remaining satellites. If it finds a suitable polar satellite, it connects its network interfaces to that satellite's uplink and downlink channels, and restarts the handoff timer. If it does not find a suitable satellite, it restarts the timer and tries again later. If any link changes occur, the routing agent is notified.

The elevation mask and handoff timer interval are settable via OTcl:

```
HandoffManager/Term set elevation_mask_ 10; # degrees
HandoffManager/Term set term_handoff_int_ 10; # seconds
```

In addition, handoffs may be randomized to avoid phase effects by setting the following variable:

```
HandoffManager set handoff_randomization_ 0; # 0 is false, 1 is true
```



If `handoff_randomization_` is true, then the next handoff interval is a random variate picked from a uniform distribution across  $(0.5 * \text{term\_handoff\_int\_}, 1.5 * \text{term\_handoff\_int\_})$ .

Crossseam ISLs are the only type of ISLs that are handed off. The criteria for handing off a crossseam ISL is whether or not there exists a satellite in the neighboring plane that is closer to the given satellite than the one to which it is currently connected. Again, a handoff timer running within the handoff manager on the polar satellite determines when the constellation is checked for handoff opportunities. Crossseam ISL handoffs are initiated by satellites in the lower-numbered plane of the two. It is therefore possible for a transient condition to arise in which a polar satellite has two crossseam ISLs (to different satellites). The satellite handoff interval is again settable from OTcl and may also be randomized:

```
HandoffManager/Sat set sat_handoff_int_ 10; # seconds
```

Interplane and crossseam ISLs are deactivated near the poles, because the pointing requirements for the links are too severe as the satellite draw close to one another. Shutdown of these links is governed by a parameter:

```
HandoffManager/Sat set latitude_threshold_ 70; # degrees
```

The values for this parameter in the example scripts are speculative; the exact value is dependent upon the satellite hardware. The handoff manager checks the latitude of itself and its peer satellite upon a handoff timeout; if either or both of the satellites is above `latitude_threshold_` degrees latitude (north or south), the link is deactivated until both satellites drop below this threshold.

Finally, if crossseam ISLs exist, there are certain situations in which the satellites draw too close to one another in the mid-latitudes (if the orbits are not close to being pure polar orbits). We check for the occurrence of this orbital overlap with the following parameter:

```
HandoffManager/Sat set longitude_threshold_ 10; # degrees
```

Again, the values for this parameter in the example scripts are speculative. If the two satellites are closer together in longitude than `longitude_threshold_` degrees, the link between them is deactivated. This parameter is disabled (set to 0) by default— all defaults for satellite-related bound variables can be found in `~ns/tcl/lib/ns-sat.tcl`.

## 17.2.4 Routing

The current status of routing is that it is incomplete. Ideally, one should be able to run all existing *ns* routing protocols over satellite links. However, many of the existing routing protocols implemented in OTcl require that the conventional *ns* links be used. Contributions in this area are welcome, but unfortunately it is not a trivial change.

With that being said, the current routing implementation is similar to Session routing described in Chapter 28, except that it is implemented entirely in C++. Upon each topology change, a centralized routing genie determines the global network topology, computes new routes for all nodes, and uses the routes to build a forwarding table on each node. Currently, the slot table is kept by a routing agent on each node, and packets not destined for agents on the node are sent by default to this routing agent. For each destination for which the node has a route, the forwarding table contains a pointer to the head of the corresponding outgoing link. As noted in Chapter 28, the user is cautioned that this type of centralized routing can lead to minor causality violations.

The routing genie is a `class` `SatRouteObject` and is created and invoked with the following OTcl commands:

```
set satrouteobject_ [new SatRouteObject]
```

```
$satrouteobject_ compute_routes
```

where the call to `compute_routes` is performed after all of the links and nodes in the simulator have been instantiated. Like the `Scheduler`, there is one instance of a `SatRouteObject` in the simulation, and it is accessed by means of an instance variable in C++. For example, the call to recompute routes after a topology change is:

```
SatRouteObject::instance().recompute();
```

Despite the current use of centralized routing, the design of having a routing agent on each node was mainly done with distributed routing in mind. Routing packets can be sent to port 255 of each node. The key to distributed routing working correctly is for the routing agent to be able to determine from which link a packet arrived. This is accomplished by the inclusion of a `class NetworkInterface` object in each link, which uniquely labels the link on which the packet arrived. A helper function `NsObject* intf_to_target(int label)` can be used to return the head of the link corresponding to a given label. The use of routing agents parallels that of the mobility extensions, and the interested reader can turn to those examples to see how to implement distributed routing protocols in this framework.

The shortest-path route computations use the current propagation delay of a link as the cost metric. It is possible to compute routes using only the hop count and not the propagation delays; in order to do so, set the following default variable to "false":

```
SatRouteObject set metric_delay_ "true"
```

Finally, for very large topologies (such as the Teledesic example), the centralized routing code will produce a very slow runtime because it executes an all-pairs shortest path algorithm upon each topology change even if there is no data currently being sent. To speed up simulations in which there is not much data transfer but there are lots of satellites and ISLs, one can disable *handoff-driven* and enable *data-driven* route computations. With data-driven computations, routes are computed only when there is a packet to send, and furthermore, a single-source shortest-path algorithm (only for the node with a packet to send) is executed instead of an all-pairs shortest path algorithm. The following OTcl variable can configure this option (which is set to "false" by default):

```
SatRouteObject set data_driven_computation_ "false"
```

### 17.2.5 Trace support

Tracefiles using satellite nodes and links are very similar to conventional *ns* tracing described in Chapter 25. Special `SatTrace` objects (`class SatTrace` derives from `class Trace`) are used to log the geographic latitude and longitude of the node logging the trace (in the case of a satellite node, the latitude and longitude correspond to the nadir point of the satellite).

For example, a packet on a link from node 66 to node 26 might normally be logged as:

```
+ 1.0000 66 26 cbr 210 ----- 0 66.0 67.0 0 0
```

but in the satellite simulation, the position information is appended:

```
+ 1.0000 66 26 cbr 210 ----- 0 66.0 67.0 0 0 37.90 -122.30 48.90 -120.94
```

In this case, node 66 is at latitude 37.90 degrees, longitude -122.30 degrees, while node 26 is a LEO satellite whose subsatellite point is at 48.90 degrees latitude, -120.94 degrees longitude (negative latitude corresponds to south, while negative longitude corresponds to west).

One addition is the Class `Trace/Sat/Error`, which traces any packets that are errored by an error model. The error trace logs packets dropped due to errors as follows, for example:

```
e 1.2404 12 13 cbr 210 ----- 0 12.0 13.0 0 0 -0.00 10.20 -0.00 -10.00
```

It may happen that a satellite node generates a packet that it cannot forward (such as in `sat-mixed.tcl`). This will show up as a drop in the trace file with a destination field set to -2, and the coordinates set to -999.00:

```
d 848.0000 14 -2 cbr 210 ----- 1 14.0 15.0 6 21 0.00 10.00 -999.00 -999.00
```

This indicates that node 14, in trying to send a packet to node 15, could not find an available route.

To enable tracing of all satellite links in the simulator, use the following commands *before* instantiating nodes and links:

```
set f [open out.trace w]
$ns trace-all $f
```

Then use the following line after all node and link creation (and all error model insertion, if any) to enable tracing of all satellite links:

```
$ns trace-all-satlinks $f
```

Specifically, this will put tracing around the link layer queues in all satellite links, and will put a receive trace between the mac and the link layer for received packets. To enable tracing only on a specific link on a specific node, one may use the command:

```
$node trace-inlink-queue $f $i
$node trace-outlink-queue $f $i
```

where *i* is the index of the interface to be traced.

The implementations of the satellite trace objects can be found in `~ns/tcl/lib/ns-sat.tcl` and `~ns/sattrace.{cc,h}`.

## 17.2.6 Error models

*ns* error models are described in Chapter 13. These error models can be set to cause packets to be errored according to various probability distributions. These error models are simple and don't necessarily correspond to what would be experienced on an actual satellite channel (particularly a LEO channel). Users are free to define more sophisticated error models that more closely match a particular satellite environment.

The following code provides an example of how to add an error model to a link:

```
set em_ [new ErrorModel]
$em_ unit pkt
$em_ set rate_ 0.02
$em_ ranvar [new RandomVariable/Uniform]
$node interface-errormodel $em_
```

This will add an error model to the receive path of the first interface created on node `$node` (specifically, between the MAC and link layer)— this first interface generally corresponds to the uplink and downlink interface for a satellite or a terminal (if only one uplink and/or downlink exists). To add the error model to a different stack (indexed by  $i$ ), use the following code:

```
$node interface-errormodel $em_ $i
```

### 17.2.7 Other configuration options

Given an initial configuration of satellites specified for time 0, it is possible to start the satellite configuration from any arbitrary point in time through the use of the `time_advance_` parameter (this is really only useful for LEO simulations). During the simulation run, this will set the position of the object to the position at time `Scheduler::instance().clock + time_advance_seconds`.

```
Position/Sat set time_advance_ 0; # seconds
```

### 17.2.8 *nam* support

*nam* is not currently supported. Addition of *nam* for satellite is open to interested contributors.

### 17.2.9 Integration with wired and wireless code

Recently (November 2001), support has been added to connect traditional OTcl-based wired nodes with the satellite nodes. This section describes the capabilities and limitations of that code.

The satellite code (and the wireless code) normally performs all routing in C++, while the traditional ns code uses a mix of OTcl and C++ code. For backward compatibility reasons, it is difficult to fully integrate both the wired and wireless code. The strategy for integrating wireless and wired code has been to define a special gateway node (called a "basestation"), to use hierarchical routing, and to locate a single basestation node in the wireless network with a network stack located in both the wireless and the wired subnet. Because routing is not fully integrated, the topology of the simulation is limited to only one gateway node per wireless subnet (i.e., a packet cannot enter the wireless network from one wired gateway and leave via another).

The satellite/wired code integration takes a different strategy. By selecting the node configuration `$ns node-config -wiredRouting ON` option, the C++ routing in the satellite code is turned off, and instead, all satellite topology changes lead to upcalls into the OTcl code. As a result, the `link_` array in OTcl is manipulated according to all topology changes, and OTcl-based routing can occur. The penalty for doing this is a much longer execution time for larger simulations (such as Teledesic), but for smaller simulations, the difference is not as noticeable.

An example script detailing the use of this new option is shown in `~ns/tcl/ex/sat-wired.tcl`, and a similar test in the satellite test suite exercises this code. Additionally, all of the satellite example scripts in `~ns/tcl/ex` directory can be converted to OTcl routing by using the `$ns node-config -wiredRouting ON` option. However, there are a few caveats:

- The wired routing option for satellite has only been tested with (the default) static routing: `$ns rtProto Static`. The code triggers a global routing table update upon any satellite topology change.
- The option `data_driven_computation_` can not be set to "true" when `wiredRouting` is ON. Note that the enabling or disabling of `data_driven_computation_` can give subtle differences in simulation output since routes

are computed at different times (while propagation delays are continuously changing). This effect can be seen by toggling this parameter in the Iridium example script `~ns/tcl/ex/sat-iridium.tcl`.

- In the trace file, when a packet is dropped due to “no route to host” (such as when there is a topology change), the trace looks a bit different depending on whether `wiredRouting` is turned OFF or ON. In the former case, there is one line per drop, with the destination labelled as “-2”. In the latter case, there are three events (enqueue “+”, deque “-”, and drop “d”) corresponding to the same packet, and the destination is shown as “-1”.
- In rare cases, there may be warning messages during the execution indicating “node out of range.” This can occur if a node becomes disconnected in the topology and then another node tries to send a packet to it. For example, try enabling `wiredRouting` in the file `~ns/tcl/ex/sat-mixed.tcl`. This occurs because the routing table is dynamically sized upon topology change, and if a node becomes disconnected it may not have any entries inserted in the routing table (and hence the routing table is not grown to accommodate its node number). This warning should not affect actual trace output.
- There has been no attempt to interoperate with wireless or mobile-IP code.

### 17.2.10 Example scripts

Example scripts can be found in the `~ns/tcl/ex` directory, including:

- `sat-mixed.tcl` A simulation with a mixture of polar and geostationary satellites.
- `sat-wired.tcl` Similar to the previous script, but shows how to connect wired nodes to a satellite simulation.
- `sat-repeater.tcl` Demonstrates the use of a simple bent-pipe geostationary satellite, and also error models.
- `sat-aloha.tcl` Simulates one hundred terminals in a mesh-VSAT configuration using an unslotted Aloha MAC protocol with a “bent-pipe” geostationary satellite. Terminals listen to their own transmissions (after a delay), and if they do not successfully receive their own packet within a timeout interval, they perform exponential backoff and then retransmit the packet. Three variants exist: `basic`, `basic_tracing`, and `poisson`. These variants are described further in the header comments of the script.
- `sat-iridium.tcl` Simulates a broadband LEO constellation with parameters similar to that of the Iridium constellation (with supporting scripts `sat-iridium-links.tcl`, `sat-iridium-linkswithcross.tcl`, and `sat-iridium-nodes.tcl`).
- `sat-teledesic.tcl` Simulates a broadband LEO constellation with parameters similar to those proposed for the 288 satellite Teledesic constellation (with supporting scripts `sat-teledesic-links.tcl` and `sat-teledesic-nodes.tcl`).

In addition, there is a test suite script that tries to exercise a lot of features simultaneously, it can be found at `~ns/tcl/test/test-suite-sat.tcl`.

## 17.3 Implementation

The code for the implementation of satellite extensions can be found in `~ns/{sat.h, sathandoff.{cc,h}, satlink.{cc,h}, satnode.{cc,h}, satposition.{cc,h}, satroute.{cc,h}, sattrace.{cc,h}}`, and `~ns/tcl/lib/ns-sat.tcl`. Almost all of the mechanism is implemented in C++.

In this section, we focus on some of the key components of the implementation; namely, the use of linked lists, the node structure, and a detailed look at the satellite link structure.

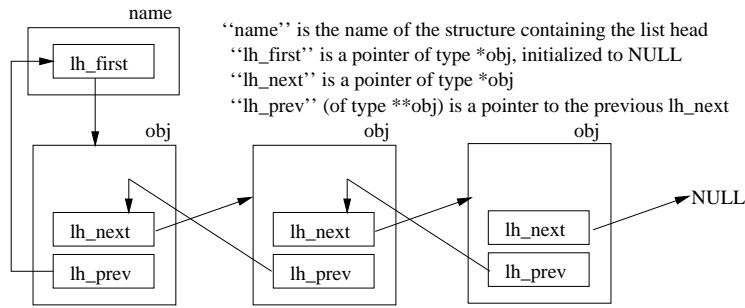


Figure 17.4: Linked list implementation in *ns*.

### 17.3.1 Use of linked lists

There are a number of linked lists used heavily in the implementation:

- `class Node` maintains a (static) list of all objects of class `Node` in the simulator. The variable `Node::nodehead_` stores the head of the list. The linked list of nodes is used for centralized routing, for finding satellites to hand off to, and for tracing.
- `class Node` maintains a list of all (satellite) links on the node. Specifically, the list is a list of objects of class `LinkHead`. The variable `linklisthead_` stores the head of the list. The linked list of `LinkHeads` is used for checking whether or not to handoff links, and to discover topology adjacencies.
- `class Channel` maintains a list of all objects of class `Phy` on the channel. The head of the list is stored in the variable `if_head_`. This list is used to determine the set of interfaces on a channel that should receive a copy of a packet.

Figure 17.4 provides a schematic of how the linked list is organized. Each object in the list is linked through a “LINK\_ENTRY” that is a protected member of the class. This entry contains a pointer to the next item in the list and also a pointer to the address of the previous “next” pointer in the preceding object. Various macros found in `~ns/list.h` can be used to manipulate the list; the implementation of linked-lists in *ns* is similar to the `queue` implementation found in some variants of BSD UNIX.

### 17.3.2 Node structure

Figure 17.5 is a schematic of the main components of a `SatNode`. The structure bears resemblance to the `MobileNode` in the wireless extensions, but there are several differences. Like all *ns* nodes, the `SatNode` has an “entry” point to a series of classifiers. The address classifier contains a slot table for forwarding packets to foreign nodes, but since OTcl routing is not used, all packets not destined for this node (and hence forwarded to the port classifier), are sent to the default target, which points to a routing agent. Packets destined on the node for port 255 are classified as routing packets and are also forwarded to the routing agent.

Each node contains one or more “network stacks” that include a generic `SatLinkHead` at the entry point of the link. The `SatLinkHead` is intended to serve as an API to get at other objects in the link structure, so it contains a number of pointers (although the API here has not been finalized). Packets leaving the network stack are sent to the node’s entry. An important feature is that each packet leaving a network stack has its `iface_` field in the common packet header coded with the unique `NetworkInterface` index corresponding to the link. This value can be used to support distributed routing as described below.

The base class routing agent is `class SatRouteAgent`; it can be used in conjunction with centralized routing. `SatRouteAgents` contain a forwarding table that resolves a packet’s address to a particular `LinkHead` target—it is the job of the `SatRouteObject`

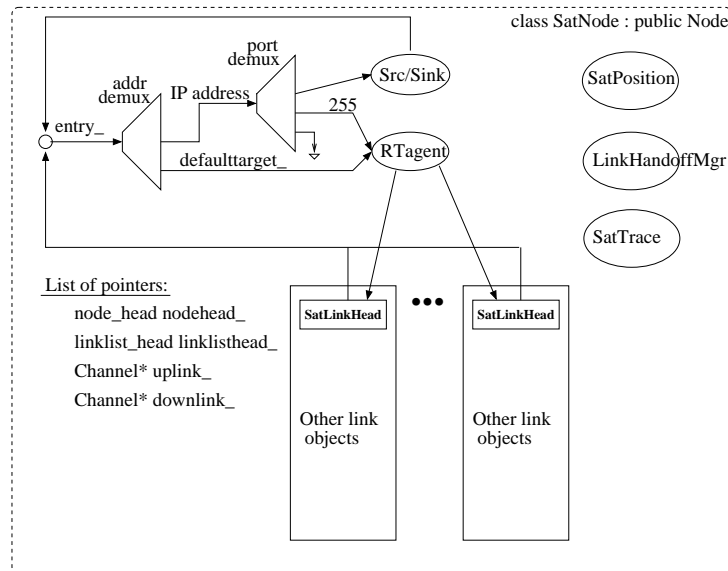


Figure 17.5: Structure of class `SatNode`.

to populate this table correctly. The `SatRouteAgent` populates certain fields in the header and then sends the packet down to the appropriate link. To implement a distributed routing protocol, a new `SatRouteAgent` could be defined—this would learn about topology by noting the interface index marked in each packet as it came up the stack—a helper function in the node `intf_to_target()` allows it to resolve an index value to a particular `LinkHead`.

There are pointers to three additional objects in a `SatNode`. First, each `SatNode` contains a position object, discussed in the previous section. Second, each `SatNode` contains a `LinkHandoffMgr` that monitors for opportunities to hand links off and coordinates the handoffs. Satellite nodes and terminal nodes each have their specialized version of a `LinkHandoffMgr`.

Finally, a number of pointers to objects are contained in a `SatNode`. We discussed `linklisthead_` and `nodehead_` in the previous subsection. The `uplink_` and `downlink_` pointers are used for convenience under the assumption that, in most simulations, a satellite or a terminal has only one uplink and downlink channel.

### 17.3.3 Detailed look at satellite links

Figure 17.6 provides a more detailed look at how satellite links are composed. In this section, we describe how packets move up and down the stack, and the key things to note at each layer. The file `~ns/tcl/lib/ns-sat.tcl` contains the various OTcl instprocs that assemble links according to Figure 17.6. We describe the composite structure herein as a “network stack.” Most of the code for the various link components is in `~ns/satlink.{cc,h}`.

The entry point to a network stack is the `SatLinkHead` object. The `SatLinkHead` object derives from `Class LinkHead`; the aim of link head objects is to provide a uniform API for all network stacks.<sup>4</sup> The `SatLinkHead` object contains pointers to the LL, Queue, MAC, Error model, and both Phy objects. The `SatLinkHead` object can also be queried as to what type of network stack it is—e.g., GSL, interplane ISL, crossseam ISL, etc.. Valid codes for the `type_` field are currently found in `~ns/sat.h`. Finally, the `SatLinkHead` stores a boolean variable `linkup_` that indicates whether the link to at least one other node on the channel is up. The C++ implementation of `SatLinkHead` is found in `~ns/satlink.{cc,h}`.

Packets leaving a node pass through the `SatLinkHead` transparently to the `class SatLL` object. The `SatLL` class derives from LL (link layer). Link layer protocols (like ARQ protocols) can be defined here. The current `SatLL` assigns a MAC

<sup>4</sup>In the author’s opinion, all network stacks in *ns* should eventually have a `LinkHead` object at the front—the class `SatLinkHead` would then disappear.

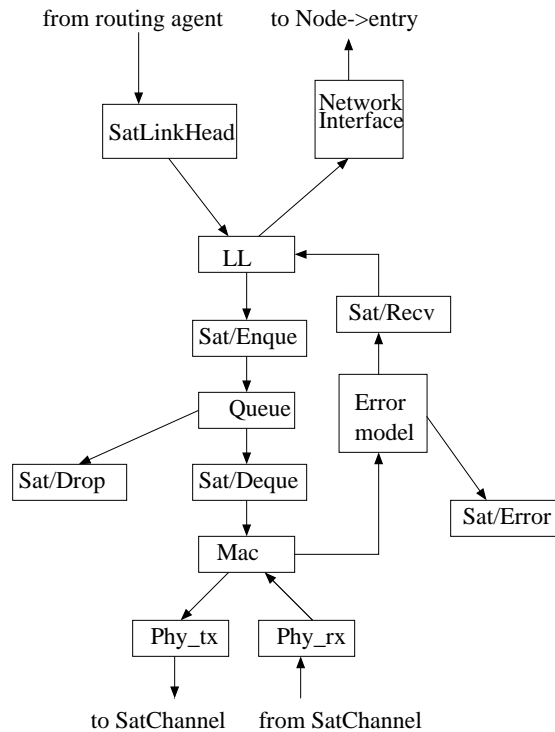


Figure 17.6: Detailed look at network interface stack.

address to the packet. Note that in the satellite case, we do not use an Address Resolution Protocol (ARP); instead, we simply use the `MAC index_` variable as its address, and we use a helper function to find the MAC address of the corresponding interface of the next-hop node. Since `class LL` derives from `class LinkDelay`, the `delay_` parameter of `LinkDelay` can be used to model any processing delay in the link layer; by default this delay is zero.

The next object an outgoing packet encounters is the interface queue. However, if tracing is enabled, tracing elements may surround the queue, as shown in Figure 17.6. This part of a satellite link functions like a conventional *ns* link.

The next layer down is the MAC layer. The MAC layer draws packets from the queue (or deque trace) object— a handshaking between the MAC and the queue allows the MAC to draw packets out of the queue as it needs them. The transmission time of a packet is modelled in the MAC also— the MAC computes the transmission delay of the packet (based on the sum of the `LINK_HDRSIZE` field defined in `satlink.h` and the `size` field in the common packet header), and does not call up for another packet until the current one has been “sent” to the next layer down. Therefore, it is important to set the bandwidth of the link correctly at this layer. For convenience, the transmit time is encoded in the mac header; this information can be used at the receiving MAC to calculate how long it must wait to detect a collision on a packet, for example.

Next, the packet is sent to a transmitting interface (`Phy_tx`) of class `SatPhy`. This object just sends the packet to the attached channel. We noted earlier in this chapter that all interfaces attached to a channel are part of the linked list for that channel. This is not true for transmit interfaces, however. Only receive interfaces attached to a channel comprise this linked list, since only receive interfaces should get a copy of transmitted packets. The use of separate transmit and receive interfaces mirrors the real world where full-duplex satellite links are made up of RF channels at different frequencies.

The outgoing packet is next sent to a `SatChannel`, which copies the packet to every receiving interface (of class `SatPhy`) on the channel. The `Phy_rx` sends the packet to the MAC layer. At the MAC layer, the packet is held for the duration of its transmission time (and any appropriate collision detection is performed if the MAC, such as the Aloha MAC, supports it). If the packet is determined to have arrived safely at the MAC, it next passes to an `ErrorModel` object, if it exists. If not, the packet moves through any receive tracing objects to the `SatLL` object. The `SatLL` object passes the packet up after a processing delay (again, by default, the value for `delay_` is zero).



The final object that a received packet passes through is an object of class `NetworkInterface`. This object stamps the `iface_` field in the common header with the network stack's unique index value. This is used to keep track of which network stack a packet arrived on. The packet then goes to the entry of the `SatNode` (usually, an address classifier).

Finally, “geo-repeater” satellites exist, as described earlier in this chapter. Geo-repeater network stacks are very simple— they only contain a `Phy_tx` and a `Phy_rx` of class `RepeaterPhy`, and a `SatLinkHead`. Packets received by a `Phy_rx` are sent to the `Phy_tx` without delay. The geo-repeater satellite is a degenerate satellite node, in that it does not contain things like tracing elements, handoff managers, routing agents, or any other link interfaces other than repeater interfaces.

## 17.4 Commands at a glance

Following is a list of commands related to satellite networking:

```
$ns_ node-config -satNodeType <type>
```

This node configuration declares that the subsequent new nodes created will be of type `<type>`, where `<type>` can be one of the following: `geo`, `geo-repeater`, `polar`, `terminal`. Other required fields for satellite nodes (for setting up initial links and channels) are as follows (see Section 5.3): `$ns_ node-config -llType <type>`

```
$ns_ node-config -ifqType <type>
```

```
$ns_ node-config -ifqLen <length>
```

```
$ns_ node-config -macType <type>
```

```
$ns_ node-config -channelType <type>
```

```
$ns_ node-config -downlinkBW <value>
```

(note— `satNodeType` `geo-repeater` only requires specifying the `channelType`— all other options are disregarded. See `tcl/ex/sat-repeater.tcl` for an example.)

```
$ns_ satnode-polar <alt> <inc> <lon> <alpha> <plane> <linkargs> <chan>
```

This is a simulator wrapper method for creating a polar satellite node. Two links, uplink and downlink, are created along with two channels, uplink channel and downlink channel. `<alt>` is the polar satellite altitude, `<inc>` is orbit inclination w.r.t equator, `<lon>` is the longitude of ascending node, `<alpha>` gives the initial position of the satellite along this orbit, `<plane>` defines the plane of the polar satellite. `<linkargs>` is a list of link argument options that defines the network interface (like `LL`, `Qtype`, `Qlim`, `PHY`, `MAC` etc).

```
$ns_ satnode-geo <lon> <linkargs> <chan>
```

This is a wrapper method for creating a geo satellite node that first creates a `satnode` plus two link interfaces (uplink and downlink) plus two satellite channels (uplink and downlink). `<chan>` defines the type of channel.

```
$ns_ satnode-geo-repeater <lon> <chan>
```

This is a wrapper method for making a geo satellite repeater node that first creates a `satnode` plus two link interfaces (uplink and downlink) plus two satellite channels (uplink and downlink).

```
$ns_ satnode-terminal <lat> <lon>
```

This is a wrapper method that simply creates a terminal node. The `<lat>` and `<lon>` defines the latitude and longitude respectively of the terminal.

```
$ns_ satnode <type> <args>
```

This is a more primitive method for creating `satnodes` of type `<type>` which can be `polar`, `geo` or `terminal`.

```
$satnode add-interface <type> <ll> <qtype> <qlim> <mac_bw> <phy>
```

This is an internal method of `Node/SatNode` that sets up link layer, mac layer, interface queue and physical layer structures for the satellite nodes.

```
$satnode add-isl <ltype> <node1> <node2> <bw> <qtype> <qlim>
```

This method creates an ISL (inter-satellite link) between the two nodes. The link type (inter, intra or cross-seam), BW of the link, the queue-type and queue-limit are all specified.

```
$satnode add-gsl <ltype> <opt_ll> <opt_ifq> <opt_qlim> <opt_mac> <opt_bw> <opt_phy>  
<opt_inlink> <opt_outlink>
```

This method creates a GSL (ground to satellite link). First a network stack is created that is defined by LL, IfQ, Qlim, MAC, BW and PHY layers. Next the node is attached to the channel inlink and outlink.

# Chapter 18

## Radio Propagation Models

This chapter describes the radio propagation models implemented in *ns*. These models are used to predict the received signal power of each packet. At the physical layer of each wireless node, there is a receiving threshold. When a packet is received, if its signal power is below the receiving threshold, it is marked as error and dropped by the MAC layer.

Up to now there are three propagation models in *ns*, which are the free space model<sup>1</sup>, two-ray ground reflection model<sup>2</sup> and the shadowing model<sup>3</sup>. Their implementation can be found in *~ns/propagation.{cc,h}*, *~ns/tworayground.{cc,h}* and *~ns/shadowing.{cc,h}*. This documentation reflects the APIs in ns-2.1b7.

### 18.1 Free space model

The free space propagation model assumes the ideal propagation condition that there is only one clear line-of-sight path between the transmitter and receiver. H. T. Friis presented the following equation to calculate the received signal power in free space at distance  $d$  from the transmitter [12].

$$P_r(d) = \frac{P_t G_t G_r \lambda^2}{(4\pi)^2 d^2 L} \quad (18.1)$$

where  $P_t$  is the transmitted signal power.  $G_t$  and  $G_r$  are the antenna gains of the transmitter and the receiver respectively.  $L(L \geq 1)$  is the system loss, and  $\lambda$  is the wavelength. It is common to select  $G_t = G_r = 1$  and  $L = 1$  in *ns* simulations.

The free space model basically represents the communication range as a circle around the transmitter. If a receiver is within the circle, it receives all packets. Otherwise, it loses all packets

The OTcl interface for utilizing a propagation model is the `node-config` command. One way to use it here is

```
$ns_ node-config -propType Propagation/FreeSpace
```

Another way is

---

<sup>1</sup>Based on the code contributed to *ns* from the CMU Monarch project.

<sup>2</sup>Contributed to *ns* from the CMU Monarch project.

<sup>3</sup>Implemented in *ns* by Wei Ye at USC/ISI

```
set prop [new Propagation/FreeSpace]
$ns_ node-config -propInstance $prop
```

## 18.2 Two-ray ground reflection model

A single line-of-sight path between two mobile nodes is seldom the only means of propagation. The two-ray ground reflection model considers both the direct path and a ground reflection path. It is shown [29] that this model gives more accurate prediction at a long distance than the free space model. The received power at distance  $d$  is predicted by

$$P_r(d) = \frac{P_t G_t G_r h_t^2 h_r^2}{d^4 L} \quad (18.2)$$

where  $h_t$  and  $h_r$  are the heights of the transmit and receive antennas respectively. Note that the original equation in [29] assumes  $L = 1$ . To be consistent with the free space model,  $L$  is added here.

The above equation shows a faster power loss than Eqn. (18.1) as distance increases. However, The two-ray model does not give a good result for a short distance due to the oscillation caused by the constructive and destructive combination of the two rays. Instead, the free space model is still used when  $d$  is small.

Therefore, a cross-over distance  $d_c$  is calculated in this model. When  $d < d_c$ , Eqn. (18.1) is used. When  $d > d_c$ , Eqn. (18.2) is used. At the cross-over distance, Eqns. (18.1) and (18.2) give the same result. So  $d_c$  can be calculated as

$$d_c = (4\pi h_t h_r) / \lambda \quad (18.3)$$

Similarly, the OTcl interface for utilizing the two-ray ground reflection model is as follows.

```
$ns_ node-config -propType Propagation/TwoRayGround
```

Alternatively, the user can use

```
set prop [new Propagation/TwoRayGround]
$ns_ node-config -propInstance $prop
```

## 18.3 Shadowing model

### 18.3.1 Background

The free space model and the two-ray model predict the received power as a deterministic function of distance. They both represent the communication range as an ideal circle. In reality, the received power at certain distance is a random variable due to multipath propagation effects, which is also known as fading effects. In fact, the above two models predicts the mean received power at distance  $d$ . A more general and widely-used model is called the shadowing model [29].

Environment		$\beta$
Outdoor	Free space	2
	Shadowed urban area	2.7 to 5
In building	Line-of-sight	1.6 to 1.8
	Obstructed	4 to 6

Table 18.1: Some typical values of path loss exponent  $\beta$

Environment	$\sigma_{dB}$ (dB)
Outdoor	4 to 12
Office, hard partition	7
Office, soft partition	9.6
Factory, line-of-sight	3 to 6
Factory, obstructed	6.8

Table 18.2: Some typical values of shadowing deviation  $\sigma_{dB}$

The shadowing model consists of two parts. The first one is known as path loss model, which also predicts the mean received power at distance  $d$ , denoted by  $\overline{P_r(d)}$ . It uses a close-in distance  $d_0$  as a reference.  $\overline{P_r(d)}$  is computed relative to  $P_r(d_0)$  as follows.

$$\frac{P_r(d_0)}{P_r(d)} = \left( \frac{d}{d_0} \right)^\beta \quad (18.4)$$

$\beta$  is called the path loss exponent, and is usually empirically determined by field measurement. From Eqn. (18.1) we know that  $\beta = 2$  for free space propagation. Table 18.1 gives some typical values of  $\beta$ . Larger values correspond to more obstructions and hence faster decrease in average received power as distance becomes larger.  $P_r(d_0)$  can be computed from Eqn. (18.1).

The path loss is usually measured in dB. So from Eqn. (18.4) we have

$$\left[ \frac{\overline{P_r(d)}}{P_r(d_0)} \right]_{dB} = -10\beta \log \left( \frac{d}{d_0} \right) \quad (18.5)$$

The second part of the shadowing model reflects the variation of the received power at certain distance. It is a log-normal random variable, that is, it is of Gaussian distribution if measured in dB. The overall shadowing model is represented by

$$\left[ \frac{P_r(d)}{P_r(d_0)} \right]_{dB} = -10\beta \log \left( \frac{d}{d_0} \right) + X_{dB} \quad (18.6)$$

where  $X_{dB}$  is a Gaussian random variable with zero mean and standard deviation  $\sigma_{dB}$ .  $\sigma_{dB}$  is called the shadowing deviation, and is also obtained by measurement. Table 18.2 shows some typical values of  $\sigma_{dB}$ . Eqn. (18.6) is also known as a log-normal shadowing model.

The shadowing model extends the ideal circle model to a richer statistic model: nodes can only probabilistically communicate when near the edge of the communication range.

### 18.3.2 Using shadowing model

Before using the model, the user should select the values of the path loss exponent  $\beta$  and the shadowing deviation  $\sigma_{dB}$  according to the simulated environment.

The OTcl interface is still the `node-config` command. One way to use it is as follows, and the values for these parameters are just examples.

```
# first set values of shadowing model
Propagation/Shadowing set pathlossExp_ 2.0    ;# path loss exponent
Propagation/Shadowing set std_db_ 4.0          ;# shadowing deviation (dB)
Propagation/Shadowing set dist0_ 1.0           ;# reference distance (m)
Propagation/Shadowing set seed_ 0              ;# seed for RNG
```

```
$ns_ node-config -propType Propagation/Shadowing
```

The shadowing model creates a random number generator (RNG) object. The RNG has three types of seeds: raw seed, pre-defined seed (a set of known good seeds) and the heuristic seed (details in Section 24.1). The above API only uses the pre-defined seed. If a user want different seeding method, the following API can be used.

```
set prop [new Propagation/Shadowing]
$prop set pathlossExp_ 2.0
$prop set std_db_ 4.0
$prop set dist0_ 1.0
$prop seed <seed-type> 0          ;# user can specify seeding method

$ns_ node-config -propInstance $prop
```

The `<seed-type>` above can be raw, predef or heuristic.

## 18.4 Communication range

In some applications, a user may want to specify the communication range of wireless nodes. This can be done by set an appropriate value of the receiving threshold in the network interface, *i.e.*,

```
Phy/WirelessPhy set RXThresh_ <value>
```

A separate C program is provided at `~ns/indep-utils/propagation/threshold.cc` to compute the receiving threshold. It can be used for all the propagation models discussed in this chapter. Assume you have compiled it and get the executable named as `threshold`. You can use it to compute the threshold as follows

```
threshold -m <propagation-model> [other-options] distance
```

where `<propagation-model>` is either `FreeSpace`, `TwoRayGround` or `Shadowing`, and the distance is the communication range in meter.

[other-options] are used to specify parameters other than their default values. For the shadowing model there is a necessary parameter, -r <receive-rate>, which specifies the rate of correct reception at the distance. Because the communication range in the shadowing model is not an ideal circle, an inverse Q-function [29] is used to calculate the receiving threshold. For example, if you want 95% of packets can be correctly received at the distance of 50m, you can compute the threshold by

```
threshold -m Shadowing -r 0.95 50
```

Other available values of [other-options] are shown below

```
-pl <path-loss-exponent> -std <shadowing-deviation> -Pt <transmit-power>  
-fr <frequency> -Gt <transmit-antenna-gain> -Gr <receive-antenna-gain>  
-L <system-loss> -ht <transmit-antenna-height> -hr <receive-antenna-height>  
-d0 <reference-distance>
```

## 18.5 Commands at a glance

Following is a list of commands for propagation models.

```
$ns_ node-config -propType <propagation-model>
```

This command selects <propagation-model> in the simulation. the <propagation model> can be Propagation/FreeSpace, Propagation/TwoRayGround or Propagation/Shadowing

```
$ns_ node-config -propInstance $prop
```

This command is another way to utilize a propagation model. \$prop is an instance of the <propagation-model>.

```
$sprop_ seed <seed-type> <value>
```

This command seeds the RNG. \$sprop\_ is an instance of the shadowing model.

```
threshold -m <propagation-model> [other-options] distance
```

This is a separate program at `~ns/indep-utils/propagation/threshold.cc`, which is used to compute the receiving threshold for a specified communication range.

# Chapter 19

## Energy Model in ns

Energy Model, as implemented in *ns*, is a node attribute. The energy model represents level of energy in a mobile host. The energy model in a node has a initial value which is the level of energy the node has at the beginning of the simulation. This is known as `initialEnergy_`. It also has a given energy usage for every packet it transmits and receives. These are called `txPower_` and `rxPower_`. The files where the energy model is defined are `ns/energymodel.cc` and `h`. Other functions/methods described in this chapter may be found in `ns/wireless-phy.cc`, `ns/cmu-trace.cc`, `ns/tcl/lib[ns-lib.tcl, ns-node.tcl, ns-mobilenode.tcl]`.

### 19.1 The C++ EnergyModel Class

The basic energy model is very simple and is defined by class `EnergyModel` as shown below:

```
class EnergyModel : public TclObject
public:
    EnergyModel(double energy)  energy_ = energy;
    inline double energy()  return energy_;
    inline void setenergy(double e)  energy_ = e;
    virtual void DecrTxEnergy(double txtime, double P_tx)
        energy_ -= (P_tx * txtime);

    virtual void DecrRcvEnergy(double rcvtime, double P_rcv)
        energy_ -= (P_rcv * rcvtime);

protected:
    double energy_;
;
```

As seen from the `EnergyModel` Class definition above, there is only a single class variable `energy_` which represents the level of energy in the node at any given time. The constructor `EnergyModel(energy)` requires the initial-energy to be passed along as a parameter. The other class methods are used to decrease the energy level of the node for every packet transmitted (`DecrTxEnergy(txtime, P_tx)`) and every packet received (`DecrRcvEnergy(rcvtime, P_rcv)`) by the node. `P_tx` and `P_rcv` are the transmitting and receiving power (respectively) required by the node's interface or PHY. At the beginning of simulation, `energy_` is set to `initialEnergy_` which is then decremented for every transmission and



reception of packets at the node. When the energy level at the node goes down to zero, no more packets can be received or transmitted by the node. If tracing is turned on, line `DEBUG: node <node-id> dropping pkts due to energy = 0` is printed in the tracefile.

## 19.2 The OTcl interface

Since the energy model is a node attribute, it may be defined by the following node configuration APIs:

```
$ns_ node-config -energyModel $energymodel \  
-rxPower $p_rx \  
-txPower $p_tx \  
-initialEnergy $initialenergy
```

Optional values for above configuration parameters of the energy model are given in the following table:

Attribute	optional values	default values
-energyModel	"EnergyModel"	none
-rxPower	receiving power in watts (e.g 0.3)	281.8mW
-txPower	transmitting power in watts (e.g 0.4)	281.8mW
-initialEnergy	energy in joules (e.g 0.1)	0.0

## Chapter 20

# Directed Diffusion

The directed diffusion module in ns has been ported from SCADDS group's implementation of directed diffusion at USC/ISI. There is an older version of diffusion in ns that was implemented several years back and has become relatively old and outdated. This older version can be found under directory diffusion. And the newer version of diffusion resides under `~ns//diffusion3`. This chapter talks about the newer diffusion model in ns. The module and methods described here can be found under `~ns/tcl/lib/ns-diffusion.tcl`, `ns-lib.tcl` and all relevant C++ code can be found under `~ns/diffusion3`. Visit the SCADDS group webpage at <http://www.isi.edu/scadds> for details about their implementation.

### 20.1 What is Directed Diffusion?

Directed Diffusion is a method of data dissemination especially suitable in distributed sensing scenarios. It differs from IP method of communication. For IP “nodes are identified by their end-points, and inter-node communication is layered on an end-to-end delivery service provided within the network”. Directed diffusion, on the other hand is data-centric. Data generated by sensor nodes are identified by their attribute-value pair. Sinks or nodes that request data send out “interest”s into the network. Data generated by “source” nodes that match these interests, “flow” towards the sinks. Intermediate nodes are capable of caching and transforming data. For details on directed diffusion, see “Directed Diffusion: A Scalable and Robust Communication Paradigm for Sensor Networks”, authored by Chalermek Intanagonwiwat, Ramesh Govindan and Deborah Estrin that appeared in MobiCOM, August 2000, Boston, Massachusetts. This and other diffusion related papers can be viewed at <http://www.isi.edu/scadds/publications.html> under publications section.

### 20.2 The diffusion model in ns

The directed diffusion model consists of a core diffusion layer, a diffusion library provides an application programming interface for overlying diffusion applications and finally the application layer which includes both diffusion applications and filters. The core diffusion layer is used to receive/send out packets from/into the network. The library provides a interface for the overlying application classes for publishing/subscribing etc. These APIs have been described in details in a document called Network Routing API 8.0 and can be found at <http://www.isi.edu/scadds/publications.html> under APIs section. In the following paragraphs we are going to describe how the diffusion model looks like in ns.

First we start with a brief description of the diffusion3 directory structure. If the reader wishes to examine the C++ code related to NS Diffusion that underpins the OTcl script commands, it may be found in `~ns/ns/diffusion3`. Here is a summary by subdirectory:

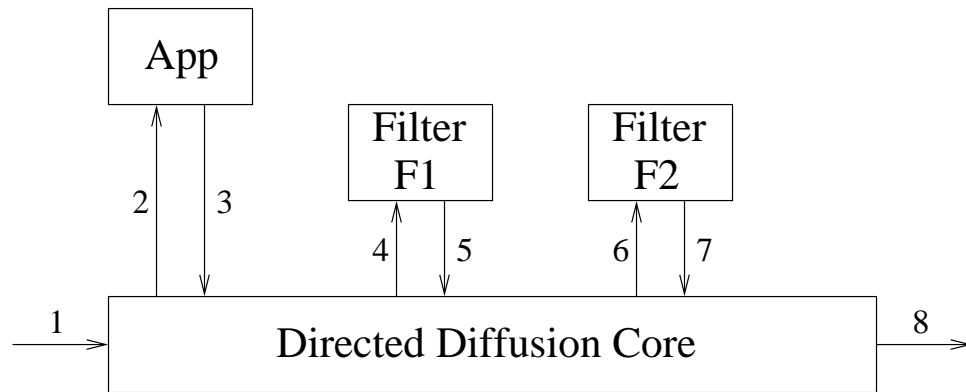


Figure 20.1: Message flow in directed diffusion

**apps** contains sample source and sink applications like gear, ping and rmst.

**lib** has DiffusionRouting class definitions and definitions of diffusion application class. In addition there are sub-dirs called main and nr. main houses misc diffusion utility code. nr includes attribute definition and the class NR which is an abstract factory for the API (so that either ISI or MIT implementations may derive from it).

**ns** contains ns wrappers for diffusion code. These wrapper classes allow the core diffusion code and the diffusion API to be seamlessly incorporated into the NS class hierarchy. The DiffRoutingAgent is a wrapper for the Core Diffusion code, and DiffAppAgent is a wrapper for the DiffusionRouting (API) code.

**filter\_core** has the core diffusion agent.

**filters** holds the different filters supported by diffusion implementation including two-phase-pull, one-phase-pull, gear, rmst, log, tag and srcrt (as of 10/03/03).

The above Figure 20.1 is from SCADDS' network routing API document available from their homepage (URL given earlier). The document describes attribute factories, matching rules for attributes, how applications interface with the core diffusion layer, and filter/timer APIs. All messages coming from/going out in the network is received at/sent out from the core diffusion layer. Messages can also come to core-diffusion layer from local applications and/or filters that might be connected to the node. The applications use the publish/subscribe/send interface to send interest and data messages to the network.

The core diffusion agent and diffusion application agent are attached to two well-known ports defined in `~ns/tcl/lib/ns-default.tcl`. Diffusion applications attached to the node call the underlying diffusion application agent for publishing/subscribing/sending data.

## 20.3 Some mac issues for diffusion in ns

In the shim layer that sits between diffusion and ns, (see `diffusion3/ns` dir for code implementing this layer) all diffusion packets are encapsulated within ns packets and are marked to be broadcasted. In previous versions all diffusion packets were marked to be broadcast in ns. This is now changed. Now all diffusion pkts in ns uses the diffusion next\_hop info thus allowing both broadcast and unicast.

So previously this only-brdcast feature supported for diffusion packets resulted in some problems at the mac layer. The mac-802.11 does not try to re-transmit a broadcast packet in case there is a collision and the packet is dropped. Coupled to this was the fact that mac-802.11 didn't do random selection of slots in the contention window before it transmitted a packet (a broadcast data or rts for unicast pkts). As a result there were a high number of collisions at the mac layer and a lot of packets were lost. This was fixed by adding random selection of slots before mac tx'ed a broadcast packet (or a rts packet).

However if we have a large and dense topology, there is a chance that two or more nodes may select the same slot in the mac contention window (the contention window size varies from 31 to 1023 for DSSS PHY MIB specifications). Thus now we need to add some extra jitter at the higher application layer. Diffusion has a provision to do this by compiling ns with the macro `USE_BROADCAST_MAC`. What this does is it in addition to delaying certain messages (to avoid collisions), when run with a `BROADCAST` MAC layer, diffusion will use a different set of values for delays and jitters. These different delay/jitter values are defined under `diffusion3/lib/main/config.hh`. Since this might increase the latency you might want to fine-tune the delay values by hand.

## 20.4 APIs for using filters in diffusion

As described earlier (see figure 20.1), filters can be attached to a diffusion node for various reasons. There can be basic diffusion filters providing two-phase-pull (`GradientFilter`) and one-phase-pull (`OnePhasePullFilter`) diffusion routing algorithms. There is the `GeoRoutingFilter` or gear that provides a certain location (co-ordinate) based routing algorithm. There is also other filters for RMST routing algorithm (`RmstFilter`), logging (`LogFilter`), source routing (`SourceRouteFilter`) and tagging (`TagFilter`). See Commands at a glance section for details on APIs for adding filters to a diffusion node.

## 20.5 Ping: an example diffusion application implementation

There is a ping application implemented under `diffusion3/apps/ping` subdir. The application consists of a ping sender and receiver. The receiver requests for data by sending out “interest”s in the network. The interests get diffused through the network. The ping-sender on receiving matching interests, sends out data.

### 20.5.1 Ping Application as implemented in C++

The ping-sender and -receiver classes, namely `PingSenderApp` and `PingReceiverApp` both derive from `DiffApp`, the parent class for all diffusion based applications. See `diffusion3/lib/diffapp{.cc,.hh}` for detailed implementation of the `DiffApp` class.

The ping-sender uses `MySenderReceive` object that handles all callbacks for it. Also the ping-sender defines two functions `setupSubscription()` and `setupPublication()`. The first function creates interest attributes that matches with data attributes it (the sender) has to offer. Next it calls the dr-library function `subscribe()`. The subscription is used by the ping-sender to create an internal state against which attributes for interests received from the network are matched against. In case of a match, the matching data is sent out into the network. Function `setupPublication()` create attributes for the data it has to offer and calls the library function `publish()` which in turn returns a publish handle. The ping-sender uses this handle to periodically send out data which is forwarded by the gradient to core-diffusion to be sent out into the network only if it finds a matching interest.

The ping-receiver object uses a similar callback object called `MyReceiverReceive`. And it defines a function `setupSubscription()` that creates attributes for the interest the receiver will be sending. Next it calls the dr library supported `subscribe()` which sends the interest out into the network. The `recv()` function is used to recv matching data and the receiver then calculates the latency for each data packet received from the ping-sender. The ping sender can be found under `ping_sender.cc,.h`. And the `ping_receiver` is implemented under `ping_receiver.cc,.h`. Some common defines and attribute factories for data/interest attributes are defined in `ping.hh` and `ping_common.cc`.

## 20.5.2 Tcl APIs for the ping application

An example script for the ping application is under `~ns/tcl/ex/diffusion3/simple-diffusion.tcl`. The example scenario consists of 3 nodes of which one is a ping-sender and another is a ping-receiver. The source and sink nodes are far away from one another and communicate only through a third node. The option `adhocRouting` is defined as `Directed_Diffusion`. This enables a core-diffusion agent to be created during the time of node creation. Also it creates a `diffusionApplication` agent if one is not present already. The option `diffusionFilter` needs to be provided at the time of node configuration that defines the one or more filters that would be added to the node. There is also an option for specifying `stopTime` which is the time the simulation ends. At this time there is a callback to a function that prints out all statistical data into `/tmp/diffusion-*.out`.

Node configuration is done as follows:

```
$ns_ node-config -adhocRouting $opt(adhocRouting) -llType $opt(ll)
                  -diffusionFilter $opt(filters)      -stopTime $opt(ptime)
```

The ping sender application is created in the following way:

```
set src_(0) [new Application/DiffApp/PingSender]
$ns_ attach-diffapp $node_(0) $src_(0)
$ns_ at 0.123 "$src_(0) publish"
```

The first line creates a ping-sender object. Simulator class method `attach-diffapp` basically attaches the application to the underlying `diffusionApplication` agent for that given node. The command `publish` essentially “starts” the sender application.

Similarly the ping sink is created as follows:

```
#Diffusion sink application
set snk_(0) [new Application/DiffApp/PingReceiver]
$ns_ attach-diffapp $node_(2) $snk_(0)
$ns_ at 1.456 "$snk_(0) subscribe"
```

The command `subscribe` starts the ping-receiver application.

Thus in order to create your own application, you need to :

1. define attribute-factories and attributes for application interest/data.
2. create the application class (using dr-library APIs)
3. add tcl commands to start the application

See `~ns/tcl/lib/ns-lib.tcl`, `ns-diffusion.tcl` for implementations of OTcl hooks for directed diffusion. Also see chapter on Mobility in this manual for details on mobility model and wireless simulations in *ns*.

## 20.6 Changes required to add yr diffusion application to ns

Let's say you have an application (it might even be a certain filter, which also is by class hierarchy, a diffusion application and it would derive from class `DiffApp`) that runs on the test-bed version. Now you want to run diffusion on *ns* and so want to use

yr application in the ns context. The few lines describe the changes/additions you need to make for yr diffusion application to work in ns environment.

We will consider onePhasePullFilter object (under diffusion3/filters/misc/log.\*) as an example. As a first step you need to create a split object out of the application class object, presumably defined as a pure c++ object. A split object is one that is created by the user in the interpreter (in OTcl space) and which is also has a shadow object in the compiled hierarchy (in c++ space). In ns, a split object is derived from class TclClass as shown below:

```
#ifdef NS_DIFFUSION
static class LogFilterClass : public TclClass
public:
    LogFilterClass() : TclClass("Application/DiffApp/LogFilter")
    TclObject * create(int argc, const char*const* argv)
        return(new LogFilter());

    class_log_filter;
#endif //DIFFUSION
```

Note that all filter objects specifically have a handle to the DiffAppAgent (the diffusion routing object) passed in the constructor call. Filter objects get created from function create-diffusionApp-agent diffFilters defined in ns-diffusion.tcl. Users need not specifically call the OTcl function create-diffusionApp-agent as it is called during node creation based on the node-configuration parameters. See how filters are defined in node-config under commands at a glance section. However application objects which are not filter objects (like ping\_sender, push\_receiver etc) are created by users directly from user scripts. And in that case the handle to DiffAppAgent is passed using \$ns attach-diffapp \$node \$app where the application \$app is attached to the node object \$node.

So for the reasons explained above the constructors are different in non NS\_DIFFUSION context as shown below.

```
#ifdef NS_DIFFUSION
LogFilter::LogFilter()
#else
LogFilter::LogFilter(int argc, char **argv)
#endif // NS_DIFFUSION

    // Create Diffusion Routing class
#ifdef NS_DIFFUSION
    parseCommandLine(argc, argv);
    dr_ = NR::createNR(diffusion_port_);
#endif // !NS_DIFFUSION

    filter_callback_ = new LogFilterReceive(this);

#ifdef NS_DIFFUSION
    // Set up the filter
    filter_handle_ = setupFilter();
    ...
    ...
#endif // !NS_DIFFUSION
```

Next you need to add the c++ function `command(..)` that allows execution of tcl commands through the compiled shadow object. For example the otcl command `start` is used to start a filter application as follows `$app start`. While commands `publish` and `subscribe` are used to start sender and receiver applications respectively. The command function is added, again with the `NS_DIFFUSION` scope using `ifdef` statements, as follows:

```
#ifdef NS_DIFFUSION
int LogFilter::command(int argc, const char*const* argv)
{
    if (argc == 2)
        if (strcmp(argv[1], "start") == 0)
            run();
        return TCL_OK;

    return DiffApp::command(argc, argv);
}

#endif // NS_DIFFUSION
```

Note how the parent class `command` function is invoked in case the command string is not found. Look into `lib/diffapp.*` to see all otcl commands supported for the `DiffApp` class.

Once these changes made to your c++ code, you would also need to write a tcl script (see the section on test-suite for example tcl scripts) that uses your diffusion application using the right tcl APIs.

## 20.7 Test-suites for diffusion

we start with a simple testcase of 3 nodes with 1 ping source and 1 ping sender. There are other tests for 2 phase-pull(2pp), 1 phase-pull(1pp), push and gear (with 2pp and push) scenarios. In future we plan to extend the test-suite for testing different components/functionalities of directed diffusion. All diffusion3 related test cases can be found under `~ns/tcl/test/test-suite-diffusion3.tcl`.

## 20.8 Commands at a glance

Following is a list of commands used for diffusion related simulation in ns.

```
$ns_ node-config -adhocRouting $opt(adhocRouting)
               -llType $opt(ll)
               ...
               -diffusionFilter $opt(filters)
               -stopTime $(pre-stop)
               ...
```

where,

value of `opt(adhocRouting)` is set to `Directed_Diffusion`

This command is used to enable directed diffusion in wireless nodes.

value of `opt(filters)` can be a list of filters that is required to be attached to diffusion nodes. This command allows adding filter objects to diffusion-enabled nodes.

value of `opt(pre-stop)` is usually the time simulation stops when all statistical data is dumped. This command allows dumping of statistical data into an output file after running a diffusion simulation.

```
set src [new Application/DiffApp/PingSender]
```

This command is used to create ping-sender application.

```
set snk [new Application/DiffApp/PingReceiver]
```

This command is used to create ping-receiver application.

```
set src [new Application/DiffApp/PushSender]
```

This command is used to create push-sender application.

```
set snk [new Application/DiffApp/PushReceiver]
```

This command is used to create push-receiver application.

```
set src [new Application/DiffApp/GearSenderApp]
```

This command is used to create gear-sender application.

```
set snk [new Application/DiffApp/GearReceiverApp]
```

This command is used to create gear-receiver application.

```
$gearApp push-pull-options <push/pull> <point/region> <co-ordinatesX1> <X2> <Y1> <Y2>
```

This command defines the type of routing algorithm gear is using. In case the second option is defined as region, all four co-ordinates should be defined. While if point is chosen, only X1 and Y1 maybe defined.

```
$ns_ attach-diffapp $node_ $src_
```

where the diffusion application `$src_` gets attached to the given `$node_`.

```
$src_(0) publish
```

Command to start a ping source (sender).

```
$snk_(0) subscribe
```

Command to start a ping sink (receiver).



## Chapter 21

# XCP: eXplicit Congestion control Protocol

XCP is a feedback-based congestion control system that uses direct, explicit, router feedback to avoid congestion in the network. It is designed for both scalability and generality. It was developed by Dina Katabi, starting from a suggestion by Mark Handley (refer to [31] for detailed descriptions). The *ns* code for XCP which was originally developed by Dina Katabi was modified, extended and integrated into ns-2.28 at USC/ISI. It still continues to evolve as of today. If you are interested in looking at Dina's original source code please go to her website at <http://www.ana.lcs.mit.edu/dina/XCP/>

### 21.1 What is XCP?

XCP is a congestion control protocol that can be applied to any transport protocol. It performs especially well in very high delay-bandwidth product networks. Typically in large bandwidth-delay product networks, efficiency of TCP goes down in the event of multiple of packet losses and becomes unstable irrespective of queueing schemes used. However in similar environments, XCP, using a control theory based feedback model, achieves much more efficiency, stability and fairness by sending feedback from the network to the sender for setting the data flow into the network.

XCP's scalability results from the fact that it requires no per-flow state in the router to calculate the feedback. Most router-assisted congestion control systems maintain per-flow information used to allocate the resources. XCP keeps very little information in the router, and this information is chosen to minimize both the amount of router state and the per-packet operations needed to update that state.

For generality, XCP divides the resource allocation function between two controllers: a congestion controller that ensures that flows use all available capacity, and a fairness controller that ensures that flows are allocated the capacity fairly. Most congestion control systems fail to make this division, much less to implement as two conceptually distinct systems. This division allows a clear exposition and implementation of two basic resource allocation functions in XCP. XCP sources send additional information about their current round-trip times and router-assigned throughput in each packet. XCP routers insert feedback into the packets that is interpreted by the sources.

Although XCP may be implemented for any transport protocol, however as an initial test it has been implemented in TCP. The next section gives details of the way XCP is implemented in *ns*.

## 21.2 Implementation of XCP in NS

In *ns*, the XCP implementation can be found under *~ns/xcp* directory. The protocol needs to be deployed in the TCP end points (source and receiver) as well within the intermediate nodes which is mostly the router and may sometimes be a link-layer switch as well. The end-point part of XCP code may be found under *xcp-end-sys.cc,h* and the router portion of the code may be found under *xcp.cc,h* and *xcpq.cc,h*.

### 21.2.1 Endpoints in XCP

The end points consist of TCP source and sink agents using XCP as their congestion control mechanism. The intermediate node or router writes feedback in each packet header as the *delta\_throughput* value, about the data capacity that may be incremented if feedback is positive and should be decreased if negative. When this packet reaches the receiver this *delta\_throughput* value is returned to the sender in a *reverse\_feedback* field of a congestion header in the returning packet, which is an ACK packet in case of TCP.

The sender upon receiving this *reverse\_feedback* value adjusts its sending rate by increasing or decreasing its congestion window size as the case maybe.

The packet header that is used by XCP is implemented as a structure called *hdr\_xcp* in *ns* and is defined as follows:

```
double x_; // idealized inter-packet time
double rtt_;
enum {
    XCP_DISABLED = 0,
    XCP_ENABLED,
    XCP_ACK,
} xcp_enabled_; // to indicate that the flow is XCP enabled
int xcpId_; // Sender's ID (debugging only)
double cwnd_; // The current window (debugging only)
double reverse_feedback_;

// --- Initialized by source and Updated by Router
double delta_throughput_;
unsigned int controlling_hop_;// router ID (debugging only)
```

The xcp receiver is responsible for copying the *delta\_throughput* value into the *reverse\_feedback* field of the ack packets. In some cases where delayed ack is used, the receiver calculates the sum of the *delta\_throughput* values in arriving packets for copying into the *reverse\_feedback* field of the outgoing ack packet.

The *controlling\_hop\_* field that carries the address of the router who has last updated the feedback is used for debugging purposes only.

In case of a packet loss in the network, TCP's Van Jacobson congestion control should most likely override XCP. However in *ns* this happens a little differently. With receiving of duplicate acks indicating packet loss, the *cwnd* gets halved and fast retransmit and fast recovery algorithms come into play. However xcp routers continue to send feedback to the source based on which the source tries to open its *cwnd*. So it seems to be a mish-mash of VJCC and XCP algorithms. However for most cases this issue doesnot arise as XCP router very rarely experiences a pkt drop as the queue is continuously regulated and drained by XCP. Understanding the correct behaviour of XCP in face of pkt loss is an area of current study and shall be implemented in the future.

### 21.2.2 XCP Router

The XCP router consists of a wrapper class that holds virtual queues for XCP, TCP and OTHER traffic flows. OTHER flow maybe anything other than XCP and TCP. In the current implementation, the XCP queue is a drop-tail queue while those for TCP and OTHER use RED.

These underlying queues are bundled in a wrapper class XCPWrapQ that provides necessary interface to the XCP router. The XCP/TCP/OTHER queues are serviced in a Weighted Round-Robin manner. Each queue has a weight that determines the percentage of the service it receives. The current queue weights of 0.5 each for the XCP and TCP allows equal service between the two. The third queue reserved for OTHER flows has not been used as yet and has a weight of 0.0.

OTCL Class Queue/XCP has a flag named `tcp_xcp_on_` which is set to a default value of 0. This should be set to 1 for those simulations that use both XCP and TCP flows. This flag is used to split the link capacity of the router between the XCP and TCP queues in simulations that use both flow types. This is supposed to be a temporary fix and should go away once the dynamic queue weights come into effect. A caveat for the `tcp_xcp` flag is that it is set as an OTcl class variable and not per instance variable. This might cause some problems in topologies that uses mix of intermittent xcp and tcp flows for which some router would require to support both TCP and XCP and some wouldn't.

Every packet received by the wrapper queue class is first marked or assigned a code point depending on the type of the packet. Packets, for the current TCP implementation, are marked for XCP, TCP/TCP-ACK and OTHER packets. This code point is used to enqueue packets in the right queue. The wrapper class is implemented in `xcp.cc,h`.

### 21.2.3 XCP queue

The XCP queue is responsible for sending back feedback in every packet header which is used by the sender to control rate of sending data packets into the network. XCP uses two control algorithms, the congestion controller and the fairness controller that are executed once every estimation control interval,  $T_e$ .

In *ns* the `estimation_timer` is used to maintain this interval which is based on the average rtt values of the (xcp) flows seen through the router. However there may be better ways of defining this interval. The outstanding queue in the router is measured at a separate interval  $T_q$ , for which a `queue_timer` is used. Finally an `rtt_timer` is used to measure certain parameters in the router like packet drops, queue size, utilization etc for a given interval  $T_r$ , that may either be set by user from tcl scripts or it may use the highest rtt value seen for all flows in the router.

The `rtt_timer` interval value,  $T_r$  maybe set from tcl using the following API:

```
$queue queue-sample-everyrtt $rtt_value
```

where `$queue` is a handle to the xcp router and `$rtt_value` is the interval for which xcp queue parameters like packet drop, queue size etc shall be measured. See example script under `~ns/tcl/ex/xcp/parking_lot_topo/parking_lot_topo.tcl` for use of this API.

On packet arrival the total input traffic seen at the XCP queue is incremented by the size of the packet received. The sum of inversed throughput and sum of rtt by throughput is incremented as well. Values for throughput and rtt are read from the xcp header as set by the TCP source. Each value is normalised by the packet size.

On the event of the estimation timer going off, average rtt of all flows is estimated using the above two sums as follows

```
avg_rtt = sum_rtt_by_throughput / sum_inv_throughput
```

The aggregate feedback is calculated based on the available bandwidth capacity of the router, arriving traffic bandwidth and

the persistent queue length in the router. Further detailed explanation of calculations used by the XCP router algorithm can be found in XCP specification available from XCP's website at <http://www.isi.edu/isi-xcp/docs/draft-falk-xcp-spec-00.txt>

Each packet carries the current throughput value of the flow and a throughput adjustment or the `delta_throughput` in its header. The XCP router based on the feedback values it calculates in the estimation control timeout, calculates a per-packet throughput adjustment feedback for every packet. Positive feedback is applied equally per-flow while negative feedback is made proportional to each flow's capacity. Also a downstream router can change the `delta_throughput` value in a packet's header only if the feedback value calculated is less than that in the header (written by an less congested upstream router). The implementation of XCP queue in *ns* may be found in `xcpq.{cc,h}`.

## 21.3 XCP example script

Let's walk through a simple `xcp` script that is similar to `~ns/tcl/ex/xcp/xcp_test.tcl` The example uses a small dumbbell topology having 3 `xcp` sources running over a bottleneck link.

The topology is setup using the node and link creation APIs. The bottleneck is a duplex link that has a `xcp` router in both directions. For details on creating nodes, links etc in *ns* see Marc Greis' NS tutorial at <http://www.isi.edu/nsnam/ns/tutorial>.

The bottleneck link having a XCP queue is created as follows:

```
set R0 [$ns node]           ;# create Bottleneck between nodes R0 and R1
set R1 [$ns node]
$ns duplex-link $R0 $R1 <BW>Mb <delay>ms XCP
```

The side links connecting source nodes to the bottleneck link have XCP queues as well. The API `queue-limit` allows users to set the buffer size in the queue.

The `xcp` source and sink is created as follows (very similar to `tcp`):

```
set xcp [new Agent/TCP/Reno/XCP]
$ns attach-agent $src_node $xcp
set xcp_sink [new Agent/TCPSink/XCPSink]
$ns attach-agent $rcvr_node $xcp_sink
$ns connect $xcp $xcp_sink
...
...
```

There is a tcl class `GeneralSender` used in the example script that sets up `xcp` agents in the source nodes and then connects them to the `xcp` receiver in the destination node. An FTP source is used in all the 3 sources.

Note that once the topology is set up the link bandwidth information needs to be propagated to the `xcp` queue as this is used by the `xcp` router for feedback calculation. So for every `xcp` queue use the following tcl command:

`$xcp_queue set-link-capacity <bandwidth_in_bits_per_sec>` Next we need to trace variables in the `xcp` router and `xcp` sources. The `GeneralSender` class procedure `trace-xcp` sets up tracing for `xcp` sources using variable-tracing in *ns*.

```

GeneralSender instproc trace-xcp parameters {
    $self instvar tcp_id_ tcpTrace_
    global ftracetcpsid_
    set ftracetcpsid_ [open xcp$id_.tr w]
    set tcpTrace_ [set ftracetcpsid_]
    $tcp_ attach-trace [set ftracetcpsid_]
    if { -1 < [lsearch $parameters cwnd] } { $tcp_ tracevar cwnd_ }
    if { -1 < [lsearch $parameters seqno] } { $tcp_ tracevar t_seqno_ }
}

```

For tracing xcp queue it is required to attach a file descriptor to the xcp queue.

```
$xcpq attach <file-descriptor>
```

This is an example of how the trace at an xcp source looks like:

```

0.00000 2 0 1 0 cwnd_ 1.000
0.00000 2 0 1 0 t_seqno_ 0
0.079 x x x x throughput 0.1
0.07900 2 0 1 0 t_seqno_ 1
0.119064 x x x x reverse_feedback_ 0
0.119064 x x x x controlling_hop_ 0
0.119064 x x x x newcwnd 1
0.11906 2 0 1 0 cwnd_ 2.000
0.119064 x x x x throughput 50000
0.11906 2 0 1 0 t_seqno_ 2
0.119064 x x x x throughput 50000
0.11906 2 0 1 0 t_seqno_ 3

```

The first field gives the timestamp; the next 4 fields give the source id (node/port) and destination id (node/port) for the xcp flow. The next field gives the name of the variable being traced followed by the value of the variable. Note that variables like `cwnd_`, `t_seqno_` are using variable tracing which is a function supported by the OTcl lib. While variables like `throughput`, `reverse_feedback` use the XCPAgent class function `trace_var` defined in `xcp-end-sys.cc`. For more on variable tracing in *ns* please read section 3.4.3 in the *ns* manual at <http://www.isi.edu/nsnam/ns/doc/index.html>

And example of trace output at a xcp bottleneck router looks like below:

```

Tq_ 0.0472859 0.025
queue_bytes_ 0.0472859 0
routerId_ 0.0472859 0
pos_fbk 0.053544 0
neg_fbk 0.053544 0
delta_throughput 0.053544 0
Thruput2 0.053544 60000
pos_fbk 0.054024 0
neg_fbk 0.054024 0
delta_throughput 0.054024 0

```

```

Thruput2 0.054024 60000
residue_pos_fbk_not_allocated 0.0638023 0
residue_neg_fbk_not_allocated 0.0638023 0
input_traffic_bytes_ 0.0638023 2480
avg_rtt_ 0.0638023 0.04

```

Here the first field describes the name of the variable, the second field gives the timestamp and the third field gives the value of the variable. The XCPQueue class function `trace_var()` is used to trace variables in the xcp queue.

Additionally packet traces may be created in *ns* using the following tcl APIs:

```

set f_all [open out.tr w]
$ns trace-all $f_all

```

First open a file and then attach the file descriptor to the *ns* trace object such that a trace of each packet as it travels through the network is logged and dumped into the output file.

An example of such a file would look like:

```

+ 0.003 4 0 xcp 40 ----- 2 4.0 1.2 0 0
- 0.003 4 0 xcp 40 ----- 2 4.0 1.2 0 0
r 0.013016 4 0 xcp 40 ----- 2 4.0 1.2 0 0
+ 0.013016 0 1 xcp 40 ----- 2 4.0 1.2 0 0
- 0.013016 0 1 xcp 40 ----- 2 4.0 1.2 0 0
r 0.023032 0 1 xcp 40 ----- 2 4.0 1.2 0 0
+ 0.023032 1 0 ack 40 ----- 2 1.2 4.0 0 1
- 0.023032 1 0 ack 40 ----- 2 1.2 4.0 0 1
r 0.033048 1 0 ack 40 ----- 2 1.2 4.0 0 1
+ 0.033048 0 4 ack 40 ----- 2 1.2 4.0 0 1
- 0.033048 0 4 ack 40 ----- 2 1.2 4.0 0 1
r 0.043064 0 4 ack 40 ----- 2 1.2 4.0 0 1
+ 0.043064 4 0 xcp 1200 ----- 2 4.0 1.2 1 2
- 0.043064 4 0 xcp 1200 ----- 2 4.0 1.2 1 2
+ 0.043064 4 0 xcp 1200 ----- 2 4.0 1.2 2 3
- 0.043544 4 0 xcp 1200 ----- 2 4.0 1.2 2 3

```

Lets try to read the first line:

```

+ 0.003 4 0 xcp 40 ----- 2 4.0 1.2 0 0

```

+ means a packet is enqueued in the queue (in node 4) as it hopped between node 4 to node 0. You'll find traces showing packets enqueued (+) and then dequeued (-) at the queue, after which it is transmitted over the link to be received by the next node. packet type is xcp and it is of size 40 bytes. The xcp flow has an id of 2 and the packet header has a source node/port id of 4.0 and dest node/port id of 1.2 and the unique packet id is 0.

## 21.4 Test-suites for XCP

The xcp test-suite uses 3 tests. The first one is similar to the one we discussed in the earlier section. It consists of a dumb-bell topology where 3 xcp flows share a bottleneck link. The second test has a similar topology having 3 xcp and 1 tcp flow sharing the same bottleneck. And finally the last test is built on Dina Katabi's parking-lot experiment referred in her SIGCOMM'02 paper. It is a downsized version of Dina's example. The test uses a 9-hop link string topology. It has 10 long XCP flows that flow through the entire length of the chain topology. Then it has 10 XCP flows that run through each hop, starting at (n-1)th hop and ending at nth hop and so on creating the intermittent flows. And finally there are long XCP flows in the reverse direction, starting from the last (10th) node and ending in the first (1st) node. There is a bottleneck at the middle of the chain topology. Thus the third test employs a large and complex topology and shows the utilization, queue length and packet drop values at every link.

## 21.5 Commands at a glance

Following is a list of commands used for xcp related simulation in ns.

```
set xcp_src [new Agent/TCP/Reno/XCP]
```

This command creates an XCP source agent.

```
set xcp_dst [new Agent/TCPSink/XCPSink]
```

This command creates an XCP sink.

```
$ns duplex-link $R0 $R1 <BW>Mb <delay>ms XCP
```

This code creates a duplex link with specified bandwidth and link delay using an XCP router between node R0 and R1.

```
$xcp_queue set-link-capacity <bandwidth_in_bits_per_sec>
```

This command propagates the link bandwidth information to the xcp queue which uses it for the router feedback calculation.

```
set tfile [open tfile w]
```

```
$xcp_queue attach $tfile
```

This Tcl command allows a file to be attached for tracing xcp queue parameters.

```
$xcp_src attach-trace <file-descriptor> $xcp_src tracevar <var-to-traced>
```

This command allows the xcp sources to trace variables.

```
$queue queue-sample-everyrtt $rtt_value
```

This command allows the user to set rtt interval values at which variables like packet\_drops and queue lengths are measured at the xcp queue.

```
Queue/XCP set tcp_xcp_on_ 1
```

This flag lets tcp and xcp flows to use the same xcp router. This flag is a temporary fix and should go away when dynamic queue weights come into effect.

## Chapter 22

# DelayBox: Per-Flow Delay and Loss

DelayBox is an ns node that should be placed in between the source and destination nodes. With Delaybox, packets from a TCP flow can be delayed, dropped, and/or forced through a bottleneck link before being passed on to the next node. A distribution can be used to specify delay, loss, and/or bottleneck link speed for a source - destination pair. Each flow between that source - destination pair draws from the distribution to determine its characteristics. Delays in DelayBox are per-flow, rather than per-packet. Since DelayBox distinguishes between flows, the `fid_` variable (flow identifier) should be set for each flow in the simulation. DelayBox can be used with both Tcp and FullTcp agents.

### 22.1 Implementation Details

DelayBox maintains two tables: a rule table and a flow table. Entries in the rule table are added by the user in the OTcl simulation script and give an outline of how flows from a source to a destination should be treated. The fields are source, destination, delay Random Variable (in ms), loss rate Random Variable (in fraction of packets dropped), and bottleneck link speed Random Variable (in Mbps). The bottleneck link speed field is optional. Entries in the flow table are created internally and specify exactly how each flow should be handled. Its values are obtained by sampling from the distributions given in the rule table. The fields are source, destination, flow ID, delay, loss, and bottleneck link speed (if applicable). Full-TCP flows are defined as beginning at the receipt of the first SYN of a new flow ID and ending after the sending of the first FIN. Packets after the first FIN are forwarded immediately (i.e., they are neither delayed nor dropped by DelayBox). For TcpAgent, flows are defined as beginning at the receipt of the first 40 byte packet of a new flow ID. Since there are no FIN packets in TcpAgent, TcpAgent flows are never considered finished nor are they removed from the flow table.

DelayBox also maintains a set of queues to handle delaying packets. There is one queue per entry in the flow table. These queues are implemented as delta queues, in that the time to transfer the packet is kept only for the head packet. All other packets are stored with the difference between the time they should be transferred and the time the previous packet should be transferred. The actual time the previous packet should be transferred is stored in the variable `deltasum_`, named so because it is the sum of all delta values in the queue (including the head packet's transfer time). If the bottleneck link speed has been specified for the flow, a processing delay is computed for each packet by dividing the size of the packet by the flow's bottleneck link speed.

When a packet is received, its transfer time (current time + delay) is calculated. (This transfer time is the time that the first bit of the packet will begin transfer. Packets that wait in the queue behind this packet must be delayed by the amount of time to transfer all bits of the packet over the bottleneck link.) There are two scenarios to consider in deciding how to set the packet's delta:



1. If the packet is due to be transferred before the last bit of the last packet in the queue, its delta (the time between transferring the previous packet and transferring this packet) is set to the previous packet's processing delay. This packet has to queue behind the previous packet, but will be ready to be transmitted as soon as the previous packet has completed its transfer.
2. If the packet is due to be transferred after the last bit of the last packet in the queue, its delta is difference between this packet's transfer time and the previous packet's transfer time.

If the current packet is the only packet in the queue, DelayBox schedules a timer for the receipt of the packet. When this timer expires, DelayBox will pass the packet on to the standard packet forwarder for processing. Once a packet has been passed up, DelayBox will look for the next packet in the queue to be processed and schedule a timer for its transfer. All packets, both data and ACKs, are delayed in this manner.

Packets that should be dropped are neither queued nor passed on. All packets in a queue are from the same connection and are delayed the same amount (except for delays due to packet size) and are dropped with the same probability. **Note:** Drops at DelayBox are not recorded in the trace-queue file.

## 22.2 Example

More examples are available in the `tcl/ex/delaybox/` directory of the ns source code. The validation script `test-suite-delaybox` is in `tcl/test/` and can be run with the command `test-all-delaybox` from that directory.

```
# test-delaybox.tcl - NS file transfer with DelayBox

# setup ns
remove-all-packet-headers;           # removes all packet headers
add-packet-header IP TCP;             # adds TCP/IP headers
set ns [new Simulator];               # instantiate the simulator

global defaultRNG
$defaultRNG seed 999

# create nodes
set n_src [$ns node]
set db(0) [$ns DelayBox]
set db(1) [$ns DelayBox]
set n_sink [$ns node]

# setup links
$ns duplex-link $db(0) $db(1) 100Mb 1ms DropTail
$ns duplex-link $n_src $db(0) 100Mb 1ms DropTail
$ns duplex-link $n_sink $db(1) 100Mb 1ms DropTail

set src [new Agent/TCP/FullTcp]
set sink [new Agent/TCP/FullTcp]
$src set fid_ 1
$sink set fid_ 1

# attach agents to nodes
$ns attach-agent $n_src $src
```

```

$ns attach-agent $n_sink $sink

# make the connection
$ns connect $src $sink
$sink listen

# create random variables
set recvr_delay [new RandomVariable/Uniform];      # delay 1-20 ms
$recvr_delay set min_ 1
$recvr_delay set max_ 20
set sender_delay [new RandomVariable/Uniform];     # delay 20-100 ms
$sender_delay set min_ 20
$sender_delay set max_ 100
set recvr_bw [new RandomVariable/Constant];        # bw 100 Mbps
$recvr_bw set val_ 100
set sender_bw [new RandomVariable/Uniform];        # bw 1-20 Mbps
$sender_bw set min_ 1
$sender_bw set max_ 20
set loss_rate [new RandomVariable/Uniform];        # loss 0-1% loss
$loss_rate set min_ 0
$loss_rate set max_ 0.01

# setup rules for DelayBoxes
$db(0) add-rule [$n_src id] [$n_sink id] $recvr_delay $loss_rate $recvr_bw
$db(1) add-rule [$n_src id] [$n_sink id] $sender_delay $loss_rate $sender_bw

# output delays to files
$db(0) set-delay-file "db0.out"
$db(1) set-delay-file "db1.out"

# schedule traffic
$ns at 0.5 "$src advance 10000"
$ns at 1000.0 "$db(0) close-delay-file; $db(1) close-delay-file; exit 0"

# start the simulation
$ns run

```

## 22.3 Commands at a Glance

The following commands on the DelayBox class can be accessed from OTcl:

```

[$ns DelayBox]
Creates a new DelayBox node.

```

```

$delaybox add-rule <srcNodeID> <dstNodeID> <delayRV> [<lossRV>] [<linkSpeedRV>]
Add a rule to the rule table, specifying delay, loss rate, and bottleneck link speed RandomVariables for packets flowing from srcNode to dstNode. Delay is required, but loss rate and link speed are optional.

```

```

$delaybox list-rules
List all rules in the rule table

```

`$delaybox list-flows`

List all flows in the flow table

`$delaybox set-asymmetric`

Specifies that the delay should be only on the data path rather than applied to both the data and ACK paths

`$delaybox set-delay-file <filename>`

Output delays for each flow to filename. Format: `srcNode dstNode fid delay`

`$delaybox close-delay-file`

Closes the file where delays are written

`$delaybox set-debug <int>`

Set the debugging level

- 1: Output when packets are dropped at DelayBox
- 2: Level 1 +  
Contents of the queue at each queue operation

## **Part III**

# **Support**

# Chapter 23

## Debugging ns

*ns* is a simulator engine built in C++ and has an OTcl (Object-oriented Tcl) interface that is used for configuration and commands. Thus in order to debug *ns* we will have to deal with debugging issues involving both OTcl and C++. This chapter gives debugging tips at Tcl and C++ level and shows how to move to-and-fro the Tcl and C++ boundaries. It also briefly covers memory debugging and memory conservation in *ns*.

### 23.1 Tcl-level Debugging

*Ns* supports Don Libs' Tcl debugger ( see its Postscript documentation at <http://expect.nist.gov/tcl-debug/tcl-debug.ps.Z> and its source code at <http://expect.nist.gov/tcl-debug/tcl-debug.tar.gz> ). Install the program or leave the source code in a directory parallel to *ns-2* and it will be built. Unlike *expect*, described in the *tcl-debug* documentation, we do not support the *-D* flag. To enter the debugger, add the lines "debug 1" to your script at the appropriate location. In order to build *ns* with the debugging flag turned on, configure *ns* with configuration option "*--enable-debug*" and incase the Tcl debugger has been installed in a directory not parallel to *ns-2*, provide the path with configuration option "*--with-tcldebug=<give/your/path/to/tcl-debug/library>*".

An useful debugging command is `$ns_ gen-map` which may be used to list all OTcl objects in a raw form. This is useful to correlate the position and function of an object given its name. The name of the object is the OTcl handle, usually of the form `_o###`. For TclObjects, this is also available in a C++ debugger, such as *gdb*, as `this->name_`.

### 23.2 C++-Level Debugging

Debugging at C++ level can be done using any standard debugger. The following macro for *gdb* makes it easier to see what happens in Tcl-based subroutines:

```
## for Tcl code
define pargvc
set $i=0
while $i < argc
    p argv[$i]
    set $i=$i+1
```

```

    end
end
document pargvc
Print out argc argv[i]'s common in Tcl code.
(presumes that argc and argv are defined)
end

```

## 23.3 Mixing Tcl and C debugging

It is a painful reality that when looking at the Tcl code and debugging Tcl level stuff, one wants to get at the C-level classes, and vice versa. This is a smallish hint on how one can make that task easier. If you are running ns through gdb, then the following incantation (shown below) gets you access to the Tcl debugger. Notes on how you can then use this debugger and what you can do with it are documented earlier in this chapter and at this URL (<http://expect.nist.gov/tcl-debug/tcl-debug.ps.Z>).

```

(gdb) run
Starting program: /nfs/prot/kannan/PhD/simulators/ns/ns-2/ns
...

Breakpoint 1, AddressClassifier::AddressClassifier (this=0x12fbd8)
    at classifier-addr.cc:47
(gdb) p this->name_
$1 = 0x2711e8 "_o73"
(gdb) call Tcl::instance().eval("debug 1")
15: lappend auto_path $dbg_library
dbg15.3> w
*0: application
    15: lappend auto_path /usr/local/lib/dbg
dbg15.4> Simulator info instances
_o1
dbg15.5> _o1 now
0
dbg15.6> # and other fun stuff
dbg15.7> _o73 info class
Classifier/Addr
dbg15.8> _o73 info vars
slots_ shift_ off_ip_ offset_ off_flags_ mask_ off_cmn_
dbg15.9> c
(gdb) w
Ambiguous command "w": while, whatis, where, watch.
(gdb) where
#0  AddressClassifier::AddressClassifier (this=0x12fbd8)
    at classifier-addr.cc:47
#1  0x5c68 in AddressClassifierClass::create (this=0x10d6c8, argc=4,
    argv=0xefffcdc0) at classifier-addr.cc:63
...
(gdb)

```

In a like manner, if you have started ns through gdb, then you can always get gdb's attention by sending an interrupt, usually a <Ctrl-C> on berkeloidrones. However, note that these do tamper with the stack frame, and on occasion, may (sometimes can (and rarely, does)) screw up the stack so that, you may not be in a position to resume execution. To its credit, gdb appears to be smart enough to warn you about such instances when you should tread softly, and carry a big stick.

## 23.4 Memory Debugging

The first thing to do if you run out of memory is to make sure you can use all the memory on your system. Some systems by default limit the memory available for individual programs to something less than all available memory. To relax this, use the `limit` or `ulimit` command. These are shell functions—see the manual page for your shell for details. `Limit` is for `csh`, `ulimit` is for `sh/bash`.

Simulations of large networks can consume a lot of memory. Ns-2.0b17 supports Gray Watson's `dmalloc` library (see its web documentation at <http://www.letters.com/dmalloc/> and get the source code from <ftp://ftp.letters.com/src/dmalloc/dmalloc.tar.gz>). To add it, install it on your system or leave its source in a directory parallel to `ns-2` and specify `--with-dmalloc` when configuring `ns`. Then build all components of `ns` for which you want memory information with debugging symbols (this should include at least `ns-2`, possibly `tlcl` and `otcl` and maybe also `tcl`).

### 23.4.1 Using dmalloc

In order to use `dmalloc` do the following:

- Define an alias

```
for csh: alias dmalloc 'eval `dmalloc -C \!*`',
for bash: function dmalloc { eval `command dmalloc -b $*` }%$
```

- Next turn debugging on by typing `dmalloc -l logfile low`
- Run your program (which was configured and built with `dmalloc` as described above).
- Interpret logfile by running `dmalloc_summarize ns <logfile`. (You need to download `dmalloc_summarize` separately.)

On some platforms you may need to link things statically to get `dmalloc` to work. On Solaris this is done by linking with these options: `"-Xlinker -B -Xlinker -static libraries -Xlinker -B -Xlinker -dynamic -ldl -lX11 -lXext"`. (You'll need to change `Makefile`. Thanks to Haobo Yu and Doug Smith for working this out.)

We can interpret a sample summary produced from this process on `ns-2/tcl/ex/newmcast/cmcast-100.tcl` with an exit statement after the 200'th `duplex-link-of-interfaces` statement:

Ns allocates 6MB of memory.

1MB is due to `TclObject::bind`

900KB is `StringCreate`, all in 32-byte chunks

700KB is `NewVar`, mostly in 24-byte chunks

`Dmalloc_summarize` must map function names to and from their addresses. It often can't resolve addresses for shared libraries, so if you see lots of memory allocated by things beginning with "ra=", that's what it is. The best way to avoid this problem is to build `ns` statically (if not all, then as much as possible).

`Dmalloc`'s memory allocation scheme is somewhat expensive, plus there's bookkeeping costs. Programs linked against `dmalloc` will consume more memory than against most standard mallocs.

`Dmalloc` can also diagnose other memory errors (duplicate frees, buffer overruns, etc.). See its documentation for details.

### 23.4.2 Memory Conservation Tips

Some tips to saving memory (some of these use examples from the `cmcast-100.tcl` script): If you have many links or nodes:

**Avoid `trace-all`** : `$ns trace-all $f` causes trace objects to be pushed on all links. If you only want to trace one link, there's no need for this overhead. Saving is about 14 KB/link.

**Use arrays for sequences of variables** : Each variable, say `n$i` in `set n$i [$ns node]`, has a certain overhead. If a sequence of nodes are created as an array, i.e. `n($i)`, then only one variable is created, consuming much less memory. Saving is about 40+ Byte/variable.

**Avoid unnecessary variables** : If an object will not be referred to later on, avoid naming the object. E.g. `set cmcast(1) [new CtrMcast $ns $n(1) $ctrmcastcomp [list 1 1]]` would be better if replaced by `new CtrMcast $ns $n(1) $ctrmcastcomp [list 1 1]`. Saving is about 80 Byte/variable.

**Run on top of FreeBSD** : `malloc()` overhead on FreeBSD is less than on some other systems. We will eventually port that allocator to other platforms.

**Dynamic binding** : Using `bind()` in C++ consumes memory for each object you create. This approach can be very expensive if you create many identical objects. Changing `bind()` to `delay_bind()` changes this memory requirement to per-class. See `ns/object.cc` for an example of how to do binding, either way.

**Disabling packet headers** : For packet-intensive simulations, disabling all packet headers that you will not use in your simulation may significantly reduce memory usage. See Section 12.1 for detail.

### 23.4.3 Some statistics collected by `dmalloc`

A memory consumption problem occurred in recent simulations (`cmcast-[150,200,250].tcl`), so we decided to take a closer look at scaling issue. See page <http://www-mash.cs.berkeley.edu/ns/ns-scaling.html> which demonstrates the efforts in finding the bottleneck.

The following table summarises the results of investigating the bottleneck:

<i>KBytes</i>	<i>cmcast-50.tcl(217 Links)</i>	<i>cmcast-100.tcl(950 Links)</i>
trace-all	8,084	28,541
turn off trace-all	5,095	15,465
use array	5,091	15,459
remove unnecessary variables	5,087	15,451
on SunOS	5,105	15,484

## 23.5 Memory Leaks

This section deals with memory leak problems in *ns*, both in Tcl as well as C/C++.



### 23.5.1 OTcl

OTcl, especially TclCL, provides a way to allocate new objects. However, it does not accordingly provide a garbage collection mechanism for these allocated objects. This can easily lead to unintentional memory leaks. Important: tools such as `dmalloc` and `purify` are unable to detect this kind of memory leaks. For example, consider this simple piece of OTcl script:

```
set ns [new Simulator]
for set i 0 $i < 500 incr i
    set a [new RandomVariable/Constant]
```

One would expect that the memory usage should stay the same after the first `RandomVariable` is allocated. However, because OTcl does not have garbage collection, when the second `RandomVariable` is allocated, the previous one is not freed and hence results in memory leak. Unfortunately, there is no easy fix for this, because garbage collection of allocated objects is essentially incompatible with the spirit of Tcl. The only way to fix this now is to always explicitly free every allocated OTcl object in your script, in the same way that you take care of malloc-ed object in C/C++.

### 23.5.2 C/C++

Another source of memory leak is in C/C++. This is much easier to track given tools that are specifically designed for this task, e.g., `dmalloc` and `purify`. `ns` has a special target `ns-pure` to build purified `ns` executable. First make sure that the macro `PURIFY` in the `ns Makefile` contains the right `-collector` for your linker (check `purify` man page if you don't know what this is). Then simply type `make ns-pure`. See earlier sections in this chapter on how to use `ns` with `libdmalloc`.

## Chapter 24

# Mathematical Support

The simulator includes a small collection of mathematical functions used to implement random variate generation and integration. This area of the simulator is currently undergoing some changes.

The procedures and functions described in this chapter can be found in `~ns/tools/rng.{cc, h}`, `~ns/tools/random.{cc, h}`, `~ns/tools/ranvar.{cc, h}`, `~ns/tools/pareto.{cc, h}`, `~ns/tools/expoo.{cc, h}`, `~ns/tools/integrator.{cc, h}`, and `~ns/tcl/lib/ns-random.tcl`.

### 24.1 Random Number Generation

The RNG class contains an implementation of the combined multiple recursive generator MRG32k3a proposed by L'Ecuyer [16]. The C++ code was adapted from [18]. This replaces the previous implementation of RNG, which used the minimal standard multiplicative linear congruential generator of Park and Miller [27]. The newer (MRG32k3a) RNG is used in ns versions 2.1b9 and later.

The MRG32k3a generator provides  $1.8 \times 10^{19}$  independent streams of random numbers, each of which consists of  $2.3 \times 10^{15}$  substreams. Each substream has a period (*i.e.*, the number of random numbers before overlap) of  $7.6 \times 10^{22}$ . The period of the entire generator is  $3.1 \times 10^{57}$ . Figure 24.1 provides a graphical idea of how the streams and substreams fit together.

A default RNG (`defaultRNG`), created at simulator initialization time, is provided. If multiple random variables are used in a simulation, each random variable should use a separate RNG object. When a new RNG object is created, it is automatically seeded to the beginning of the next independent stream of random numbers. Used in this manner, the implementation allows for a maximum of  $1.8 \times 10^{19}$  random variables.

Often, multiple independent replications of a simulation are needed (*i.e.*, to perform statistical analysis given multiple runs with fixed parameters). For each replication, a different substream should be used to ensure that the random number streams are independent. (This process is given as an OTcl example later.) This implementation allows for a maximum of  $2.3 \times 10^{15}$  independent replications. Each random variable in a single replication can produce up to  $7.6 \times 10^{22}$  random numbers before overlapping.

**Note:** Only the most common functions are described here. For more information, see [18] and the source code found in `tools/rng.h` and `tools/rng.cc`. For a comparison of this RNG to the more common LCG16807 RNG (and why LCG16807 is not a good RNG), see [17].

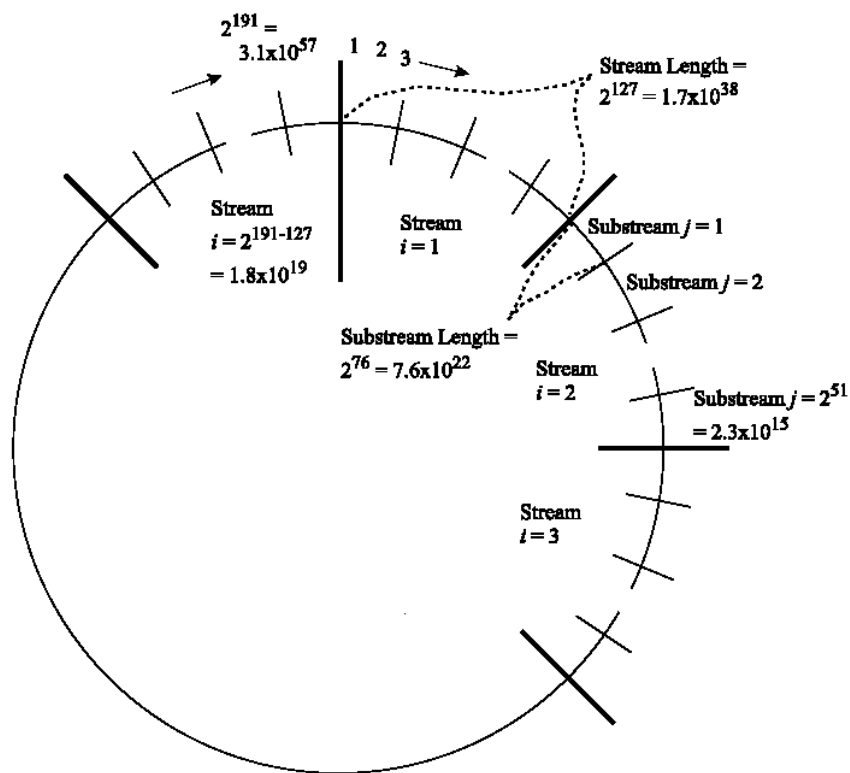


Figure 24.1: Overall arrangement of streams and substreams. [18]

### 24.1.1 Seeding The RNG

Due to the nature of the RNG and its implementation, it is not necessary to set a seed (the default is 12345). If you wish to change the seed, functions are available. You should only set the seed of the default RNG. Any other RNGs you create are automatically seeded such that they produce independent streams. The range of valid seeds is 1 to MAXINT.

To get non-deterministic behavior, set the seed of the default RNG to 0. This will set the seed based on the current time of day and a counter. **This method should not be used to set seeds for independent replications.** There is no guarantee that the streams produced by two random seeds will not overlap. The only way to guarantee that two streams do not overlap is to use the substream capability provided by the RNG implementation.

#### Example

```
# Usage: ns rng-test.tcl [replication number]

if {$argc > 1} {
    puts "Usage: ns rng-test.tcl \[replication number\]"
    exit
}
set run 1
if {$argc == 1} {
```

```

    set run [lindex $argv 0]
}
if {$run < 1} {
    set run 1
}

# seed the default RNG
global defaultRNG
$defaultRNG seed 9999

# create the RNGs and set them to the correct substream
set arrivalRNG [new RNG]
set sizeRNG [new RNG]
for {set j 1} {$j < $run} {incr j} {
    $arrivalRNG next-substream
    $sizeRNG next-substream
}

# arrival_ is a exponential random variable describing the time
# between consecutive packet arrivals
set arrival_ [new RandomVariable/Exponential]
$arrival_ set avg_ 5
$arrival_ use-rng $arrivalRNG

# size_ is a uniform random variable describing packet sizes
set size_ [new RandomVariable/Uniform]
$size_ set min_ 100
$size_ set max_ 5000
$size_ use-rng $sizeRNG

# print the first 5 arrival times and sizes
for {set j 0} {$j < 5} {incr j} {
    puts [format "%-8.3f  %-4d" [$arrival_ value] \
        [expr round([$size_ value])]]
}

```

## Output

```

% ns rng-test.tcl 1
6.358      4783
5.828      1732
1.469      2188
0.732      3076
4.002      626

% ns rng-test.tcl 5
0.691      1187
0.204      4924
8.849      857
2.111      4505
3.200      1143

```

## 24.1.2 OTcl Support

### Commands

The following commands on the RNG class can be accessed from OTcl and are found in `tools/rng.cc`:

`seed n` – seed the RNG to *n*, if *n* == 0, the seed is set according to the current time and a counter

`next-random` – return the next random number

`seed` – return the current value of the seed

`next-substream` – advance to the next substream

`reset-start-substream` – reset the stream to the beginning of the current substream

`normal avg std` – return a number sampled from a normal distribution with the given average and standard deviation

`lognormal avg std` – return a number sampled from a lognormal distribution with the given average and standard deviation

The following commands on the RNG class can be accessed from OTcl and are found in `tcl/lib/ns-random.tcl`:

`exponential mu` – return a number sampled from an exponential distribution with mean *mu*

`uniform min max` – return an integer sampled from a uniform distribution on [*min*, *max*]

`integer k` – return an integer sampled from a uniform distribution on [0, *k*-1]

### Example

```
# Usage: ns rng-test2.tcl [replication number]

if {$argc > 1} {
    puts "Usage: ns rng-test2.tcl \[replication number\]"
    exit
}
set run 1
if {$argc == 1} {
    set run [lindex $argv 0]
}
if {$run < 1} {
    set run 1
}

# the default RNG is seeded with 12345

# create the RNGs and set them to the correct substream
set arrivalRNG [new RNG]
set sizeRNG [new RNG]
for {set j 1} {$j < $run} {incr j} {
    $arrivalRNG next-substream
}
```

```

    $sizeRNG next-substream
}

# print the first 5 arrival times and sizes
for {set j 0} {$j < 5} {incr j} {
    puts [format "%-8.3f  %-4d" [$arrivalRNG lognormal 5 0.1] \
        [expr round([$sizeRNG normal 5000 100])]]
}

```

## Output

```

% ns rng-test2.tcl 1
142.776    5038
174.365    5024
147.160    4984
169.693    4981
187.972    4982

% ns rng-test2.tcl 5
160.993    4907
119.895    4956
149.468    5131
137.678    4985
158.936    4871

```

## 24.1.3 C++ Support

### Member Functions

The random number generator is implemented by the RNG class and is defined in `tools/rng.h`.

**Note:** The `Random` class in `tools/random.h` is an older interface to the standard random number stream.

Member functions provide the following operations:

```

void set_seed (long seed) – set the seed of the RNG, if seed == 0, the seed is set according to the current time
and a counter

long seed (void) – return the current seed

long next (void) – return the next random number as an integer on [0, MAXINT]

double next_double (void) – return the next random number on [0, 1]

void reset_start_substream (void) – reset the stream to the beginning of the current substream

void reset_next_substream (void) – advance to the next substream

int uniform (int k) – return an integer sampled from a uniform distribution on [0, k-1]

double uniform (double r) – return a number sampled from a uniform distribution on [0, r]

```

double uniform (double a, double b) – return a number sampled from a uniform distribution on [a, b]

double exponential (void) – return a number sampled from an exponential distribution with mean 1.0

double exponential (double k) – return a number sampled from an exponential distribution with mean k

double normal (double avg, double std) – return a number sampled from a normal distribution with the given average and standard deviation

double lognormal (double avg, double std) – return a number sampled from a lognormal distribution with the given average and standard deviation

### Example

```
/* create new RNGs */
RNG arrival (23456);
RNG size;

/* set the RNGs to the appropriate substream */
for (int i = 1; i < 3; i++) {
    arrival.reset_next_substream();
    size.reset_next_substream();
}

/* print the first 5 arrival times and sizes */
for (int j = 0; j < 5; j++) {
    printf ("%8.3f  %-4d\n", arrival.lognormal(5, 0.1),
            int(size.normal(500, 10)));
}
```

### Output

```
161.826  506
160.591  503
157.145  509
137.715  507
118.573  496
```

## 24.2 Random Variables

The class `RandomVariable` provides a thin layer of functionality on top of the base random number generator and the default random number stream. It is defined in `~ns/ranvar.h`:

```
class RandomVariable : public TclObject {
public:
    virtual double value() = 0;
    int command(int argc, const char*const* argv);
    RandomVariable();
```

```
protected:
    RNG* rng_;
};
```

Classes derived from this abstract class implement specific distributions. Each distribution is parameterized with the values of appropriate parameters. The value method is used to return a value from the distribution.

The currently defined distributions, and their associated parameters are:

```
class UniformRandomVariable min_, max_
class ExponentialRandomVariable avg_
class ParetoRandomVariable avg_, shape_
class ParetoIIRandomVariable avg_, shape_
class ConstantRandomVariable val_
class HyperExponentialRandomVariable avg_, cov_
class NormalRandomVariable avg_, std_
class LogNormalRandomVariable avg_, std_
```

The RandomVariable class is available in OTcl. For instance, to create a random variable that generates number uniformly on [10, 20]:

```
set u [new RandomVariable/Uniform]
$u set min_ 10
$u set max_ 20
$u value
```

By default, RandomVariable objects use the default random number generator described in the previous section. The use-rng method can be used to associate a RandomVariable with a non-default RNG:

```
set rng [new RNG]
$rng seed 0

set e [new RandomVariable/Exponential]
$e use-rng $rng
```

## 24.3 Integrals

The class Integrator supports the approximation of (continuous) integration by (discrete) sums; it is defined in *~ns/integrator.h* as

From integrator.h:

```
class Integrator : public TclObject {
public:
    Integrator();
    void set(double x, double y);
    void newPoint(double x, double y);
    int command(int argc, const char*const* argv);
```



```
protected:
    double lastx_;
    double lasty_;
    double sum_;
};
```

From integrator.cc:

```
Integrator::Integrator() : lastx_(0.), lasty_(0.), sum_(0.)
{
    bind("lastx_", &lastx_);
    bind("lasty_", &lasty_);
    bind("sum_", &sum_);
}

void Integrator::set(double x, double y)
{
    lastx_ = x;
    lasty_ = y;
    sum_ = 0.;
}

void Integrator::newPoint(double x, double y)
{
    sum_ += (x - lastx_) * lasty_;
    lastx_ = x;
    lasty_ = y;
}

int Integrator::command(int argc, const char*const* argv)
{
    if (argc == 4) {
        if (strcmp(argv[1], "newpoint") == 0) {
            double x = atof(argv[2]);
            double y = atof(argv[3]);
            newPoint(x, y);
            return (TCL_OK);
        }
    }
    return (TclObject::command(argc, argv));
}
```

This class provides a base class used by other classes such as QueueMonitor that keep running sums. Each new element of the running sum is added by the newPoint(x, y) function. After the  $k$ th execution of newPoint, the running sum is equal to  $\sum_{i=1}^k y_{i-1}(x_i - x_{i-1})$  where  $x_0 = y_0 = 0$  unless lastx\_, lasty\_, or sum\_ are reset via OTcl. Note that a new point in the sum can be added either by the C++ member newPoint or the OTcl member newpoint. The use of integrals to compute certain types of averages (e.g. mean queue lengths) is given in (pp. 429–430, [15]).

## 24.4 ns-random

**ns-random** is an obsolete way to generate random numbers. This information is provided only for backward compatibility.

`ns-random` is implemented in `~ns/misc.{cc,h}`. When called with no argument, it generates a random number with uniform distribution between 0 and `MAXINT`. When an integer argument is provided, it seeds the random generator with the given number. A special case is when `ns-random 0` is called, it randomly seeds the generator based on current time. This feature is useful to produce non-deterministic results across runs.

## 24.5 Some mathematical-support related objects

**INTEGRATOR OBJECTS** Integrator Objects support the approximate computation of continuous integrals using discrete sums. The running sum(integral) is computed as: `sum_ += [lasty_ * (x - lastx_)]` where (x, y) is the last element entered and (lastx\_, lasty\_) was the element previous to that added to the sum. lastx\_ and lasty\_ are updated as new elements are added. The first sample point defaults to (0,0) that can be changed by changing the values of (lastx\_,lasty\_). `$integrator newpoint <x> <y>`

Add the point (x,y) to the sum. Note that it does not make sense for x to be less than lastx\_.

There are no configuration parameters specific to this object.

State Variables are:

**lastx\_** x-coordinate of the last sample point.

**lasty\_** y-coordinate of the last sample point.

**sum\_** Running sum (i.e. the integral) of the sample points.

**SAMPLES OBJECT** Samples Objects support the computation of mean and variance statistics for a given sample.

`$samples mean`

Returns mean of the sample.

`$samples variance`

Returns variance of the sample.

`$samples cnt`

Returns a count of the sample points considered.

`$samples reset`

Reset the Samples object to monitor a fresh set of samples.

There are no configuration parameters or state variables specific to this object.

## 24.6 Commands at a glance

Following is a list of mathematical support related commands commonly used in simulation scripts:

`set rng [new RNG]`

This creates a new random number generator.

`$rng seed <0 or n>`

This command seeds the RNG. If 0 is specified, the RNG is seeded heuristically. Otherwise the RNG is seeded with the value <n>.

`$rng next-random`

This returns the next random number from RNG.

`$rng uniform <a> <b>`

This returns a number uniformly distributed on <a> and <b>.

`$rng integer <k>`

This returns an integer uniformly distributed on 0 and k-1.

`$rng exponential`

This returns a number that has exponential distribution with average 1.

`set rv [new Randomvariable/<type of random-variable>]`

This creates an instance of a random variable object that generates random variables with specific distribution. The different types of random variables derived from the base class are:

RandomVariable/Uniform, RandomVariable/Exponential, RandomVariable/Pareto, RandomVariable/Constant, RandomVariable/HyperExponential. Each of these distribution types are parameterized with values of appropriate parameters. For details see section 24.2 of this chapter.

`$rv use-rng <rng>`

This method is used to associated a random variable object with a non-default RNG. Otherwise by default, the random variable object is associated with the default random number generator.

## Chapter 25

# Trace and Monitoring Support

The procedures and functions described in this chapter can be found in `~ns/trace.{cc, h}`, `~ns/tcl/lib/ns-trace.tcl`, `~ns/queue-monitor.{cc, h}`, `~ns/tcl/lib/ns-link.tcl`, `~ns/packet.h`, `~ns/flowmon.cc`, and `~ns/classifier-hash.cc`.

There are a number of ways of collecting output or trace data on a simulation. Generally, trace data is either displayed directly during execution of the simulation, or (more commonly) stored in a file to be post-processed and analyzed. There are two primary but distinct types of monitoring capabilities currently supported by the simulator. The first, called *traces*, record each individual packet as it arrives, departs, or is dropped at a link or queue. Trace objects are configured into a simulation as nodes in the network topology, usually with a Tcl “Channel” object hooked to them, representing the destination of collected data (typically a trace file in the current directory). The other types of objects, called *monitors*, record counts of various interesting quantities such as packet and byte arrivals, departures, etc. Monitors can monitor counts associated with all packets, or on a per-flow basis using a *flow monitor* below (Section 25.7).

To support traces, there is a special *common* header included in each packet (this format is defined in `~ns/packet.h` as `hdr_cmn`). It presently includes a unique identifier on each packet, a packet type field (set by agents when they generate packets), a packet size field (in bytes, used to determine the transmission time for packets), and an interface label (used for computing multicast distribution trees).

Monitors are supported by a separate set of objects that are created and inserted into the network topology around queues. They provide a place where arrival statistics and times are gathered and make use of the `class Integrator` (Section 24.3) to compute statistics over time intervals.

### 25.1 Trace Support

The trace support in OTcl consists of a number of specialized classes visible in OTcl but implemented in C++, combined with a set of Tcl helper procedures and classes defined in the *ns* library.

All following OTcl classes are supported by underlying C++ classes defined in `~ns/trace.cc`. Objects of the following types are inserted directly in-line in the network topology:

Trace/Hop	trace a “hop” (XXX what does this mean exactly; it is not really used XXX)
Trace/Enque	a packet arrival (usually at a queue)
Trace/Deque	a packet departure (usually at a queue)
Trace/Drop	packet drop (packet delivered to drop-target)
Trace/Recv	packet receive event at the destination node of a link
SnoopQueue/In	on input, collect a time/size sample (pass packet on)
SnoopQueue/Out	on output, collect a time/size sample (pass packet on)
SnoopQueue/Drop	on drop, collect a time/size sample (pass packet on)
SnoopQueue/EDrop	on an “early” drop, collect a time/size sample (pass packet on)

Objects of the following types are added in the simulation and are referenced by the objects listed above. They are used to aggregate statistics collected by the SnoopQueue objects:

QueueMonitor	receive and aggregate collected samples from snoopers
QueueMonitor/ED	queue-monitor capable of distinguishing between “early” and standard packet drops
QueueMonitor/ED/Flowmon	per-flow statistics monitor (manager)
QueueMonitor/ED/Flow	per-flow statistics container
QueueMonitor/Compat	a replacement for a standard QueueMonitor when <i>ns</i> v1 compatibility is in use

### 25.1.1 OTcl Helper Functions

The following helper functions may be used within simulation scripts to help in attaching trace elements (see `~ns/tcl/lib/ns-lib.tcl`); they are instance procedures of the class Simulator:

<code>flush-trace { }</code>	flush buffers for all trace objects in simulation
<code>create-trace { type file src dst }</code>	create a trace object of type <i>type</i> between the given <i>src</i> and <i>dst</i> nodes. If <i>file</i> is non-null, it is interpreted as a Tcl channel and is attached to the newly-created trace object. The procedure returns the handle to the newly created trace object.
<code>trace-queue { n1 n2 file }</code>	arrange for tracing on the link between nodes <i>n1</i> and <i>n2</i> . This function calls <code>create-trace</code> , so the same rules apply with respect to the <i>file</i> argument.
<code>trace-callback{ ns command }</code>	arranges to call <i>command</i> when a line is to be traced. The procedure treats <i>command</i> as a string and evaluates it for every line traced. See <code>~ns/tcl/ex/callback_demo.tcl</code> for additional details on usage.
<code>monitor-queue { n1 n2 }</code>	this function calls the <code>init-monitor</code> function on the link between nodes <i>n1</i> and <i>n2</i> .
<code>drop-trace { n1 n2 trace }</code>	the given <i>trace</i> object is made the drop-target of the queue associated with the link between nodes <i>n1</i> and <i>n2</i> .

The `create-trace{ }` procedure is used to create a new Trace object of the appropriate kind and attach an Tcl I/O channel to it (typically a file handle). The `src_` and `dst_` fields are used by the underlying C++ object for producing the trace output file so that trace output can include the node addresses defining the endpoints of the link which is being traced. Note that they are not used for *matching*. Specifically, these values in no way relate to the packet header `src` and `dst` fields, which are also displayed when tracing. See the description of the Trace class below (Section 25.3).

The `trace-queue` function enables Enque, Deque, and Drop tracing on the link between nodes *n1* and *n2*. The Link trace procedure is described below (Section 25.2).

The `monitor-queue` function is constructed similarly to `trace-queue`. By calling the link's `init-monitor` procedure, it arranges for the creation of objects (`SnoopQueue` and `QueueMonitor` objects) which can, in turn, be used to ascertain time-aggregated queue statistics.

The `drop-trace` function provides a way to specify a `Queue`'s drop target without having a direct handle of the queue.

## 25.2 Library support and examples

The Simulator procedures described above require the `trace` and `init-monitor` methods associated with the `OTcl Link` class. Several subclasses of link are defined, the most common of which is called `SimpleLink`. Thus, the `trace` and `init-monitor` methods are actually part of the `SimpleLink` class rather than the `Link` base class. The `trace` function is defined as follows (in `ns-link.tcl`):

```
#
# Build trace objects for this link and
# update the object linkage
#
SimpleLink instproc trace { ns f } {
    $self instvar enqT_ deqT_ drpT_ queue_ link_ head_ fromNode_ toNode_
    $self instvar drophead_

    set enqT_ [$ns create-trace Enque $f $fromNode_ $toNode_]
    set deqT_ [$ns create-trace Deque $f $fromNode_ $toNode_]
    set drpT_ [$ns create-trace Drop $f $fromNode_ $toNode_]

    $drpT_ target [$drophead_ target]
    $drophead_ target $drpT_
    $queue_ drop-target $drpT_

    $deqT_ target [$queue_ target]
    $queue_ target $deqT_

    if { [$head_ info class] == "networkinterface" } {
        $enqT_ target [$head_ target]
        $head_ target $enqT_
        # puts "head is i/f"
    } else {
        $enqT_ target $head_
        set head_ $enqT_
        # puts "head is not i/f"
    }
}
```

This function establishes `Enque`, `Deque`, and `Drop` traces in the simulator `$ns` and directs their output to I/O handle `$f`. The function assumes a queue has been associated with the link. It operates by first creating three new trace objects and inserting the `Enque` object before the queue, the `Deque` object after the queue, and the `Drop` object between the queue and its previous drop target. Note that all trace output is directed to the same I/O handle.

This function performs one other additional task. It checks to see if a link contains a network interface, and if so, leaves it as the first object in the chain of objects in the link, but otherwise inserts the `Enque` object as the first one.

The following functions, `init-monitor` and `attach-monitor`, are used to create a set of objects used to monitor queue sizes of a queue associated with a link. They are defined as follows:

```
SimpleLink instproc attach-monitors { insnoop outsnoop dropsnoop qmon } {
    $self instvar queue_ head_ snoopIn_ snoopOut_ snoopDrop_
    $self instvar drophead_ qMonitor_

    set snoopIn_ $insnoop
    set snoopOut_ $outsnoop
    set snoopDrop_ $dropsnoop

    $snoopIn_ target $head_
    set head_ $snoopIn_

    $snoopOut_ target [$queue_ target]
    $queue_ target $snoopOut_

    $snoopDrop_ target [$drophead_ target]
    $drophead_ target $snoopDrop_

    $snoopIn_ set-monitor $qmon
    $snoopOut_ set-monitor $qmon
    $snoopDrop_ set-monitor $qmon
    set qMonitor_ $qmon
}

#
# Insert objects that allow us to monitor the queue size
# of this link. Return the name of the object that
# can be queried to determine the average queue size.
#
SimpleLink instproc init-monitor { ns qtrace sampleInterval } {
    $self instvar qMonitor_ ns_ qtrace_ sampleInterval_

    set ns_ $ns
    set qtrace_ $qtrace
    set sampleInterval_ $sampleInterval
    set qMonitor_ [new QueueMonitor]

    $self attach-monitors [new SnoopQueue/In] \
        [new SnoopQueue/Out] [new SnoopQueue/Drop] $qMonitor_

    set bytesInt_ [new Integrator]
    $qMonitor_ set-bytes-integrator $bytesInt_
    set pktsInt_ [new Integrator]
    $qMonitor_ set-pkts-integrator $pktsInt_
    return $qMonitor_
}
```

These functions establish queue monitoring on the `SimpleLink` object in the simulator `ns`. Queue monitoring is implemented by constructing three `SnoopQueue` objects and one `QueueMonitor` object. The `SnoopQueue` objects are linked in around a `Queue` in a way similar to `Trace` objects. The `SnoopQueue/In(Out)` object monitors packet arrivals(departures) and reports them to an associated `QueueMonitor` agent. In addition, a `SnoopQueue/Out` object is

also used to accumulate packet drop statistics to an associated QueueMonitor object. For init-monitor the same QueueMonitor object is used in all cases. The C++ definitions of the SnoopQueue and QueueMonitor classes are described below.

## 25.3 The C++ Trace Class

Underlying C++ objects are created in support of the interface specified in Section 25.3 and are linked into the network topology as network elements. The single C++ Trace class is used to implement the OTcl classes Trace/Hop, Trace/Enque, Trace/Deque, and Trace/Drop. The type\_ field is used to differentiate among the various types of traces any particular Trace object might implement. Currently, this field may contain one of the following symbolic characters: **e** for enqueue, **d** for deque, **h** for hop, and **l** for drop. The overall class is defined as follows in `~ns/trace.cc`:

```
class Trace : public Connector {
protected:
    int type_;
    nsaddr_t src_;
    nsaddr_t dst_;
    Tcl_Channel channel_;
    int callback_;
    char wrk_[256];
    void format(int tt, int s, int d, Packet* p);
    void annotate(const char* s);
    int show_tcphdr_; // bool flags; backward compat
public:
    Trace(int type);
    ~Trace();
    int command(int argc, const char*const* argv);
    void recv(Packet* p, Handler*);
    void dump();
    inline char* buffer() { return (wrk_); }
};
```

The src\_, and dst\_ internal state is used to label trace output and is independent of the corresponding field names in packet headers. The main recv() method is defined as follows:

```
void Trace::recv(Packet* p, Handler* h)
{
    format(type_, src_, dst_, p);
    dump();
    /* hack: if trace object not attached to anything free packet */
    if (target_ == 0)
        Packet::free(p);
    else
        send(p, h); /* Connector::send() */
}
```

The function merely formats a trace entry using the source, destination, and particular trace type character. The dump function writes the formatted entry out to the I/O handle associated with channel\_. The format function, in effect, dictates the trace file format.



## 25.4 Trace File Format

The `Trace::format()` method defines the trace file format used in trace files produced by the `Trace` class. It is constructed to maintain backward compatibility with output files in earlier versions of the simulator (*i.e.*, `ns v1`) so that `ns v1` post-processing scripts continue to operate. The important pieces of its implementation are as follows:

```
// this function should retain some backward-compatibility, so that
// scripts don't break.
void Trace::format(int tt, int s, int d, Packet* p)
{
    hdr_cmh *th = (hdr_cmh*)p->access(off_cmh_);
    hdr_ip *iph = (hdr_ip*)p->access(off_ip_);
    hdr_tcp *tcph = (hdr_tcp*)p->access(off_tcp_);
    hdr_rtp *rh = (hdr_rtp*)p->access(off_rtp_);
    packet_t t = th->ptype();
    const char* name = packet_info.name(t);

    if (name == 0)
        abort();

    int seqno;
    /* XXX */
    /* CBR's now have seqno's too */
    if (t == PT_RTP || t == PT_CBR)
        seqno = rh->seqno();
    else if (t == PT_TCP || t == PT_ACK)
        seqno = tcph->seqno();
    else
        seqno = -1;

    ...

    if (!show_tcphdr_) {
        sprintf(wrk_, "%c %g %d %d %s %d %s %d %d.%d %d.%d %d %d",
            tt,
            Scheduler::instance().clock(),
            s,
            d,
            name,
            th->size(),
            flags,
            iph->flowid() /* was p->class_ */,
            iph->src() >> 8, iph->src() & 0xff, // XXX
            iph->dst() >> 8, iph->dst() & 0xff, // XXX
            seqno,
            th->uid() /* was p->uid_ */);
    } else {
        sprintf(wrk_,
            "%c %g %d %d %s %d %s %d %d.%d %d.%d %d %d %d 0x%x %d",
            tt,
            Scheduler::instance().clock(),
            s,
            d,
```

```

        name,
        th->size(),
        flags,
        iph->flowid() /* was p->class_ */,
        iph->src() >> 8, iph->src() & 0xff,      // XXX
        iph->dst() >> 8, iph->dst() & 0xff,      // XXX
        seqno,
        th->uid(), /* was p->uid_ */
        tcph->ackno(),
        tcph->flags(),
        tcph->hlen());
    }

```

This function is somewhat unelegant, primarily due to the desire to maintain backward compatibility. It formats the source, destination, and type fields defined in the trace object (*not in the packet headers*), the current time, along with various packet header fields including, type of packet (as a name), size, flags (symbolically), flow identifier, source and destination packet header fields, sequence number (if present), and unique identifier. The `show_tcphdr_` variable indicates whether the trace output should append tcp header information (ack number, flags, header length) at the end of each output line. This is especially useful for simulations using FullTCP agents (Section 33.3). An example of a trace file (without the tcp header fields) might appear as follows:

```

+ 1.84375 0 2 cbr 210 ----- 0 0.0 3.1 225 610
- 1.84375 0 2 cbr 210 ----- 0 0.0 3.1 225 610
r 1.84471 2 1 cbr 210 ----- 1 3.0 1.0 195 600
r 1.84566 2 0 ack 40 ----- 2 3.2 0.1 82 602
+ 1.84566 0 2 tcp 1000 ----- 2 0.1 3.2 102 611
- 1.84566 0 2 tcp 1000 ----- 2 0.1 3.2 102 611
r 1.84609 0 2 cbr 210 ----- 0 0.0 3.1 225 610
+ 1.84609 2 3 cbr 210 ----- 0 0.0 3.1 225 610
d 1.84609 2 3 cbr 210 ----- 0 0.0 3.1 225 610
- 1.8461 2 3 cbr 210 ----- 0 0.0 3.1 192 511
r 1.84612 3 2 cbr 210 ----- 1 3.0 1.0 196 603
+ 1.84612 2 1 cbr 210 ----- 1 3.0 1.0 196 603
- 1.84612 2 1 cbr 210 ----- 1 3.0 1.0 196 603
+ 1.84625 3 2 cbr 210 ----- 1 3.0 1.0 199 612

```

Here we see 14 trace entries, five enqueue operations (indicated by “+” in the first column), four deque operations (indicated by “-”), four receive events (indicated by “r”), and one drop event. (this had better be a trace fragment, or some packets would have just vanished!). The simulated time (in seconds) at which each event occurred is listed in the second column. The next two fields indicate between which two nodes tracing is happening. The next field is a descriptive name for the type of packet seen (Section 25.5). The next field is the packet’s size, as encoded in its IP header.

The next field contains the flags, which not used in this example. The flags are defined in the `flags[]` array in `trace.cc`. Four of the flags are used for ECN: “E” for Congestion Experienced (CE) and “N” for ECN-Capable-Transport (ECT) indications in the IP header, and “C” for ECN-Echo and “A” for Congestion Window Reduced (CWR) in the TCP header. For the other flags, “P” is for priority, and “F” is for TCP Fast Start.

The next field gives the *IP flow identifier* field as defined for IP version 6. The subsequent two fields indicate the packet’s source and destination node addresses, respectively. The following field indicates the sequence number.<sup>2</sup> The last field is a

<sup>1</sup>In *ns v1*, each packet included a `class` field, which was used by CBQ to classify packets. It then found additional use to differentiate between “flows” at one trace point. In *ns v2*, the flow ID field is available for this purpose, but any additional information (which was commonly overloaded into the class field in *ns v1*) should be placed in its own separate field, possibly in some other header

<sup>2</sup>In *ns v1*, all packets contained a sequence number, whereas in *ns v2* only those Agents interested in providing sequencing will generate sequence numbers. Thus, this field may not be useful in *ns v2* for packets generated by agents that have not filled in a sequence number. It is used here to remain backward compatible with *ns v1*.

unique packet identifier. Each new packet created in the simulation is assigned a new, unique identifier.

## 25.5 Packet Types

Each packet contains a packet type field used by `Trace::format` to print out the type of packet encountered. The type field is defined in the `TraceHeader` class, and is considered to be part of the trace support; it is not interpreted elsewhere in the simulator. Initialization of the type field in packets is performed by the `Agent::allocpkt(void)` method. The type field is set to integer values associated with the definition passed to the `Agent` constructor (Section 10.6.3). The currently-supported definitions, their values, and their associated symbolic names are as follows (defined in `~ns/packet.h`):

```
enum packet_t {
PT_TCP,
PT_UDP,
PT_CBR,
PT_AUDIO,
PT_VIDEO,
PT_ACK,
PT_START,
PT_STOP,
PT_PRUNE,
PT_GRAFT,
PT_GRAFTACK,
PT_JOIN,
PT_ASSERT,
PT_MESSAGE,
PT_RTCP,
PT_RTP,
PT_RTPROTO_DV,
PT_CtrMcast_Encap,
PT_CtrMcast-Decap,
PT_SRM,
/* simple signalling messages */
PT_REQUEST,
PT_ACCEPT,
PT_CONFIRM,
PT_TEARDOWN,
PT_LIVE, // packet from live network
PT_REJECT,

PT_TELNET, // not needed: telnet use TCP
PT_FTP,
PT_PARETO,
PT_EXP,
PT_INVALID,
PT_HTTP,
/* new encapsulator */
PT_ENCAPSULATED,
PT_MFTP,
/* CMU/Monarch's extensions */
PT_ARP,
PT_MAC,
```

```

PT_TORA,
PT_DSR,
PT_AODV,

// insert new packet types here

PT_NTTYPE // This MUST be the LAST one
};

```

The constructor of class `p_info` glues these constants with their string values:

```

p_info() {
name_[PT_TCP]= "tcp";
name_[PT_UDP]= "udp";
name_[PT_CBR]= "cbr";
name_[PT_AUDIO]= "audio";
...
name_[PT_NTTYPE]= "undefined";
}

```

See also section 12.2.2 for more details.

## 25.6 Queue Monitoring

Queue monitoring refers to the capability of tracking the dynamics of packets at a queue (or other object). A queue monitor tracks packet arrival/departure/drop statistics, and may optionally compute averages of these values. Monitoring may be applied to all packets (aggregate statistics), or per-flow statistics (using a Flow Monitor).

Several classes are used in supporting queue monitoring. When a packet arrives at a link where queue monitoring is enabled, it generally passes through a `SnoopQueue` object when it arrives and leaves (or is dropped). These objects contain a reference to a `QueueMonitor` object.

A `QueueMonitor` is defined as follows (*~ns/queue-monitor.cc*):

```

class QueueMonitor : public TclObject {
public:
    QueueMonitor() : bytesInt_(NULL), pktsInt_(NULL), delaySamp_(NULL),
        size_(0), pkts_(0),
        parrivals_(0), barrivals_(0),
        pdepartures_(0), bdepartures_(0),
        pdrops_(0), bdrops_(0),
        srcId_(0), dstId_(0), channel_(0) {

        bind("size_", &size_);
        bind("pkts_", &pkts_);
        bind("parrivals_", &parrivals_);
        bind("barrivals_", &barrivals_);
        bind("pdepartures_", &pdepartures_);
        bind("bdepartures_", &bdepartures_);
    }

```

```

        bind("pdrops_", &pdrops_);
        bind("bdrops_", &bdrops_);
        bind("off_cmn_", &off_cmn_);
    };

    int size() const { return (size_); }
    int pkts() const { return (pkts_); }
    int parrivals() const { return (parrivals_); }
    int barrivals() const { return (barrivals_); }
    int pdepartures() const { return (pdepartures_); }
    int bdepartures() const { return (bdepartures_); }
    int pdrops() const { return (pdrops_); }
    int bdrops() const { return (bdrops_); }
    void printStats();
    virtual void in(Packet*);
    virtual void out(Packet*);
    virtual void drop(Packet*);
    virtual void edrop(Packet*) { abort(); }; // not here
    virtual int command(int argc, const char*const* argv);
    ...

// packet arrival to a queue
void QueueMonitor::in(Packet* p)
{
    hdr_cmn* hdr = (hdr_cmn*)p->access(off_cmn_);
    double now = Scheduler::instance().clock();
    int pktsz = hdr->size();

    barrivals_ += pktsz;
    parrivals_++;
    size_ += pktsz;
    pkts_++;
    if (bytesInt_)
        bytesInt_->newPoint(now, double(size_));
    if (pktsInt_)
        pktsInt_->newPoint(now, double(pkts_));
    if (delaySamp_)
        hdr->timestamp() = now;
    if (channel_)
        printStats();
}

... in(), out(), drop() are all defined similarly ...

```

In addition to the packet and byte counters, a queue monitor may optionally refer to objects that keep an integral of the queue size over time using `Integrator` objects, which are defined in Section 24.3. The `Integrator` class provides a simple implementation of integral approximation by discrete sums.

All bound variables beginning with **p** refer to packet counts, and all variables beginning with **b** refer to byte counts. The variable `size_` records the instantaneous queue size in bytes, and the variable `pkts_` records the same value in packets. When a `QueueMonitor` is configured to include the integral functions (on bytes or packets or both), it computes the approximate integral of the queue size (in bytes) with respect to time over the interval  $[t_0, now]$ , where  $t_0$  is either the start of the simulation or the last time the `sum_` field of the underlying `Integrator` class was reset.

The `QueueMonitor` class is not derived from `Connector`, and is not linked directly into the network topology. Rather, objects of the `SnoopQueue` class (or its derived classes) are inserted into the network topology, and these objects contain references to an associated queue monitor. Ordinarily, multiple `SnoopQueue` objects will refer to the same queue monitor. Objects constructed out of these classes are linked in the simulation topology as described above and call `QueueMonitor` out, in, or drop procedures, depending on the particular type of snoopy queue.

## 25.7 Per-Flow Monitoring

A collection of specialized classes are used to implement per-flow statistics gathering. These classes include: `QueueMonitor/ED/Flowmon`, `QueueMonitor/ED/Flow`, and `Classifier/Hash`. Typically, an arriving packet is inspected to determine to which flow it belongs. This inspection and flow mapping is performed by a *classifier* object (described in section 25.7.1). Once the correct flow is determined, the packet is passed to a *flow monitor*, which is responsible for collecting per-flow state. Per-flow state is contained in *flow* objects in a one-to-one relationship to the flows known by the flow monitor. Typically, a flow monitor will create flow objects on-demand when packets arrive that cannot be mapped to an already-known flow.

### 25.7.1 The Flow Monitor

The `QueueMonitor/ED/Flowmon` class is responsible for managing the creation of new flow objects when packets arrive on previously unknown flows and for updating existing flow objects. Because it is a subclass of `QueueMonitor`, each flow monitor contains an aggregate count of packet and byte arrivals, departures, and drops. Thus, it is not necessary to create a separate queue monitor to record aggregate statistics. It provides the following OTcl interface:

<code>classifier</code>	get(set) classifier to map packets to flows
<code>attach</code>	attach a Tcl I/O channel to this monitor
<code>dump</code>	dump contents of flow monitor to Tcl channel
<code>flows</code>	return string of flow object names known to this monitor

The `classifier` function sets or gets the name of the previously-allocated object which will perform packet-to-flow mapping for the flow monitor. Typically, the type of classifier used will have to do with the notion of “flow” held by the user. One of the hash based classifiers that inspect various IP-level header fields is typically used here (e.g. `fi d`, `src/dst`, `src/dst/fi d`). Note that while classifiers usually receive packets and forward them on to downstream objects, the flow monitor uses the classifier only for its packet mapping capability, so the flow monitor acts as a passive monitor only and does not actively forward packets.

The `attach` and `dump` functions are used to associate a Tcl I/O stream with the flow monitor, and dump its contents on-demand. The file format used by the `dump` command is described below.

The `flows` function returns a list of the names of flows known by the flow monitor in a way understandable to Tcl. This allows tcl code to interrogate a flow monitor in order to obtain handles to the individual flows it maintains.

### 25.7.2 Flow Monitor Trace Format

The flow monitor defines a trace format which may be used by post-processing scripts to determine various counts on a per-flow basis. The format is defined by the following code in `~ns/flowmon.cc`:



`$ns_ trace-all <tracefile>`

This is the command used to setup tracing in ns. All traces are written in the `<tracefile>`.

`$ns_ namtrace-all <namtracefile>`

This command sets up nam tracing in ns. All nam traces are written in to the `<namtracefile>`.

`$ns_ namtrace-all-wireless <namtracefile> <X> <Y>`

This command sets up wireless nam tracing. `<X>` and `<Y>` are the x-y co-ordinates for the wireless topology and all wireless nam traces are written into the `<namtracefile>`.

`$ns_ nam-end-wireless <stoptime>`

This tells nam the simulation stop time given in `<stoptime>`.

`$ns_ trace-all-satlinks <tracefile>`

This is a method to trace satellite links and write traces into `<tracefile>`.

`$ns_ flush-trace`

This command flushes the trace buffer and is typically called before the simulation run ends.

`$ns_ get-nam-traceall`

Returns the namtrace file descriptor stored as the Simulator instance variable called `namtraceAllFile_`.

`$ns_ get-ns-traceall`

Similar to `get-nam-traceall`. This returns the file descriptor for ns tracefile which is stored as the Simulator instance called `traceAllFile_`.

`$ns_ create-trace <type> <file> <src> <dst> <optional:op>`

This command creates a trace object of type `<type>` between the `<src>` and `<dst>` nodes. The traces are written into the `<file>`. `<op>` is the argument that may be used to specify the type of trace, like nam. if `<op>` is not defined, the default trace object created is for nstraces.

`$ns_ trace-queue <n1> <n2> <optional:file>`

This is a wrapper method for `create-trace`. This command creates a trace object for tracing events on the link represented by the nodes `<n1>` and `<n2>`.

`$ns_ namtrace-queue <n1> <n2> <optional:file>`

This is used to create a trace object for namtracing on the link between nodes `<n1>` and `<n2>`. This method is very similar to and is the namtrace counterpart of method `trace-queue`.

`$ns_ drop-trace <n1> <n2> <trace>`

This command makes the given `<trace>` object a drop-target for the queue associated with the link between nodes `<n1>` and `<n2>`.

`$ns_ monitor-queue <n1> <n2> <qtrace> <optional:sampleinterval>`

This sets up a monitor that keeps track of average queue length of the queue on the link between nodes `<n1>` and `<n2>`. The default value of `sampleinterval` is 0.1.

`$link trace-dynamics <ns> <fileID>` Trace the dynamics of this link and write the output to `fileID` filehandle. `ns` is an instance of the Simulator or MultiSim object that was created to invoke the simulation.

The tracefile format is backward compatible with the output files in the ns version 1 simulator so that ns-1 postprocessing scripts can still be used. Trace records of traffic for link objects with Enqueue, Dequeue, receive or Drop Tracing have the following form:

`<code> <time> <hsrc> <hdst> <packet>`



where

```
<code> := [hd+-] h=hop d=drop +=enqueue -=dequeue r=receive <time> :=  
simulation time in seconds  
<hsrc> := first node address of hop/queuing link  
<hdst> := second node address of hop/queuing link  
<packet> := <type> <size> <flags> <flowID> <src.sport> <dst.dport> <seq>  
<pktID>  
<type> := tcp|telnet|cbr|ack etc.  
<size> := packet size in bytes  
<flags> := [CP] C=congestion, P=priority  
<flowID> := flow identifier field as defined for IPv6  
<src.sport> := transport address (src=node,sport=agent)  
<dst.sport> := transport address (dst=node,dport=agent)  
<seq> := packet sequence number  
<pktID> := unique identifier for every new packet
```

Only those agents interested in providing sequencing will generate sequence numbers and hence this field may not be useful for packets generated by some agents. For links that use RED gateways, there are additional trace records as follows:

```
<code> <time> <value>
```

where

```
<code> := [Qap] Q=queue size, a=average queue size, p=packet dropping  
probability  
<time> := simulation time in seconds  
<value> := value
```

Trace records for link dynamics are of the form:

```
<code> <time> <state> <src> <dst>
```

where

```
<code> := [v]  
<time> := simulation time in seconds  
<state> := [link-up | link-down]  
<src> := first node address of link  
<dst> := second node address of link
```

# Chapter 26

## Test Suite Support

The ns distribution contains many test suites under `~ns/tcl/test`, which used by validation programs (`~ns/validate`, `validate-wired`, `validate-wireless`, and `validate.win32`) to verify that the installation of ns is correct. If you modify or add new modules to ns, you are encouraged to run the validation programs to make sure that your changes do not affect other parts in ns.

### 26.1 Test Suite Components

Each test suite under `~ns/tcl/test` is written to verify the correctness of a certain module in ns. It has 3 components:

- A shell script (`test-all-xxx`) to start the test;
- A ns tcl script (`test-suite-xxx.tcl`) to actually run through the tests defined.
- A subdirectory (`test-output-xxx`) under `~ns/tcl/test`, which contains the correct trace files generated by the test suite. These files are used to verify if the test suite runs correctly with your ns.

(Note: xxx stands for the name of the test suite.)

### 26.2 Write a Test Suite

You can take one of the test suites under `~ns/tcl/test` as a template when you are writing your own, for example the test suite written for wireless lan (`test-all-wireless-lan`, `test-suite-wireless-lan.tcl`, and `test-output-wireless-lan`).

To write a test suite, you first need to write the shell script (`test-all-xxx`). In the shell script, you specify the module to be tested, the name of the ns tcl script and the output subdirectory. You can run this shell script in quiet mode. Below is the example (`test-all-wireless-lan`):

```
# To run in quiet mode:  "./test-all-wireless-lan quiet".  
  
f="wireless-lan" # Specify the name of the module to test.
```

```

file="test-suite-$f.tcl" # The name of the ns script.
directory="test-output-$f" # Subdirectory to hold the test results
version="v2" # Specify the ns version.

# Pass the arguments to test-all-templatel, which will run through
# all the test cases defined in test-suite-wireless-lan.tcl.
./test-all-templatel $file $directory $version $@

```

You also need to create several test cases in the ns script (test-suite-xxx.tcl) by defining a subclass of TestSuite for each different test. For example, in test-suite-wireless-lan.tcl, each test case uses a different Ad Hoc routing protocol. They are defined as:

```

Class TestSuite

# wireless model using destination sequence distance vector
Class Test/dsdv -superclass TestSuite

# wireless model using dynamic source routing
Class Test/dsr -superclass TestSuite
... ..

```

Each test case is basically a simulation scenario. In the super class TestSuite, you can define some functions, like init and finish to do the work required by each test case, for example setting up the network topology and ns trace. The test specific configurations are defined within the corresponding sub-class. Each sub-class also has a run function to start the simulation.

```

TestSuite instproc init {} {
    global opt tracefd topo chan prop
    global node_ god_
    $self instvar ns_ testName_
    set ns_ [new Simulator]
    ... ..
}

TestSuite instproc finish {} {
    $self instvar ns_
    global quiet

    $ns_ flush-trace

    puts "finishing.."
    exit 0
}

Test/dsdv instproc init {} {
    global opt node_ god_
    $self instvar ns_ testName_
    set testName_ dsdv
    ... ..

    $self next

```

```

... ..

$ns_ at $opt(stop).1 "$self finish"
}

Test/dsdv instproc run {} {
    $self instvar ns_
    puts "Starting Simulation..."
    $ns_ run
}

```

All the tests are started by the function `runtest` in the `ns` script.

```

proc runtest {arg} {
    global quiet
    set quiet 0

    set b [llength $arg]
    if {$b == 1} {
        set test $arg
    } elseif {$b == 2} {
        set test [lindex $arg 0]
        if {[lindex $arg 1] == "QUIET"} {
            set quiet 1
        }
    } else {
        usage
    }
    set t [new Test/$test]
    $t run
}

global argv arg0
runtest $argv

```

When you run the tests, trace files are generated and saved to the output subdirectory. These trace files are compared to the those correct trace coming with the test suite. If the comparison shows difference, the test is failed.

# Chapter 27

## ns Code Styles

We recommend the following coding guidelines for ns

### 27.1 Indentation style

- We recommend using the BSD Kernel Normal Form coding style located at <http://cvsweb.netbsd.org/bsdweb.cgi/sharesrc/share/misc/style?rev=HEAD&content-type=text/x-cvsweb-markup>
- Although KNF is specified for C, it also applies reasonably well to C++ and Tcl. Most of ns already follows KNF and it is also extensively used for the BSD and Linux kernels.
- The high order bit is 8-space indents. Using 8-space indents avoid confusion about what a "tab" character represents. A downside is it makes deeply nested looping structures hard to fit in 80 columns. (Some people consider this a feature. :-)

### 27.2 Variable Naming Conventions

- Instance variables of a class should all end in an underscore. This helps distinguish instance variables from global and local variables.
- C++ and Tcl bound variables should have the same names This helps identify the bound variables quickly and reduces complexity

### 27.3 Miscellaneous

- Avoid the use of C++ templates. Ns is supported on multiple platforms and templates are not very portable and are often difficult to debug. Exception: This guideline has been relaxed for some imported code, but the core of ns should build and run without templates.
- For NsObjects, use the `debug_` instance variable to enable debugging functionality. This avoids repetitive definitions of debug statements and allows debugging a particular Nsobject without recompilation.

Example: To enable debugging in Queue object include the following statement in your tcl script.

```
Queue set debug_ true
```

Debugging statments can be inserted in the classes inheriting from Queue as follows:

```
debug("This is a debug statement %d",variable_to_debug);
```

All debugging statements are sent to stdout.

- Error messages which cause the program to exit should be sent to stderr. All other messages should be sent to stdout

## **Part IV**

# **Routing**

# Chapter 28

## Unicast Routing

This section describes the structure of unicast routing in *ns*. We begin by describing the interface to the user (Section 28.1), through methods in the class `Simulator` and the class `RouteLogic`. We then describe configuration mechanisms for specialized routing (Section 28.2) such as asymmetric routing, or equal cost multipath routing. The next section describes the configuration mechanisms for individual routing strategies and protocols (Section 28.3). We conclude with a comprehensive look at the internal architecture (Section 28.4) of routing in *ns*.

The procedures and functions described in this chapter can be found in `~ns/tcl/lib/ns-route.tcl`, `~ns/tcl/rtglib/route-proto.tcl`, `~ns/tcl/mcast/McastProto.tcl`, and `~ns/rtProtoDV.{cc, h}`.

### 28.1 The Interface to the Simulation Operator (The API)

The user level simulation script requires one command: to specify the unicast routing strategy or protocols for the simulation. A routing strategy is a general mechanism by which *ns* will compute routes for the simulation. There are four routing strategies in *ns*: Static, Session, Dynamic and Manual. Conversely, a routing protocol is a realization of a specific algorithm. Currently, Static and Session routing use the Dijkstra's all-pairs SPF algorithm []; one type of dynamic routing strategy is currently implemented: the Distributed Bellman-Ford algorithm []. In *ns*, we blur the distinction between strategy and protocol for static and session routing, considering them simply as protocols<sup>1</sup>.

`rtproto{}` is the instance procedure in the class `Simulator` that specifies the unicast routing protocol to be used in the simulation. It takes multiple arguments, the first of which is mandatory; this first argument identifies the routing protocol to be used. Subsequent arguments specify the nodes that will run the instance of this protocol. The default is to run the same routing protocol on all the nodes in the topology. As an example, the following commands illustrate the use of the `rtproto{}` command.

```
$ns rtproto Static           ; # Enable static route strategy for the simulation
$ns rtproto Session         ; # Enable session routing for this simulation
$ns rtproto DV $n1 $n2 $n3  ; # Run DV agents on nodes $n1, $n2, and $n3
$ns rtproto LS $n1 $n2      ; # Run link state routing on specified nodes
```

If a simulation script does not specify any `rtproto{}` command, then *ns* will run Static routing on all the nodes in the topology.

---

<sup>1</sup>The consideration is that static and session routing strategies/protocols are implemented as agents derived from the class `Agent/rtProto` similar to how the different dynamic routing protocols are implemented; hence the blurred distinctions.



Multiple `rtproto{}` lines for the same or different routing protocols can occur in a simulation script. However, a simulation cannot use both centralized routing mechanisms such as static or session routing and detailed dynamic routing protocols such as DV.

In dynamic routing, each node can be running more than one routing protocol. In such situations, more than one routing protocol can have a route to the same destination. Therefore, each protocol affixes a preference value to each of its routes. These values are non-negative integers in the range 0 . . . 255. The lower the value, the more preferred the route. When multiple routing protocol agents have a route to the same destination, the most preferred route is chosen and installed in the node's forwarding tables. If more than one agent has the most preferred routes, the ones with the lowest metric is chosen. We call the least cost route from the most preferred protocol the "candidate" route. If there are multiple candidate routes from the same or different protocols, then, currently, one of the agent's routes is randomly chosen<sup>2</sup>.

**Preference Assignment and Control** Each protocol agent stores an array of route preferences, `rtpref_`. There is one element per destination, indexed by the node handle. The default preference values used by each protocol are derived from a class variable, `preference_`, for that protocol. The current defaults are:

```
Agent/rtProto set preference_ 200 ; # global default preference
Agent/rtProto/Direct3 set preference_ 100
Agent/rtProto/DV set preference_ 120
```

A simulation script can control routing by altering the preference for routes in one of three ways: alter the preference for a specific route learned *via* a particular protocol agent, alter the preference for all routes learned by the agent, or alter the class variables for the agent before the agent is created.

**Link Cost Assignment and Control** In the currently implemented route protocols, the metric of a route to a destination, at a node, is the cost to reach the destination from that node. It is possible to change the link costs at each of the links. The instance procedure `cost{}` is invoked as `$ns cost <node1> <node2> <cost>`, and sets the cost of the link from `<node1>` to `<node2>` to `<cost>`.

```
$ns cost $n1 $n2 10 ; # set cost of link from $n1 to $n2 to 10
$ns cost $n2 $n1 5 ; # set cost of link in reverse direction to 5
[$ns link $n1 $n2] cost? ; # query cost of link from $n1 to $n2
[$ns link $n2 $n1] cost? ; # query cost of link in reverse direction
```

Notice that the procedure sets the cost along one direction only. Similarly, the procedure `cost?{}` returns the cost of traversing the specified unidirectional link. The default cost of a link is 1.

## 28.2 Other Configuration Mechanisms for Specialised Routing

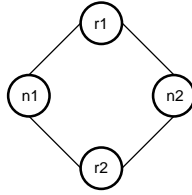
It is possible to adjust preference and cost mechanisms to get two special types of route configurations: asymmetric routing, and multipath routing.

<sup>2</sup>This really is undesirable, and may be fixed at some point. The fix will probably be to favor the agents in class preference order. A user level simulation relying on this behavior, or getting into this situation in specific topologies is not recommended.

<sup>3</sup>Direct is a special routing strategy that is used in conjunction with Dynamic routing. We will describe this in greater detail as part of the route architecture description.

**Asymmetric Routing** Asymmetric routing occurs when the path from node  $n_1$  to node  $n_2$  is different from the path from  $n_2$  to  $n_1$ . The following shows a simple topology, and cost configuration that can achieve such a result:

Nodes  $n_1$  and  $n_2$  use different paths to reach each other. All other pairs of nodes use symmetric paths to reach each other.



```

$ns cost $n1 $r1 2
$ns cost $n2 $r2 2
$ns cost $r1 $n2 3

```

Any routing protocol that uses link costs as the metric can observe such asymmetric routing if the link costs are appropriately configured<sup>4</sup>.

**MultiPath Routing** Each node can be individually configured to use multiple separate paths to a particular destination. The instance variable `multiPath_` determines whether or not that node will use multiple paths to any destination. Each node initialises its instance variable from a class variable of the same name. If multiple candidate routes to a destination are available, all of which are learned through the same protocol, then that node can use all of the different routes to the destination simultaneously. A typical configuration is as shown below:

```

Node set multiPath_ 1 ; # All new nodes in the simulation use multiPaths where applicable
or alternately
set n1 [$ns Node] ; # only enable $n1 to use multiPaths where applicable
$n1 set multiPath_ 1

```

Currently, only DV routing can generate multipath routes.

## 28.3 Protocol Specific Configuration Parameters

**Static Routing** The static route computation strategy is the default route computation mechanism in *ns*. This strategy uses the Dijkstra's all-pairs SPF algorithm []. The route computation algorithm is run exactly once prior to the start of the simulation. The routes are computed using an adjacency matrix and link costs of all the links in the topology.

(Note that static routing is static in the sense that it is computed once when the simulation starts, as opposed to session and DV routing that allow routes to change mid-simulation. An alternative to static routing is Manual routing where routes are not computed but instead are set (manually) by the user.)

**Session Routing** The static routing strategy described earlier only computes routes for the topology once in the course of a simulation. If the above static routing is used and the topology changes while the simulation is in progress, some sources and destinations may become temporarily unreachable from each other for a short time.

Session routing strategy is almost identical to static routing, in that it runs the Dijkstra all-pairs SPF algorithm prior to the start of the simulation, using the adjacency matrix and link costs of the links in the topology. However, it will also run the same algorithm to recompute routes in the event that the topology changes during the course of a simulation. In other words, route recomputation and recovery is done instantaneously and there will not be transient routing outage as in static routing.

Session routing provides complete and instantaneous routing changes in the presence of topology dynamics. If the topology is always connected, there is end-to-end connectivity at all times during the course of the simulation. However, the user should

<sup>4</sup>Link costs can also be used to favour or disregard specific links in order to achieve particular topology configurations.

note that the instantaneous route recomputation of session routing does not prevent temporary violations of causality, such as packet reordering, around the instant that the topology changes.

**DV Routing** DV routing is the implementation of Distributed Bellman-Ford (or Distance Vector) routing in *ns*. The implementation sends periodic route updates every `advertInterval`. This variable is a class variable in the class `Agent/rtProto/DV`. Its default value is 2 seconds.

In addition to periodic updates, each agent also sends triggered updates; it does this whenever the forwarding tables in the node change. This occurs either due to changes in the topology, or because an agent at the node received a route update, and recomputed and installed new routes.

Each agent employs the split horizon with poisoned reverse mechanisms to advertise its routes to adjacent peers. “Split horizon” is the mechanism by which an agent will not advertise the route to a destination out of the interface that it is using to reach that destination. In a “Split horizon with poisoned reverse” mechanism, the agent will advertise that route out of that interface with a metric of infinity.

Each DV agent uses a default `preference_` of 120. The value is determined by the class variable of the same name.

Each agent uses the class variable `INFINITY` (set at 32) to determine the validity of a route.

**Manual Routing** Manual routing is not a route computation protocol (like the others), but simply a way for users to configure the routing table by hand, much as you would with the “route” command on a workstation.

To use manual routing, enable it with `rtproto`, then set each nodes routing tables with the `add-route-to-adj-node` command. For example:

```
$ns rtproto Manual
set n1 [$ns node]
set n2 [$ns node]
$ns duplex-link $n1 $n2 10Mb 100ms DropTail
$n1 add-route-to-adj-node -default $n2
$n2 add-route-to-adj-node -default $n1
```

For a more complete example, see `tcl/ex/many_tcp.tcl`.

## 28.4 Internals and Architecture of Routing

We start with a discussion of the classes associated with unicast routing, and the code path used to configure and execute each of the different routing protocols. We conclude with a description of the interface between unicast routing and network dynamics, and that between unicast and multicast routing.

### 28.4.1 The classes

There are four main classes, the class `RouteLogic`, the class `rtObject`, the class `rtPeer`, and the base class `Agent/rtProto` for all protocols. In addition, the routing architecture extends the classes `Simulator`, `Link`, `Node` and `Classifier`.

**class RouteLogic** This class defines two methods to configure unicast routing, and one method to query it for route information. It also defines an instance procedure that is applicable when the topology is dynamic. We discuss this last procedure in conjunction with the interface to network dynamics.

- The instance procedure `register{}` is invoked by `Simulator::rtproto{}`. It takes the protocol and a list of nodes as arguments, and constructs an instance variable, `rtprotos_`, as an array; the array index is the name of the protocol, and the value is the list of nodes that will run this protocol.
- The `configure{}` reads the `rtprotos_` instance variable, and for each element in the array, invokes route protocol methods to perform the appropriate initializations. It is invoked by the simulator run procedure.

For each protocol `<rt-proto>` indexed in the `rtprotos_` array, this routine invokes `Agent/rtProto/<rt-proto> init-all rtprotos_(<rt-proto>)`.

If there are no elements in `rtprotos_`, the routine invokes Static routing, as `Agent/rtProto/Static init-all`.

- The instance procedure `lookup{}` takes two node numbers, `nodeId1` and `nodeId2`, as argument; it returns the id of the neighbor node that `nodeId1` uses to reach `nodeId2`.

The procedure is used by the static route computation procedure to query the computed routes and populate the routes at each of the nodes. It is also used by the multicast routing protocols to perform the appropriate RPF check.

Note that this procedure overloads an instproc-like of the same name. The procedure queries the appropriate `rtObject` entities if they exist (which they will if dynamic routing strategies are used in the simulation); otherwise, the procedure invokes the instproc-like to obtain the relevant information.

**class rtObject** is used in simulations that use dynamic routing. Each node has a `rtObject` associated with it, that acts as a co-ordinator for the different routing protocols that operate at a node. At any node, the `rtObject` at that node tracks each of the protocols operating at that node; it computes and installs the next route to each destination available via each of the protocols. In the event that the routing tables change, or the topology changes, the `rtObject` will alert the protocols to take the appropriate action.

The class defines the procedure `init-all{}`; this procedure takes a list of nodes as arguments, and creates a `rtObject` at each of the nodes in its argument list. It subsequently invokes its `compute-routes`.

The assumption is that the constructor for each of the new objects will instantiate the “Direct” route protocol at each of these nodes. This route protocol is responsible for computing the routes to immediately adjacent neighbors. When `compute-routes{}` is run by the `init-all{}` procedure, these direct routes are installed in the node by the appropriate route object.

The other instance procedures in this class are:

- `init{}` The constructor sets up pointers from itself to the node, in its instance variable `node_`, and from the node to itself, through the Node instance procedure `init-routing{}` and the Node instance variable `rtObject_`. It then initializes an array of `nextHop_`, `rtpref_`, `metric_`, `rtVia_`. The index of each of these arrays is the handle of the destination node.

The `nextHop_` contains the link that will be used to reach the particular destination; `rtpref_` and `metric_` are the preference and metric for the route installed in the node; `rtVia_` is the name of the agent whose route is installed in the node.

The constructor also creates the instance of the Direct route protocol, and invokes `compute-routes{}` for that protocol.

- `add-proto{}` creates an instance of the protocol, stores a reference to it in its array of protocols, `rtProtos_`. The index of the array is the name of the protocol. It also attaches the protocol object to the node, and returns the handle of the protocol object.

- `lookup{}` takes a destination node handle, and returns the id of the neighbor node that is used to reach the destination. If multiple paths are in use, then it returns a list of the neighbor nodes that will be used. If the node does not have a route to the destination, the procedure will return -1.

- `compute-routes{}` is the core procedure in this class. It first checks to see if any of the routing protocols at the node have computed any new routes. If they have, it will determine the best route to each destination from among all the protocols. If any routes have changed, the procedure will notify each of the protocols of the number of such changes, in case any of these protocols wants to send a fresh update. Finally, it will also notify any multicast protocol that new unicast route tables have been computed.

The routine checks the protocol agent's instance variable, `rtsChanged_` to see if any of the routes in that protocol have changed since the protocol was last examined. It then uses the protocol's instance variable arrays, `nextHop_`, `rtpref_`, and `metric_` to compute its own arrays. The `rtObject` will install or modify any of the routes as the changes are found.

If any of the routes at the node have changed, the `rtObject` will invoke the protocol agent's instance procedures, `send-updates{}` with the number of changes as argument. It will then invoke the multicast route object, if it exists.

The next set of routines are used to query the `rtObject` for various state information.

- `dump-routes{}` takes a output file descriptor as argument, and writes out the routing table at that node on that file descriptor.

A typical dump output is:

- `rtProto?{}` takes a route protocol as argument, and returns a handle to the instance of the protocol running at the node.
- `nextHop?{}` takes a destination node handle, and returns the link that is used to reach that destination.
- Similarly, `rtpref?{}` and `metric?{}` take a destination node handle as argument, and return the preference and metric of the route to the destination installed at the node.

**The class `rtPeer`** is a container class used by the protocol agents. Each object stores the address of the peer agent, and the metric and preference for each route advertised by that peer. A protocol agent will store one object per peer. The class maintains the instance variable `addr_`, and the instance variable arrays, `metric_` and `rtpref_`; the array indices are the destination node handles.

The class instance procedures, `metric{}` and `preference{}`, take one destination and value, and set the respective array variable. The procedures, `metric?{}` and `preference?{}`, take a destination and return the current value for that destination. The instance procedure `addr?{}` returns the address of the peer agent.

**class `Agent/rtProto`** This class is the base class from which all routing protocol agents are derived. Each protocol agent must define the procedure `init-all{}` to initialize the complete protocol, and possibly instance procedures `init{}`, `compute-routes{}`, and `send-updates{}`. In addition, if the topology is dynamic, and the protocol supports route computation to react to changes in the topology, then the protocol should define the procedure `compute-all{}`, and possibly the instance procedure `intf-changed{}`. In this section, we will briefly describe the interface for the basic procedures. We will defer the description of `compute-all{}` and `intf-changed{}` to the section on network dynamics. We also defer the description of the details of each of the protocols to their separate section at the end of the chapter.

- The procedure `init-all{}` is a global initialization procedure for the class. It may be given a list of the nodes as an argument. This is the list of nodes that should run this routing protocol. However, centralized routing protocols such as static and session routing will ignore this argument; detailed dynamic routing protocols such as DV will use this argument list to instantiate protocol agents at each of the nodes specified.

Note that derived classes in OTcl do not inherit the procedures defined in the base class. Therefore, every derived routing protocol class must define its own procedures explicitly.

- The instance procedure `init{}` is the constructor for protocol agents that are created. The base class constructor initializes the default preference for objects in this class, identifies the interfaces incident on the node and their current status. The interfaces are indexed by the neighbor handle and stored in the instance variable array, `ifs_`; the corresponding status instance variable array is `ifstat_`.

Centralized routing protocols such as static and session routing do not create separate agents per node, and therefore do not access any of these instance procedures.

- The instance procedure `compute-routes{}` computes the actual routes for the protocol. The computation is based on the routes learned by the protocol, and varies from protocol to protocol.

This routine is invoked by the `rtObject` whenever the topology changes. It is also invoked when the node receives an update for the protocol.

If the routine computes new routes, `rtObject::compute-routes{}` needs to be invoked to recompute and possibly install new routes at the node. The actual invoking of the `rtObject` is done by the procedure that invoked this routine in the first place.

- The instance procedure `send-updates{}` is invoked by the `rtObject` whenever the node routing tables have changed, and fresh updates have to be sent to all peers. The `rtObject` passes as argument the number of changes that were done. This procedure may also be invoked when there are no changes to the routes, but the topology incident on the node changes state. The number of changes is used to determine the list of peers to which a route update must be sent.

Other procedures relate to responding to topology changes and are described later (Section 28.4.2).

## Other Extensions to the Simulator, Node, Link, and Classifier

- We have discussed the methods `rtproto{}` and `cost{}` in the class `Simulator` earlier (Section 28.1). The one other method used internally is `get-routelogic{}`; this procedure returns the instance of `routelogic` in the simulation.

The method is used by the class `Simulator`, and unicast and multicast routing.

- The class `Node` contains these additional instance procedures to support dynamic unicast routing: `init-routing{}`, `add-routes{}`, `delete-routes{}`, and `rtObject?{}`.

The instance procedure `init-routing{}` is invoked by the `rtObject` at the node. It stores a pointer to the `rtObject`, in its instance variable `rtObject_`, for later manipulation or retrieval. It also checks its class variable to see if it should use `multiPath` routing, and sets up an instance variable to that effect. If `multiPath` routing could be used, the instance variable array `routes_` stores a count of the number of paths installed for each destination. This is the only array in unicast routing that is indexed by the node id, rather than the node handle.

The instance procedure `rtObject?{}` returns the `rtObject` handle for that node.

The instance procedure `add-routes{}` takes a node id, and a list of links. It will add the list of links as the routes to reach the destination identified by the node id. The realization of `multiPath` routing is done by using a separate `Classifier/multiPath`. For any given destination id *d*, if this node has multiple paths to *d*, then the main classifier points to this `multiPath` classifier instead of the link to reach the destination. Each of the multiple paths identified by the interfaces being used is installed in the `multiPath` classifier. The `multiPath` classifier will use each of the links installed in it for succeeding packets forwarded to it.

The instance procedure `delete-routes{}` takes a node id, a list of interfaces, and a `nullAgent`. It removes each of the interfaces in the list from the installed list of interfaces. If the entry did not previously use a multipath classifier, then it must have had only one route, and the route entry is set to point to the `nullAgent` specified.

**Q: WHY DOES IT NOT POINT TO NULLAGENT IF THE ENTRIES IN THE MPATHCLASSIFIER GOES TO ZERO?**

- The main extension to the class `Link` for unicast routing is to support the notion of link costs. The instance variable `cost_` contains the cost of the unidirectional link. The instance procedures `cost{}` and `cost?{}` set and get the cost on the link.

Note that `cost{}` takes the cost as argument. It is preferable to use the simulator method to set the cost variable, similar to the simulator instance procedures to set the queue or delay on a link.

- The class `Classifier` contains three new procedures, two of which overload an existing instproc-like, and the other two provide new functionality.

The instance procedure `install{}` overloads the existing instproc-like of the same name. The procedure stores the entry being installed in the instance variable `array`, `elements_`, and then invokes the instproc-like.

The instance procedure `installNext{}` also overloads the existing instproc-like of the same name. This instproc-like simply installs the entry into the next available slot.

The instance procedure `adjacents{}` returns a list of `(key, value)` pairs of all elements installed in the classifier.

## 28.4.2 Interface to Network Dynamics and Multicast

This section describes the methods applied in unicast routing to respond to changes in the topology. The complete sequence of actions that cause the changes in the topology, and fire the appropriate actions is described in a different section. The response to topology changes falls into two categories: actions taken by individual agents at each of the nodes, and actions to be taken globally for the entire protocol.

Detailed routing protocols such as the DV implementation require actions to be performed by individual protocol agents at the affected nodes. Centralized routing protocols such as static and session routing fall into the latter category exclusively. Detailed routing protocols could use such techniques to gather statistics related to the operation of the routing protocol; however, no such code is currently implemented in *ns*.

**Actions at the individual nodes** Following any change in the topology, the network dynamics models will first invoke `rtObject::intf-changed{}` at each of the affected nodes. For each of the unicast routing protocols operating at that node, `rtObject::intf-changed{}` will invoke each individual protocol's instance procedure, `intf-changed{}`, followed by that protocol's `compute-routes{}`.

After each protocol has computed its individual routes `rtObject::intf-changed{}` invokes `compute-routes{}` to possibly install new routes. If new routes were installed in the node, `rtObject::compute-routes{}` will invoke `send-updates{}` for each of the protocols operating at the node. The procedure will also flag the multicast route object of the route changes at the node, indicating the number of changes that have been executed. `rtObject::flag-multicast{}` will, in turn, notify the multicast route object to take appropriate action.

The one exception to the interface between unicast and multicast routing is the interaction between dynamic dense mode multicast and detailed unicast routing. This dynamicDM implementation in *ns* assumes neighbor nodes will send an implicit update whenever their routes change, without actually sending the update. It then uses this implicit information to compute appropriate parent-child relationships for the multicast spanning trees. Therefore, detailed unicast routing will invoke `rtObject_ flag-multicast 1` whenever it receives a route update as well, even if that update does not result in any change in its own routing tables.

**Global Actions** Once the detailed actions at each of the affected nodes is completed, the network dynamics models will notify the RouteLogic instance (`RouteLogic::notify{}`) of changes to topology. This procedure invokes the procedure `compute-all{}` for each of the protocols that were ever installed at any of the nodes. Centralized routing protocols such as session routing use this signal to recompute the routes to the topology. Finally, the `RouteLogic::notify{}` procedure notifies any instances of centralized multicast that are operating at the node.

## 28.5 Protocol Internals

In this section, we describe any leftover details of each of the routing protocol agents. Note that this is the only place where we describe the internal route protocol agent, “Direct” routing.

**Direct Routing** This protocol tracks the state of the incident links, and maintains routes to immediately adjacent neighbors only. As with the other protocols, it maintains instance variable arrays of `nextHop_`, `rtpref_`, and `metric_`, indexed by the handle of each of the possible destinations in the topology.

The instance procedure `compute-routes{}` computes routes based on the current state of the link, and the previously known state of the incident links.

No other procedures or instance procedures are defined for this protocol.

**Static Routing** The procedure `compute-routes{}` in the class `RouteLogic` first creates the adjacency matrix, and then invokes the C++ method, `compute_routes()` of the shadow object. Finally, the procedure retrieves the result of the route computation, and inserts the appropriate routes at each of the nodes in the topology.

The class only defines the procedure `init-all{}` that invokes `compute-routes{}`.

**Session Routing** The class defines the procedure `init-all{}` to compute the routes at the start of the simulation. It also defines the procedure `compute-all{}` to compute the routes when the topology changes. Each of these procedures directly invokes `compute-routes{}`.

**DV Routing** In a dynamic routing strategy, nodes send and receive messages, and compute the routes in the topology based on the messages exchanged. The procedure `init-all{}` takes a list of nodes as the argument; the default is the list of nodes in the topology. At each of the nodes in the argument, the procedure starts the class `rtObject` and a class `Agent/rtProto/DV` agents. It then determines the DV peers for each of the newly created DV agents, and creates the relevant `rtPeer` objects.

The constructor for the DV agent initializes a number of instance variables; each agent stores an array, indexed by the destination node handle, of the preference and metric, the interface (or link) to the next hop, and the remote peer incident on the interface, for the best route to each destination computed by the agent. The agent creates these instance variables, and then schedules sending its first update within the first 0.5 seconds of simulation start.

Each agent stores the list of its peers indexed by the handle of the peer node. Each peer is a separate peer structure that holds the address of the peer agent, the metric and preference of the route to each destination advertised by that peer. We discuss the `rtPeer` structure later when discuss the route architecture. The peer structures are initialized by the procedure `add-peer{}` invoked by `init-all{}`.



The routine `send-periodic-update{}` invokes `send-updates{}` to send the actual updates. It then reschedules sending the next periodic update after `advertInterval` jittered slightly to avoid possible synchronization effects.

`send-updates{}` will send updates to a select set of peers. If any of the routes at that node have changed, or for periodic updates, the procedure will send updates to all peers. Otherwise, if some incident links have just recovered, the procedure will send updates to the adjacent peers on those incident links only.

`send-updates{}` uses the procedure `send-to-peer{}` to send the actual updates. This procedure packages the update, taking the split-horizon and poison reverse mechanisms into account. It invokes the instproc-like, `send-update{}` (Note the singular case) to send the actual update. The actual route update is stored in the class variable `msg_` indexed by a non-decreasing integer as index. The instproc-like only sends the index to `msg_` to the remote peer. This eliminates the need to convert from OTcl strings to alternate formats and back.

When a peer receives a route update it first checks to determine if the update from differs from the previous ones. The agent will compute new routes if the update contains new information.

## 28.6 Unicast routing objects

RouteLogic and rtObject are two objects that are significant to unicast routing in *ns*. RouteLogic, essentially, represents the routing table that is created and maintained centrally for every unicast simulation. rtObject is the object that every node taking part in dynamic unicast routing, has an instance of. Note that nodes will not have an instance of this object if Session routing is done as a detailed routing protocol is not being simulated in this case. The methods for RouteLogic and rtObject are described in the next section.

## 28.7 Commands at a glance

Following is a list of unicast routing related commands used in simulation scripts:

```
$ns_ rtproto <routing-proto> <args>
```

where `<routing-proto>` defines the type of routing protocol to be used, like Static, Manual, Session, DV etc. `args` may define the list of nodes on which the protocol is to be run. The node list defaults to all nodes in the topology.

Internal methods:

```
$ns_ compute-routes
```

This command computes `next_hop` information for all nodes in the topology using the topology connectivity. This `next_hop` info is then used to populate the node classifiers or the routing tables. `compute-routes` calls `compute-flat-routes` or `compute-hier-routes` depending on the type of addressing being used for the simulation.

```
$ns_ get-routelogic
```

This returns a handle to the RouteLogic object (the routing table), if one has been created. Otherwise a new routing table object is created.

```
$ns_ dump-routelogic-nh
```

This dumps next hop information in the routing table.

```
$ns_ dump-routelogic-distance
```

This dumps the distance information in the routing table.

```
$node add-route <dst> <Target>
```

This is a method used to add routing entries (nexthop information) in the node's routing table. The nexthop to <dst> from this node is the <target> object and this info is added to the node's classifier.

```
$routelogic lookup <srcid> <destid>
```

Returns the id of the node that is the next hop from the node with id srcid to the node with id destid.

```
$routelogic dump <nodeid>
```

Dump the routing tables of all nodes whose id is less than nodeid. Node ids are typically assigned to nodes in ascending fashion starting from 0 by their order of creation.

```
rtobject dump-routes <fileID>
```

Dump the routing table to the output channel specified by fileID. fileID must be a file handle returned by the Tcl open command and it must have been opened for writing.

```
$rtobject rtProto? <proto>
```

Returns a handle to the routing protocol agent specified by proto if it exists at that node. Returns an empty string otherwise.

```
$rtobject nextHop? <destID>
```

Returns the id of the node that is the next hop to the destination specified by the node id, <destID>.

```
$rtobject rtpref? destID
```

Returns the preference for the route to destination node given by destid.

```
$rtobject metric? destID
```

Returns metric for the route to destid.

# Chapter 29

## Multicast Routing

This section describes the usage and the internals of multicast routing implementation in *ns*. We first describe the user interface to enable multicast routing (Section 29.1), specify the multicast routing protocol to be used and the various methods and configuration parameters specific to the protocols currently supported in *ns*. We then describe in detail the internals and the architecture of the multicast routing implementation in *ns* (Section 29.2).

The procedures and functions described in this chapter can be found in various files in the directories *~ns/tcl/mcast*, *~ns/tcl/ctr-mcast*; additional support routines are in *~ns/mcast\_ctrl.cc,h*, *~ns/tcl/lib/ns-lib.tcl*, and *~ns/tcl/lib/ns-node.tcl*.

### 29.1 Multicast API

Multicast forwarding requires enhancements to the nodes and links in the topology. Therefore, the user must specify multicast requirements to the Simulator class before creating the topology. This is done in one of two ways:

```
set ns [new Simulator -multicast on]
or
set ns [new Simulator]
$ns multicast
```

When multicast extensions are thus enabled, nodes will be created with additional classifiers and replicators for multicast forwarding, and links will contain elements to assign incoming interface labels to all packets entering a node.

A multicast routing strategy is the mechanism by which the multicast distribution tree is computed in the simulation. *ns* supports three multiast route computation strategies: centralised, dense mode(DM) or shared tree mode(ST).

The method `mrtproto{}` in the Class Simulator specifies either the route computation strategy, for centralised multicast routing, or the specific detailed multicast routing protocol that should be used.

The following are examples of valid invocations of multicast routing in *ns*:

```
set cmc [$ns mrtproto CtrMcast]           ; # specify centralized multicast for all nodes
                                           ; # cmc is the handle for multicast protocol object
$ns mrtproto DM                           ; # specify dense mode multicast for all nodes
```

```
$ns mrtproto ST
```

```
; # specify shared tree mode to run on all nodes
```

Notice in the above examples that CtrMcast returns a handle that can be used for additional configuration of centralised multicast routing. The other routing protocols will return a null string. All the nodes in the topology will run instances of the same protocol.

Multiple multicast routing protocols can be run at a node, but in this case the user must specify which protocol owns which incoming interface. For this finer control `mrtproto-iifs{}` is used.

New/unused multicast address are allocated using the procedure `allocaddr{}`.

The agents use the instance procedures `join-group{}` and `leave-group{}`, in the class `Node` to join and leave multicast groups. These procedures take two mandatory arguments. The first argument identifies the corresponding agent and second argument specifies the group address.

An example of a relatively simple multicast configuration is:

```
set ns [new Simulator -multicast on] ; # enable multicast routing
set group [Node allocaddr] ; # allocate a multicast address
set node0 [$ns node] ; # create multicast capable nodes
set node1 [$ns node]
$ns duplex-link $node0 $node1 1.5Mb 10ms DropTail

set mproto DM ; # configure multicast protocol
set mrthandle [$ns mrtproto $mproto] ; # all nodes will contain multicast protocol agents
set udp [new Agent/UDP] ; # create a source agent at node 0
$ns attach-agent $node0 $udp
set src [new Application/Traffic/CBR]
$src attach-agent $udp
$udp set dst_addr_ $group
$udp set dst_port_ 0

set rcvr [new Agent/LossMonitor] ; # create a receiver agent at node 1
$ns attach-agent $node1 $rcvr
$ns at 0.3 "$node1 join-group $rcvr $group" ; # join the group at simulation time 0.3 (sec)
```

### 29.1.1 Multicast Behavior Monitor Configuration

*ns* supports a multicast monitor module that can trace user-defined packet activity. The module counts the number of packets in transit periodically and prints the results to specified files. `attach{}` enables a monitor module to print output to a file. `trace-topo{}` insets monitor modules into all links. `filter{}` allows accounting on specified packet header, field in the header, and value for the field). Calling `filter{}` repeatedly will result in an AND effect on the filtering condition. `print-trace{}` notifies the monitor module to begin dumping data. `pctype( )` is a global array that takes a packet type name (as seen in `trace-all{}` output) and maps it into the corresponding value. A simple configuration to filter cbr packets on a particular group is:

```
set mcastmonitor [new McastMonitor]
set chan [open cbr.tr w] ; # open trace file
$mmonitor attach $chan1 ; # attach trace file to McastMonitor object
$mcastmonitor set period_ 0.02 ; # default 0.03 (sec)
```

```

$monitor trace-topo ; # trace entire topology
$monitor filter Common ptype_ $ptype(cbr) ; # filter on ptype_ in Common header
$monitor filter IP dst_ $group ; # AND filter on dst_ address in IP header
$monitor print-trace ; # begin dumping periodic traces to specified files

```

The following sample output illustrates the output file format (time, count):

```

0.16 0
0.17999999999999999 0
0.19999999999999998 0
0.21999999999999997 6
0.23999999999999996 11
0.25999999999999995 12
0.27999999999999997 12

```

### 29.1.2 Protocol Specific configuration

In this section, we briefly illustrate the protocol specific configuration mechanisms for all the protocols implemented in *ns*.

**Centralized Multicast** The centralized multicast is a sparse mode implementation of multicast similar to PIM-SM [9]. A Rendezvous Point (RP) rooted shared tree is built for a multicast group. The actual sending of prune, join messages etc. to set up state at the nodes is not simulated. A centralized computation agent is used to compute the forwarding trees and set up multicast forwarding state, (S, G) at the relevant nodes as new receivers join a group. Data packets from the senders to a group are unicast to the RP. Note that data packets from the senders are unicast to the RP even if there are no receivers for the group.

The method of enabling centralised multicast routing in a simulation is:

```

set mproto CtrMcast ; # set multicast protocol
set mrthandle [$ns mrtproto $mproto]

```

The command procedure `mrtproto{}` returns a handle to the multicast protocol object. This handle can be used to control the RP and the boot-strap-router (BSR), switch tree-types for a particular group, from shared trees to source specific trees, and recompute multicast routes.

```

$mrthandle set_c_rp $node0 $node1 ; # set the RPs
$mrthandle set_c_bsr $node0:0 $node1:1 ; # set the BSR, specified as list of node:priority

$mrthandle get_c_rp $node0 $group ; # get the current RP ???
$mrthandle get_c_bsr $node0 ; # get the current BSR

$mrthandle switch-treetype $group ; # to source specific or shared tree

$mrthandle compute-mroutes ; # recompute routes. usually invoked automatically as needed

```

Note that whenever network dynamics occur or unicast routing changes, `compute-mroutes{}` could be invoked to recompute the multicast routes. The instantaneous re-computation feature of centralised algorithms may result in causality violations during the transient periods.

**Dense Mode** The Dense Mode protocol (`DM.tcl`) is an implementation of a dense-mode-like protocol. Depending on the value of DM class variable `CacheMissMode` it can run in one of two modes. If `CacheMissMode` is set to `pimdm` (default), PIM-DM-like forwarding rules will be used. Alternatively, `CacheMissMode` can be set to `dvmrp` (loosely based on DVMRP [30]). The main difference between these two modes is that DVMRP maintains parent-child relationships among nodes to reduce the number of links over which data packets are broadcast. The implementation works on point-to-point links as well as LANs and adapts to the network dynamics (links going up and down).

Any node that receives data for a particular group for which it has no downstream receivers, send a prune upstream. A prune message causes the upstream node to initiate prune state at that node. The prune state prevents that node from sending data for that group downstream to the node that sent the original prune message while the state is active. The time duration for which a prune state is active is configured through the DM class variable, `PruneTimeout`. A typical DM configuration is shown below:

```
DM set PruneTimeout 0.3 ; # default 0.5 (sec)
DM set CacheMissMode dvmrp ; # default pimdm
$ns mrtproto DM
```

**Shared Tree Mode** Simplified sparse mode `ST.tcl` is a version of a shared-tree multicast protocol. Class variable array `RP_` indexed by group determines which node is the RP for a particular group. For example:

```
ST set RP_($group) $node0
$ns mrtproto ST
```

At the time the multicast simulation is started, the protocol will create and install encapsulator objects at nodes that have multicast senders, decapsulator objects at RPs and connect them. To join a group, a node sends a graft message towards the RP of the group. To leave a group, it sends a prune message. The protocol currently does not support dynamic changes and LANs.

**Bi-directional Shared Tree Mode** `BST.tcl` is an experimental version of a bi-directional shared tree protocol. As in shared tree mode, RPs must be configured manually by using the class array `RP_`. The protocol currently does not support dynamic changes and LANs.

## 29.2 Internals of Multicast Routing

We describe the internals in three parts: first the classes to implement and support multicast routing; second, the specific protocol implementation details; and finally, provide a list of the variables that are used in the implementations.

### 29.2.1 The classes

The main classes in the implementation are the class `mrtObject` and the class `McastProtocol`. There are also extensions to the base classes: `Simulator`, `Node`, `Classifier`, *etc.* We describe these classes and extensions in this subsection. The specific protocol implementations also use adjunct data structures for specific tasks, such as timer mechanisms by detailed dense mode, encapsulation/decapsulation agents for centralised multicast *etc.*; we defer the description of these objects to the section on the description of the particular protocol itself.

**mrtObject class** There is one mrtObject (aka Arbiter) object per multicast capable node. This object supports the ability for a node to run multiple multicast routing protocols by maintaining an array of multicast protocols indexed by the incoming interface. Thus, if there are several multicast protocols at a node, each interface is owned by just one protocol. Therefore, this object supports the ability for a node to run multiple multicast routing protocols. The node uses the arbiter to perform protocol actions, either to a specific protocol instance active at that node, or to all protocol instances at that node.

<code>addproto{instance, [iiflist]}</code>	adds the handle for a protocol instance to its array of protocols. The second optional argument is the incoming interface list controlled by the protocol. If this argument is an empty list or not specified, the protocol is assumed to run on all interfaces (just one protocol).
<code>getType{protocol}</code>	returns the handle to the protocol instance active at that node that matches the specified type (first and only argument). This function is often used to locate a protocol's peer at another node. An empty string is returned if the protocol of the given type could not be found.
<code>all-mprotos{op, args}</code>	internal routine to execute "op" with "args" on all protocol instances active at that node.
<code>start{}</code>	
<code>stop{}</code>	start/stop execution of all protocols.
<code>notify{dummy}</code>	is called when a topology change occurs. The dummy argument is currently not used.
<code>dump-mroutes{file-handle, [grp], [src]}</code>	dump multicast routes to specified file-handle.
<code>join-group{G, S}</code>	signals all protocol instances to join (S, G).
<code>leave-group{G, S}</code>	signals all protocol instances to leave (S, G).
<code>upcall{code, s, g, iif}</code>	signalled by node on forwarding errors in classifier; this routine in turn signals the protocol instance that owns the incoming interface (iif) by invoking the appropriate handle function for that particular code.
<code>drop{rep, s, g, iif}</code>	Called on packet drop, possibly to prune an interface.

In addition, the mrtObject class supports the concept of well known groups, *i.e.*, those groups that do not require explicit protocol support. Two well known groups, ALL\_ROUTERS and ALL\_PIM\_ROUTERS are predefined in *ns*.

The class mrtObject defines two class procedures to set and get information about these well known groups.

<code>registerWellKnownGroups{name}</code>	assigns name a well known group address.
<code>getWellKnownGroup{name}</code>	returns the address allocated to well known group, name. If name is not registered as a well known group, then it returns the address for ALL_ROUTERS.

**McastProtocol class** This is the base class for the implementation of all the multicast protocols. It contains basic multicast functions:

<code>start{ }, stop{ }</code>	Set the status_ of execution of this protocol instance.
<code>getStatus{ }</code>	return the status of execution of this protocol instance.
<code>getType{ }</code>	return the type of protocol executed by this instance.
<code>upcall{code args}</code>	invoked when the node classifier signals an error, either due to a cache-miss or a wrong-iif for incoming packet. This routine invokes the protocol specific handle, <code>handle-&lt;code&gt;{ }</code> with specified args to handle the signal.

A few words about interfaces. Multicast implementation in *ns* assumes duplex links *i.e.* if there is a simplex link from node 1 to node 2, there must be a reverse simplex link from node 2 to node 1. To be able to tell from which link a packet was received, multicast simulator configures links with an interface labeller at the end of each link, which labels packets with a particular

and unique label (id). Thus, “incoming interface” is referred to this label and is a number greater or equal to zero. Incoming interface value can be negative (-1) for a special case when the packet was sent by a local to the given node agent.

In contrast, an “outgoing interface” refers to an object handler, usually a head of a link which can be installed at a replicator. This distinction is important: *incoming interface is a numeric label to a packet, while outgoing interface is a handler to an object that is able to receive packets, e.g. head of a link.*

### 29.2.2 Extensions to other classes in *ns*

In the earlier chapter describing nodes in *ns* (Chapter 5), we described the internal structure of the node in *ns*. To briefly recap that description, the node entry for a multicast node is the `switch_`. It looks at the highest bit to decide if the destination is a multicast or unicast packet. Multicast packets are forwarded to a multicast classifier which maintains a list of replicators; there is one replicator per  $\langle \text{source}, \text{group} \rangle$  tuple. Replicators copy the incoming packet and forward to all outgoing interfaces.

**Class Node** Node support for multicast is realized in two primary ways: it serves as a focal point for access to the multicast protocols, in the areas of address allocation, control and management, and group membership dynamics; and secondly, it provides primitives to access and control interfaces on links incident on that node.



<code>expandaddr{}</code> , <code>allocaddr{}</code>	Class procedures for address management. <code>expandaddr{}</code> is now obsolete. <code>allocaddr{}</code> allocates the next available multicast address.
<code>start-mcast{}</code> , <code>stop-mcast{}</code>	To start and stop multicast routing at that node.
<code>notify-mcast{}</code>	<code>notify-mcast{}</code> signals the <code>mrtObject</code> at that node to recompute multicastroutes following a topology change or unicast route update from a neighbour.
<code>getArbiter{}</code>	returns a handle to <code>mrtObject</code> operating at that node.
<code>dump-routes{file-handle}</code>	to dump the multicast forwarding tables at that node.
<code>new-group{sg iif code}</code>	When a multicast data packet is received, and the multicast classifier cannot find the slot corresponding to that data packet, it invokes <code>Node instproc new-group{}</code> to establish the appropriate entry. The code indicates the reason for not finding the slot. Currently there are two possibilities, cache-miss and wrong-iif. This procedure notifies the arbiter instance to establish the new group.
<code>join-group{ag}</code>	An agent at a node that joins a particular group invokes “ <code>node join-group &lt;agent&gt; &lt;group&gt;</code> ”. The node signals the <code>mrtObject</code> to join the particular group, and adds <code>agent</code> to its list of agents at that group. It then adds <code>agent</code> to all replicators associated with <code>group</code> .
<code>leave-group{ag}</code>	<code>Node instproc leave-group</code> reverses the process described earlier. It disables the outgoing interfaces to the receiver agents for all the replicators of the group, deletes the receiver agents from the local <code>Agents_</code> list; it then invokes the arbiter instance’s <code>leave-group{}</code> .
<code>add-mfc{sg iif oiflist}</code>	<code>Node instproc add-mfc</code> adds a <i>multicast forwarding cache</i> entry for a particular $\langle \text{source}, \text{group}, \text{iif} \rangle$ . The mechanism is: <ul style="list-style-type: none"> <li>• create a new replicator (if one does not already exist),</li> <li>• update the <code>replicator_</code> instance variable array at the node,</li> <li>• add all outgoing interfaces and local agents to the appropriate replicator,</li> <li>• invoke the multicast classifier’s <code>add-rep{}</code> to create a slot for the replicator in the multicast classifier.</li> </ul>
<code>del-mfc{sg oiflist}</code>	disables each oif in <code>oiflist</code> from the replicator for $\langle s, g \rangle$ .

The list of primitives accessible at the node to control its interfaces are listed below.

<code>add-iif{ifid link},</code>	
<code>add-oif{link if}</code>	Invoked during link creation to prep the node about its incoming interface label and outgoing interface object.
<code>get-all-oifs{}</code>	Returns all oifs for this node.
<code>get-all-iifs{}</code>	Returns all iifs for this node.
<code>iif2link{ifid}</code>	Returns the link object labelled with given interface label.
<code>link2iif{link}</code>	Returns the incoming interface label for the given link.
<code>oif2link{oif}</code>	Returns the link object corresponding to the given outgoing interface.
<code>link2oif{link}</code>	Returns the outgoing interface for the link ( <i>ns</i> object that is incident to the node).
<code>rpf-nbr{src}</code>	Returns a handle to the neighbour node that is its next hop to the specified <code>src</code> .
<code>getReps{s g}</code>	Returns a handle to the replicator that matches <code>&lt;s, g&gt;</code> . Either argument can be a wildcard (*).
<code>getReps-raw{s g}</code>	As above, but returns a list of <code>&lt;key, handle&gt;</code> pairs.
<code>clearReps{s g}</code>	Removes all replicators associated with <code>&lt;s, g&gt;</code> .

**Class Link and SimpleLink** This class provides methods to check the type of link, and the label it affixes on individual packets that traverse it. There is one additional method to actually place the interface objects on this link. These methods are:

`if-label?{}` returns the interface label affixed by this link to packets that traverse it.

**Class NetworkInterface** This is a simple connector that is placed on each link. It affixes its label id to each packet that traverses it. The packet id is used by the destination node incident on that link to identify the link by which the packet reached it. The label id is configured by the Link constructor. This object is an internal object, and is not designed to be manipulated by user level simulation scripts. The object only supports two methods:

`label{ifid}` assigns `ifid` that this object will affix to each packet.  
`label{}` returns the label that this object affixes to each packet.

The global class variable, `ifacenum_`, specifies the next available `ifid` number.

**Class Multicast Classifier** Classifier/Multicast maintains a *multicast forwarding cache*. There is one multicast classifier per node. The node stores a reference to this classifier in its instance variable `multiclassifier_`. When this classifier receives a packet, it looks at the `<source, group>` information in the packet headers, and the interface that the packet arrived from (the incoming interface or `iif`); does a lookup in the MFC and identifies the slot that should be used to forward that packet. The slot will point to the appropriate replicator.

However, if the classifier does not have an entry for this `<source, group>`, or the `iif` for this entry is different, it will invoke an upcall `new-group{}` for the classifier, with one of two codes to identify the problem:

- `cache-miss` indicates that the classifier did not find any `<source, group>` entries;
- `wrong-iif` indicates that the classifier found `<source, group>` entries, but none matching the interface that this packet arrived on.

These upcalls to TCL give it a chance to correct the situation: install an appropriate MFC-entry (for `cache-miss`) or change the incoming interface for the existing MFC-entry (for `wrong-iif`). The *return value* of the upcall determines what classifier will do with the packet. If the return value is "1", it will assume that TCL upcall has appropriately modified MFC

will try to classify packet (lookup MFC) for the second time. If the return value is “0”, no further lookups will be done, and the packet will be thus dropped.

`add-rep{ }` creates a slot in the classifier and adds a replicator for  $\langle \text{source}, \text{group}, \text{iif} \rangle$  to that slot.

**Class Replicator** When a replicator receives a packet, it copies the packet to all of its slots. Each slot points to an outgoing interface for a particular  $\langle \text{source}, \text{group} \rangle$ .

If no slot is found, the C++ replicator invokes the class instance procedure `drop{ }` to trigger protocol specific actions. We will describe the protocol specific actions in the next section, when we describe the internal procedures of each of the multicast routing protocols.

There are instance procedures to control the elements in each slot:

<code>insert{oif}</code>	inserting a new outgoing interface to the next available slot.
<code>disable{oif}</code>	disable the slot pointing to the specified oif.
<code>enable{oif}</code>	enable the slot pointing to the specified oif.
<code>is-active{ }</code>	returns true if the replicator has at least one active slot.
<code>exists{oif}</code>	returns true if the slot pointing to the specified oif is active.
<code>change-iface{source, group, oldiif, newiif}</code>	modified the iif entry for the particular replicator.

### 29.2.3 Protocol Internals

We now describe the implementation of the different multicast routing protocol agents.

#### Centralized Multicast

`CtrMcast` inherits from `McastProtocol`. One `CtrMcast` agent needs to be created for each node. There is a central `CtrMcastComp` agent to compute and install multicast routes for the entire topology. Each `CtrMcast` agent processes membership dynamic commands, and redirects the `CtrMcastComp` agent to recompute the appropriate routes.

<code>join-group{ }</code>	registers the new member with the <code>CtrMcastComp</code> agent, and invokes that agent to recompute routes for that member.
<code>leave-group{ }</code>	is the inverse of <code>join-group{ }</code> .
<code>handle-cache-miss{ }</code>	called when no proper forwarding entry is found for a particular packet source. In case of centralized multicast, it means a new source has started sending data packets. Thus, the <code>CtrMcast</code> agent registers this new source with the <code>CtrMcastComp</code> agent. If there are any members in that group, compute the new multicast tree. If the group is in RPT (shared tree) mode, then

1. create an encapsulation agent at the source;
2. a corresponding decapsulation agent is created at the RP;
3. the two agents are connected by unicast; and
4. the  $\langle S, G \rangle$  entry points its outgoing interface to the encapsulation agent.

CtrMcastComp is the centralised multicast route computation agent.

<code>reset-mroutes{ }</code>	resets all multicast forwarding entries.
<code>compute-mroutes{ }</code>	(re)computes all multicast forwarding entries.
<code>compute-tree{source, group}</code>	computes a multicast tree for one source to reach all the receivers in a specific group.
<code>compute-branch{source, group, member}</code>	is executed when a receiver joins a multicast group. It could also be invoked by <code>compute-tree{ }</code> when it itself is recomputing the multicast tree, and has to re-parent all receivers. The algorithm starts at the receiver, recursively finding successive next hops, until it either reaches the source or RP, or it reaches a node that is already a part of the relevant multicast tree. During the process, several new replicators and an outgoing interface will be installed.
<code>prune-branch{source, group, member}</code>	is similar to <code>compute-branch{ }</code> except the outgoing interface is disabled; if the outgoing interface list is empty at that node, it will walk up the multicast tree, pruning at each of the intermediate nodes, until it reaches a node that has a non-empty outgoing interface list for the particular multicast tree.

## Dense Mode

<code>join-group{group}</code>	sends graft messages upstream if $\langle S, G \rangle$ does not contain any active outgoing slots ( <i>i.e.</i> , no downstream receivers). If the next hop towards the source is a LAN, increments a counter of receivers for a particular group for the LAN
<code>leave-group{group}</code>	decrements LAN counters.
<code>handle-cache-miss{srcID group iface}</code>	depending on the value of <code>CacheMissMode</code> calls either <code>handle-cache-miss-pimdm</code> or <code>handle-cache-miss-dvmrp</code> .
<code>handle-cache-miss-pimdm{srcID group iface}</code>	if the packet was received on the correct iif (from the node that is the next hop towards the source), fan out the packet on all oifs except the oif that leads back to the next-hop-neighbor and oifs that are LANs for which this node is not a forwarder. Otherwise, if the interface was incorrect, send a prune back.
<code>handle-cache-miss-dvmrp{srcID group iface}</code>	fans out the packet only to nodes for which this node is a next hop towards the source (parent).
<code>drop{replicator source group iface}</code>	sends a prune message back to the previous hop.
<code>recv-prune{from source group iface}</code>	resets the prune timer if the interface had been pruned previously; otherwise, it starts the prune timer and disables the interface; furthermore, if the outgoing interface list becomes empty, it propagates the prune message upstream.
<code>recv-graft{from source group iface}</code>	cancels any existing prune timer, and re-enables the pruned interface. If the outgoing interface list was previously empty, it forwards the graft upstream.
<code>handle-wrong-iif{srcID group iface}</code>	This is invoked when the multicast classifier drops a packet because it arrived on the wrong interface, and invokes <code>new-group{ }</code> . This routine is invoked by <code>mrtObject instproc new-group{ }</code> . When invoked, it sends a prune message back to the source.

## 29.2.4 The internal variables

### Class mrtObject

<code>protocols_</code>	An array of handles of protocol instances active at the node at which this protocol operates indexed by incoming interface.
<code>mask-wkgroups</code>	Class variable—defines the mask used to identify well known groups.
<code>wkgroups</code>	Class array variable—array of allocated well known groups addresses, indexed by the group name. <code>wkgroups(Allocd)</code> is a special variable indicating the highest currently allocated well known group.

### McastProtocol

<code>status_</code>	takes values “up” or “down”, to indicate the status of execution of the protocol instance.
<code>type_</code>	contains the type (class name) of protocol executed by this instance, <i>e.g.</i> , DM, or ST.

### Simulator

<code>multiSim_</code>	1 if multicast simulation is enabled, 0 otherwise.
<code>MrtHandle_</code>	handle to the centralised multicast simulation object.

### Node

<code>switch_</code>	handle for classifier that looks at the high bit of the destination address in each packet to determine whether it is a multicast packet (bit = 1) or a unicast packet (bit = 0).
<code>multiclassifier_</code>	handle to classifier that performs the $\langle s, g, iif \rangle$ match.
<code>replicator_</code>	array indexed by $\langle s, g \rangle$ of handles that replicate a multicast packet on to the required links.
<code>Agents_</code>	array indexed by multicast group of the list of agents at the local node that listen to the specific group.
<code>outLink_</code>	Cached list of outgoing interfaces at this node.
<code>inLink_</code>	Cached list of incoming interfaces at this node.

### Link and SimpleLink

<code>iif_</code>	handle for the NetworkInterface object placed on this link.
<code>head_</code>	first object on the link, a no-op connector. However, this object contains the instance variable, <code>link_</code> , that points to the container Link object.

### NetworkInterface

<code>ifacenum_</code>	Class variable—holds the next available interface number.
------------------------	---

## 29.3 Commands at a glance

Following is a list of commands used for multicast simulations:

```
set ns [new Simulator -mcast on]
```

This turns the multicast flag on for the the given simulation, at the time of creation of the simulator object.

```
ns_ multicast
```

This like the command above turns the multicast flag on.

```
ns_ multicast?
```

This returns true if multicast flag has been turned on for the simulation and returns false if multicast is not turned on.

`$ns_ mrtproto <mproto> <optional:nodelist>`

This command specifies the type of multicast protocol <mproto> to be used like DM, CtrMcast etc. As an additional argument, the list of nodes <nodelist> that will run an instance of detailed routing protocol (other than centralised mcast) can also be passed.

`$ns_ mrtproto-iifs <mproto> <node> <iifs>`

This command allows a finer control than mrtproto. Since multiple mcast protocols can be run at a node, this command specifies which mcast protocol <mproto> to run at which of the incoming interfaces given by <iifs> in the <node>.

`Node allocaddr`

This returns a new/unused multicast address that may be used to assign a multicast address to a group.

`Node expandaddr`

**THIS COMMAND IS OBSOLETE NOW** This command expands the address space from 16 bits to 30 bits. However this command has been replaced by "`ns_ set-address-format-expanded`".

`$node_ join-group <agent> <grp>`

This command is used when the <agent> at the node joins a particular group <grp>.

`$node_ leave-group <agent> <grp>`

This is used when the <agent> at the nodes decides to leave the group <grp>.

**Internal methods:**

`$ns_ run-mcast`

This command starts multicast routing at all nodes.

`$ns_ clear-mcast`

This stops mcast routing at all nodes.

`$node_ enable-mcast <sim>`

This allows special mcast supporting mechanisms (like a mcast classifier) to be added to the mcast-enabled node. <sim> is the handle to the simulator object.

In addition to the internal methods listed here there are other methods specific to protocols like centralized mcast (CtrMcast), dense mode (DM), shared tree mode (ST) or bi-directional shared tree mode (BST), Node and Link class methods and NetworkInterface and Multicast classifier methods specific to multicast routing. All mcast related files may be found under `ns/tcl/mcast` directory.

**Centralised Multicast** A handle to the CtrMcastComp object is returned when the protocol is specified as 'CtrMcast' in mrtproto. Ctrmcast methods are:

`$ctrmcastcomp switch-treetype group-addr`

Switch from the Rendezvous Point rooted shared tree to source-specific trees for the group specified by group-addr.

Note that this method cannot be used to switch from source-specific trees back to a shared tree for a multicast group.

`$ctrmcastcomp set_c_rp <node-list>`

This sets the RPs.

`$ctrmcastcomp set_c_bsr <node0:0> <node1:1>`

This sets the BSR, specified as list of node:priority.

`$ctrmcastcomp get_c_rp <node> <group>`

Returns the RP for the group as seen by the node node for the multicast group with address group-addr. Note that different nodes may see different RPs for the group if the network is partitioned as the nodes might be in different partitions.

`$ctrmcastcomp get_c_bsr <node>`

Returns the current BSR for the group.

`$ctrmcastcomp compute-mroutes`

This recomputes multicast routes in the event of network dynamics or a change in unicast routes.

**Dense Mode** The dense mode (DM) protocol can be run as PIM-DM (default) or DVMRP depending on the class variable `CacheMissMode`. There are no methods specific to this mcast protocol object. Class variables are:

**PruneTimeout** Timeout value for prune state at nodes. defaults to 0.5sec.

**CacheMissMode** Used to set PIM-DM or DVMRP type forwarding rules.

**Shared Tree** There are no methods for this class. Variables are:

**RP\_** `RP_` indexed by group determines which node is the RP for a particular group.

**Bidirectional Shared Tree** There are no methods for this class. Variable is same as that of Shared Tree described above.

# Chapter 30

## Network Dynamics

This chapter describes the capabilities in *ns* to make the simulation topologies dynamic. We start with the instance procedures to the class Simulator that are useful to a simulation script (Section 30.1). The next section describes the internal architecture (Section 30.2), including the different classes and instance variables and procedures; the following section describes the interaction with unicast routing (Section 30.3). This aspect of network dynamics is still somewhat experimental in *ns*. The last section of this chapter outlines some of the deficiencies in the current realization (Section 30.4) of network dynamics, some one or which may be fixed in the future.

The procedures and functions described in this chapter can be found in `~ns/tcl/rnglib/dynamics.tcl` and `~ns/tcl/lib/route-proto.tcl`.

### 30.1 The user level API

The user level interface to network dynamics is a collection of instance procedures in the class Simulator, and one procedure to trace and log the dynamics activity. Reflecting a rather poor choice of names, these procedures are `rtmodel`, `rtmodel-delete`, and `rtmodel-at`. There is one other procedure, `rtmodel-configure`, that is used internally by the class Simulator to configure the `rtmodels` just prior to simulation start. We describe this method later (Section 30.2).

- The instance procedure `rtmodel{}` defines a model to be applied to the nodes and links in the topology. Some examples of this command as it would be used in a simulation script are:

```
$ns rtmodel Exponential 0.8 1.0 1.0 $n1
$ns rtmodel Trace dynamics.trc $n2 $n3
$ns rtmodel Deterministic 20.0 20.0 $node(1) $node(5)
```

The procedure requires at least three arguments:

- The first two arguments define the model that will be used, and the parameters to configure the model. The currently implemented models in *ns* are Exponential (On/Off), Deterministic (On/Off), Trace (driven), or Manual (one-shot) models.
- The number, format, and interpretation of the configuration parameters is specific to the particular model.
  1. The exponential on/off model takes four parameters: `<[start time]>`, `<up interval>`, `<down interval>`, `<finish time>`. `<start time>` defaults to 0.5s. from the start of the simulation, `<finish time>` defaults to the end of the simulation. `<up interval>` and `<down interval>` specify the mean of the exponential distribution defining the time that the



node or link will be up and down respectively. The default up and down interval values are 10s. and 1s. respectively. Any of these values can be specified as “—” to default to the original value.

The following are example specifications of parameters to this model:

```
0.8 1.0 1.0 ; # start at 0.8s., up/down = 1.0s., finish is default
5.0 0.5 ; # start is default, up/down = 5.0s, 0.5s., finish is default
- 0.7 ; # start, up interval are default, down = 0.7s., finish is default
- - - 10 ; # start, up, down are default, finish at 10s.
```

2. The deterministic on/off model is similar to the exponential model above, and takes four parameters:  $\langle$ start time $\rangle$ , up interval, down interval,  $\langle$ finish time $\rangle$ .  $\langle$ start time $\rangle$  defaults to the start of the simulation,  $\langle$ finish time $\rangle$  defaults to the end of the simulation. Only the interpretation of the up and down interval is different;  $\langle$ up interval $\rangle$  and  $\langle$ down interval $\rangle$  specify the exact duration that the node or link will be up and down respectively. The default values for these parameters are:  $\langle$ start time $\rangle$  is 0.5s. from start of simulation,  $\langle$ up interval $\rangle$  is 2.0s.,  $\langle$ down interval $\rangle$  is 1.0s., and  $\langle$ finish time $\rangle$  is the duration of the simulation.
3. The trace driven model takes one parameter: the name of the trace file. The format of the input trace file is identical to that output by the dynamics trace modules, viz.,  $v \langle$ time $\rangle$  link- $\langle$ operation $\rangle$   $\langle$ node1 $\rangle$   $\langle$ node2 $\rangle$ . Lines that do not correspond to the node or link specified are ignored.

```
v 0.8123 link-up 3 5
v 3.5124 link-down 3 5
```

4. The manual one-shot model takes two parameters: the operation to be performed, and the time that it is to be performed.
- The rest of the arguments to the `rtmodel{}` procedure define the node or link that the model will be applied to. If only one node is specified, it is assumed that the node will fail. This is modeled by making the links incident on the node fail. If two nodes are specified, then the command assumes that the two are adjacent to each other, and the model is applied to the link incident on the two nodes. If more than two nodes are specified, only the first is considered, the subsequent arguments are ignored.
  - instance variable, `traceAllFile_` is set.

The command returns the handle to the model that was created in this call.

Internally, `rtmodel{}` stores the list of route models created in the class Simulator instance variable, `rtModel_`.

- The instance procedure `rtmodel-delete{}` takes the handle of a route model as argument, removes it from the `rtModel_` list, and deletes the route model.
- The instance procedure `rtmodel-at{}` is a special interface to the Manual model of network dynamics.

The command takes the time, operation, and node or link as arguments, and applies the operation to the node or link at the specified time. Example uses of this command are:

```
$ns rtmodel-at 3.5 up $n0
$ns rtmodel-at 3.9 up $n(3) $n(5)
$ns rtmodel-at 40 down $n4
```

Finally, the instance procedure `trace-dynamics{}` of the class `rtModel` enables tracing of the dynamics effected by this model. It is used as:

```
set fh [open "dyn.tr" w]
$rtmodel1 trace-dynamics $fh
$rtmodel2 trace-dynamics $fh
$rtmodel1 trace-dynamics stdout
```

In this example, `$rtmodel1` writes out trace entries to both `dyn.tr` and `stdout`; `$rtmodel2` only writes out trace entries to `dyn.tr`. A typical sequence of trace entries written out by either model might be:

```

v 0.8123 link-up 3 5
v 0.8123 link-up 5 3
v 3.5124 link-down 3 5
v 3.5124 link-down 5 3

```

These lines above indicate that Link  $\langle 3, 5 \rangle$  failed at 0.8123s., and recovered at time 3.5124s.

## 30.2 The Internal Architecture

Each model of network dynamics is implemented as a separate class, derived from the base class `rtModel`. We begin by describing the base class `rtModel` and the derived classes (Section 30.2.1). The network dynamics models use an internal queuing structure to ensure that simultaneous events are correctly handled, the class `rtQueue`. The next subsection (Section 30.2.2) describes the internals of this structure. Finally, we describe the extensions to the existing classes (Section 30.3.1): the Node, Link, and others.

### 30.2.1 The class `rtModel`

To use a new route model, the routine `rtmodel{}` creates an instance of the appropriate type, defines the node or link that the model will operate upon, configures the model, and possibly enables tracing; The individual instance procedures that accomplish this in pieces are:

The constructor for the base class stores a reference to the Simulator in its instance variable, `ns_`. It also initializes the `startTime_` and `finishTime_` from the class variables of the same name.

The instance procedure `set-elements` identifies the node or link that the model will operate upon. The command stores two arrays: `links_`, of the links that the model will act upon; `nodes_`, of the incident nodes that will be affected by the link failure or recovery caused by the model.

The default procedure in the base class to set the model configuration parameters is `set-parms`. It assumes a well defined start time, up interval, down interval, and a finish time, and sets up configuration parameters for some class of models. It stores these values in the instance variables: `startTime_`, `upInterval_`, `downInterval_`, `finishTime_`. The exponential and deterministic models use this default routine, the trace based and manual models define their own procedures.

The instance procedure `trace{}` enables `trace-dynamics{}` on each of the links that it affects. Additional details on `trace-dynamics{}` is discussed in the section on extensions to the class Link (Section 30.3.1).

The next sequence of configuration steps are taken just prior to the start of the simulator. `ns` invokes `rtmodel-configure{}` just before starting the simulation. This instance procedure first acquires an instance of the class `rtQueue`, and then invokes `configure{}` for each route model in its list, `rtModel_`.

The instance procedure `configure{}` makes each link that is applied to dynamic; this is the set of links stored in its instance variable array, `links_`. Then the procedure schedules its first event.

The default instance procedure `set-first-event{}` schedules the first event to take all the links “down” at `$startTime_ + upInterval_`. Individual types of route models derived from this base class should redefine this function.

Two instance procedures in the base class , `set-event{}` and `set-event-exact{}`, can be used to schedule events in the route queue.

`set-event{interval, operation}` schedules `operation` after `interval` seconds from the current time; it uses the procedure `set-event-exact{}` below.

`set-event-exact{fireTime, operation}` schedules `operation` to execute at `fireTime`.

If the time for execution is greater than the `finishTime_`, then the only possible action is to take a failed link “up”.

Finally, the base class provides the methods to take the links `up{}` or `down{}`. Each method invokes the appropriate procedure on each of the links in the instance variable, `links_`.

**Exponential** The model schedules its first event to take the links down at `startTime_ + E(upInterval_)`;

It also defines the procedures, `up{}` and `down{}`; each procedure invokes the base class procedure to perform the actual operation. This routine then reschedules the next event at `E(upInterval)` or `E(downInterval_)` respectively.

**Deterministic** The model defines the procedures, `up{}` and `down{}`; each procedure invokes the base class procedure to perform the actual operation. This routine then reschedules the next event at `upInterval` or `downInterval_` respectively.

**Trace** The model redefines the instance procedure `set-parms{}` to open a trace file, and set events based on that input.

The instance procedure `get-next-event{}` returns the next valid event from the trace file. A valid event is an event that is applicable to one of the links in this object's `links_` variable.

The instance procedure `set-trace-events{}` uses `get-next-event{}` to schedule the next valid event.

The model redefines `set-first-event{}`, `up{}`, and `down{}` to use `set-trace-events{}`.

**Manual** The model is designed to fire exactly once. The instance procedure `set-parms{}` takes an operation and the time to execute that operation as arguments. `set-first-event{}` will schedule the event at the appropriate moment.

This routine also redefines `notify{}` to delete the object instance when the operation is completed. This notion of the object deleting itself is fragile code.

Since the object only fires once and does not have to be rescheduled, it does not overload the procedures `up{}` or `down{}`.

### 30.2.2 class rtQueue

The simulator needs to co-ordinate multiple simultaneous network dynamics events, especially to ensure the right coherent behaviour. Hence, the network dynamics models use their own internal route queue to schedule dynamics events. There is one instance of this object in the simulator, in the class `Simulator` instance variable `rtq_`.

The queue object stores an array of queued operations in its instance variable, `rtq_`. The index is the time at which the event will execute. Each element is the list of operations that will execute at that time.

The instance procedures `insq{}` and `insq-i{}` can insert an element into the queue.

The first argument is the time at which this operation will execute. `insq{}` takes the exact time as argument; `insq-i{}` takes the interval as argument, and schedules the operation `interval` seconds after the current time.

The following arguments specify the object, `$obj`, the instance procedure of that object, `$iproc`, and the arguments to that procedure, `$args`.

These arguments are placed into the route queue for execution at the appropriate time.

The instance procedure `runq{}` executes `eval $obj $iproc $args` at the appropriate instant. After all the events for that instance are executed, `runq{}` will `notify{}` each object about the execution.

Finally, the instance procedure `delq{}` can remove a queued action with the time and the name of the object.

## 30.3 Interaction with Unicast Routing

In an earlier section, we had described how unicast routing reacts (Section 28.4.2) to changes to the topology. This section details the steps by which the network dynamics code will notify the nodes and routing about the changes to the topology.

1. `rtQueue::runq{}` will invoke the procedures specified by each of the route model instances. After all of the actions are completed, `runq{}` will notify each of the models.
2. `notify{}` will then invoke instance procedures at all of the nodes that were incident to the affected links. Each route model stores the list of nodes in its instance variable array, `nodes_`.

It will then notify the `RouteLogic` instance of topology changes.

3. The `rtModel` object invokes the class `Node` instance procedure `intf-changed{}` for each of the affected nodes.
4. `Node::intf-changed{}` will notify any `rtObject` at the node of the possible changes to the topology.

Recall that these route objects are created when the simulation uses detailed dynamic unicast routing.

### 30.3.1 Extensions to Other Classes

The existing classes assume that the topology is static by default. In this section, we document the necessary changes to these classes to support dynamic topologies.

We have already described the instance procedures in the class `Simulator` to create or manipulate route models, *i.e.*, `rtmodel{}`, `rtmodel-at{}`, `rtmodel-delete{}`, and `rtmodel-configure{}` in earlier sections (Section 30.2.1). Similarly, the class `Node` contains the instance procedure `intf-changed{}` that we described in the previous section (Section 30.3).

The network dynamics code operates on individual links. Each model currently translates its specification into operations on the appropriate links. The following paragraphs describe the class `Link` and related classes.

**class `DynamicLink`** This class is the only `TclObject` in the network dynamics code. The shadow class is called `class DynaLink`. The class supports one bound variable, `status_`. `status_` is 1 when the link is up, and 0 when the link is down. The shadow object's `recv()` method checks the `status_` variable, to decide whether or not a packet should be forwarded.

**class Link** This class supports the primitives: `up` and `down`, and `up?` to set and query `status_`. These primitives are instance procedures of the class.

The instance procedures `up{}` and `down{}` set `status_` to 1 and 0 respectively.

In addition, when the link fails, `down{}` will reset all connectors that make up the link. Each connector, including all queues and the delay object will flush and drop any packets that it currently stores. This emulates the packet drop due to link failure.

Both procedures then write trace entries to each file handle in the list, `dynT_`.

The instance procedure `up?{}` returns the current value of `status_`.

In addition, the class contains the instance procedure `all-connectors{}`. This procedure takes an operation as argument, and applies the operation uniformly to all of the class instance variables that are handles for TclObjects.

**class SimpleLink** The class supports two instance procedures `dynamic{}` and `trace-dynamics{}`. We have already described the latter procedure when describing the `trace{}` procedure in the class `rtModel`.

The instance procedure `dynamic{}` inserts a `DynamicLink` object (Section 6.2) at the head of the queue. It points the down-target of the object to the drop target of the link, `drpt_`, if the object is defined, or to the `nullAgent_` in the simulator. It also signals each connector in the link that the link is now dynamic.

Most connectors ignore this signal to become dynamic; the exception is `DelayLink` object. This object will normally schedule each packet it receives for reception by the destination node at the appropriate time. When the link is dynamic, the object will queue each packet internally; it schedules only one event for the next packet that will be delivered, instead of one event per packet normally. If the link fails, the route model will signal a `reset`, at which point, the shadow object will execute its `reset` instproc-like, and flush all packets in its internal queue. Additional details about the `DelayLink` can be found in another chapter (Chapter 8).

## 30.4 Deficiencies in the Current Network Dynamics API

There are a number of deficiencies in the current API that should be changed in the next iteration:

1. There is no way to specify a cluster of nodes or links that behave in lock-step dynamic synchrony.
2. Node failure should be dealt with as its own mechanism, rather than a second grade citizen of link failure. This shows up in a number of situations, such as:
  - (a) The method of emulating node failure as the failure of the incident links is broken. Ideally, node failure should cause all agents incident on the node to be reset.
  - (b) There is no tracing associated with node failure.
3. If two distinct route models are applied to two separate links incident on a common node, and the two links experience a topology change at the same instant, then the node will be notified more than once.

## 30.5 Commands at a glance

Following is a list of commands used to simulate dynamic scenarios in *ns*:

`$ns_ rtmodel <model> <model-params> <args>`

This command defines the dynamic model (currently implemented models are: Deterministic, Exponential, Manual or Trace) to be applied to nodes and links in the topology. The first two arguments consists of the rtmodel and the parameter to configure the model. <args> stands for different type of arguments expected with different dynamic model types. This returns a handle to a model object corresponding to the specified model.

- In the Deterministic model <model-params> is <start-time>, <up-interval>, <down-interval>, <finish-time>. Starting from start-time the link is made up for up-interval and down for down-interval till finish-time is reached. The default values for start-time, up-interval, downinterval are 0.5s, 2.0s, 1.0s respectively. finish-time defaults to the end of the simulation. The start-time defaults to 0.5s in order to let the routing protocol computation quiesce.
- If the Exponential model is used model-params is of the form <up-interval>, <down-interval> where the link up-time is an exponential distribution around the mean upinterval and the link down-time is an exponential distribution around the mean down-interval. Default values for up-interval and down-interval are 10s and 1s respectively.
- If the Manual distribution is used model-params is <at> <op> where at specifies the time at which the operation op should occur. op is one of up, down. The Manual distribution could be specified alternately using the rtmodel-at method described later in the section.
- If Trace is specified as the model the link/node dynamics is read from a Tracefile. The model-params argument would in this case be the file-handle of the Tracefile that has the dynamics information. The tracefile format is identical to the trace output generated by the trace-dynamics link method (see TRACE AND MONITORING METHODS SECTION).

`$ns_ rtmodel-delete <model>`

This command takes the handle of the routemodel <model> as an argument, removes it from the list of rtmodels maintained by simulator and deletes the model.

`$ns_ rtmodel-at <at> <op> <args>`

This command is a special interface to the Manual model of network dynamics. It takes the time <at>, type of operation <op> and node or link on which to apply the operation <args> as the arguments. At time <at>, the operation <op> which maybe up or down is applied to a node or link.

`$rtmodel trace <ns> <f> <optional:op>`

This enables tracing of dynamics effected by this model in the links. <ns> is an instance of the simulator, <f> the output file to write the traces to and <op> is an optional argument that may be used to define a type of operation (like nam). This is a wrapper for the class Link procedure trace-dynamics.

`$link trace-dynamics <ns> <f> <optional:op>`

This is a class link instance procedure that is used to setup tracing of dynamics in that particular link. The arguments are same as that of class rtModel's procedure trace described above.

`$link dynamic`

This command inserts a DynamicLink object at the head of the queue and signals to all connectors in the link that the link is now dynamic.

Internal procedures:

`$ns_ rtmodel-configure`

This is an internal procedure that configures all dynamic models that are present in the list of models maintained by the simulator.

# Chapter 31

## Hierarchical Routing

This chapter describes the internals of hierarchical routing implemented in *ns*. This chapter consists of two sections. In the first section we give an overview of hierarchical routing. In the second section we walk through the API's used for setting hierarchical routing and describe the architecture, internals and code path for hier rtg in the process.

The functions and procedures described in this chapter can be found in `~ns/tcl/lib/ns-hiernode.tcl`, `tcl/lib/ns-address.tcl`, `tcl/lib/ns-route.tcl` and `route.cc`.

### 31.1 Overview of Hierarchical Routing

Hierarchical routing was mainly devised, among other things, to reduce memory requirements of simulations over very large topologies. A topology is broken down into several layers of hierarchy, thus downsizing the routing table. The table size is reduced from  $n^2$ , for flat routing, to about  $\log n$  for hierarchical routing. However some overhead costs results as number of hierarchy levels are increased. Optimum results were found for 3 levels of hierarchy and the current ns implementation supports upto a maximum of 3 levels of hierarchical routing.

To be able to use hierarchical routing for the simulations, we need to define the hierarchy of the topology as well as provide the nodes with hierarchical addressing. In flat routing, every node knows about every other node in the topology, thus resulting in routing table size to the order of  $n^2$ . For hierarchical routing, each node knows only about those nodes in its level. For all other destinations outside its level it forwards the packets to the border router of its level. Thus the routing table size gets downsized to the order of about  $\log n$ .

### 31.2 Usage of Hierarchical routing

Hierarchical routing requires some additional features and mechanisms for the simulation. For example, a new node object called *HierNode* is been defined for hier rtg. Therefore the user must specify hierarchical routing requirements before creating topology. This is done as shown below:

First, the address format ( ) or the address space used for node and port address, needs to be set in the hierarchical mode. It may be done in one of the two ways:

```
set ns [new Simulator]
```

```
$ns set-address-format hierarchical
```

This sets the node address space to a 3 level hierarchy assigning 8 bits in each level.

or,

```
$ns set-address-format hierarchical <n hierarchy levels> <# bits in  
level 1> ...<# bits in nth level>
```

which creates a node address space for n levels of hierarchy assigning bits as specified for every level.

This other than creating a hierarchical address space also sets a flag called *EnableHierRt\_* and sets the Simulator class variable *node\_factory\_* to *HierNode*. Therefore when nodes are created by calling Simulator method “node” as in :

\$ns node 0.0.1, a *HierNode* is created with an address of 0.0.1;

Class *AddrParams* is used to store the topology hierarchy like number of levels of hierarchy, number of areas in each level like number of domains, number of clusters and number of nodes in each cluster.

The API for supplying these information to *AddrParams* is shown below:

```
AddrParams set domain_num_ 2  
lappend cluster_num 2 2  
AddrParams set cluster_num_ $cluster_num  
lappend eilastlevel 2 3 2 3  
AddrParams set nodes_num_ $eilastlevel
```

This defines a topology with 2 domains, say D1 and D2 with 2 clusters each (C11 & C12 in D1 and C21 & C22 in D2). Then number of nodes in each of these 4 clusters is specified as 2,3,2 and 3 respectively.

The default values used by *AddrParams* provide a topology with a single domain with 4 clusters, with each cluster consisting of 5 nodes.

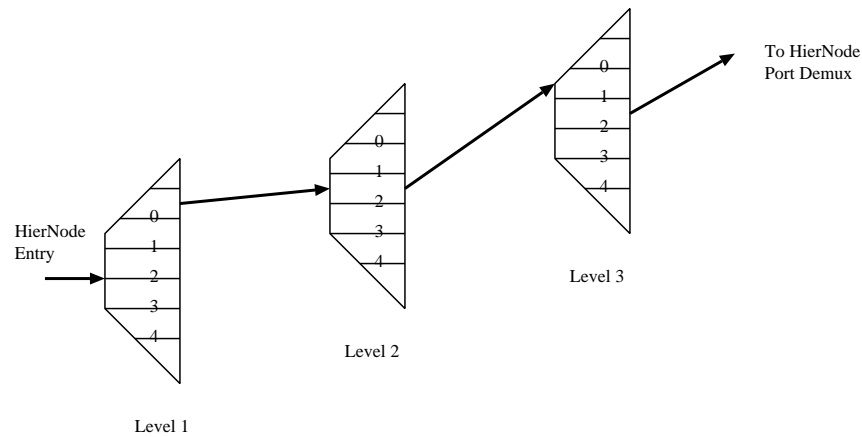
Appropriate mask and shift values are generated by *AddrParams* for the hierarchical node address space.

Each *HierNode* at the time of its creation calls the method ‘mk-default-classifier’ to setup n numbers of address classifiers for n levels of hierarchy defined in the topology.

```
HierNode instproc mk-default-classifier  
    $self instvar np_ id_ classifiers_ agents_ dmux_ neighbor_ address_  
    # puts "id=$id_"  
    set levels [AddrParams set hlevel_]  
    for set n 1 $n <= $levels incr n  
        set classifiers_($n) [new Classifier/Addr]  
        $classifiers_($n) set mask_ [AddrParams set NodeMask_($n)]  
        $classifiers_($n) set shift_ [AddrParams set NodeShift_($n)]
```

At the time of route computation, a call is made to *add-route*. *add-route* populates classifiers as shown in the *otcl* method below:





3-Level classifiers for HierNode (hier-addr:0.2.1)

Figure 31.1: Hierarchical classifiers

```
Node instproc add-route dst target
    $self instvar rtnotif_
    # Notify every module that is interested about this
    # route installation

    if $rtnotif_ != ""
    $rtnotif_ add-route $dst $target

    $self incr-rtgtable-size
```

For an example of 3 level of hierarchy, the level 1 classifier demuxes for domains, level 2 for all clusters inside the node's domain and finally classifier 3 demuxes for all nodes in the particular cluster that the node itself resides. For such a topology, a HierNode with address of 0.1.2 looks like the figure below:

Thus the size of the routing tables are considerably reduced from  $n^2$  as seen for flat routing where each node had to store the next\_hop info of all other nodes in the topology. Instead, for hierarchical routing, a given node needs to know about its neighbours in its own cluster, about the all clusters in its domain and about all the domains. This saves on memory consumption as well as run-time for the simulations using several thousands of nodes in their topology.

### 31.3 Creating large Hierarchical topologies

The previous section describes methods to create hierarchical topologies by hand. However, there is a script available in ns that converts Georgia-tech's SGB-graphs into ns compatible hierarchical topologies. Please refer to <http://www-mash.CS.Berkeley.EDU/ns/ns-topogen.html> for downloading as well as instructions on using the hierarchical converter package.

See hier-rtg-10.tcl and hier-rtg-100.tcl in `~ns/tcl/ex` for example scripts of hier routing on small and large topologies respectively.

## 31.4 Hierarchical Routing with SessionSim

Hierarchical routing may be used in conjunction with Session simulations (see Chapter 41). Session-level simulations which are used for running multicast simulations over very large topologies, gains additionally in terms of memory savings if used with hierarchical routing. See simulation script `~ns/tcl/ex/newmcast/session-hier.tcl` for an example of sessionsim over hier rtg.

## 31.5 Commands at a glance

Following is a list of hierarchical routing/addressing related commands used in simulation scripts:

```
$ns_ set-address-format hierarchical
```

This command was used to setup hierarchical addressing in *ns*. However with the recent changes in node APIs, this command has been replaced by

```
ns_ node-config -addressType hierarchical
```

This creates a default topology of 3 levels of hierarchy, assigning 8 bits to each level.

```
$ns_ set-address-format hierarchical <nlevels> <#bits in level1>....<#bits in level  
n>
```

This command creates a hierarchy of *<nlevels>* and assigns the bits in each level as specified in the arguments.

```
AddrParams set domain_num_ <n_domains>
```

```
AddrParams set cluster_num_ <n_clusters>
```

```
AddrParams set nodes_num_ <n_nodes>
```

The above APIs are used to specify the hierarchical topology, i.e the number of domains, clusters and nodes present in the topology. Default values used by AddrParams (i.e if nothing is specified) provide a topology with a single domain with 4 clusters, with each cluster consisting of 5 nodes.

Internal procedures:

```
$Node add-route <dst> <target>
```

This procedure is used to add next-hop entries of a destination *<dst>* for a given *<target>*.

```
$hiernode_ split-addrstr <str>
```

This splits up a hierarchical address string (say a.b.c) into a list of the addresses at each level (i.e, a,b and c).

## **Part V**

# **Transport**

# Chapter 32

## UDP Agents

### 32.1 UDP Agents

UDP agents are implemented in `udp.cc`. A UDP agent accepts data in variable size chunks from an application, and segments the data if needed. UDP packets also contain a monotonically increasing sequence number and an RTP timestamp. Although real UDP packets do not contain sequence numbers or timestamps, this sequence number does not incur any simulated overhead, and can be useful for traffic analysis or for simulating UDP-based applications.

The default maximum segment size (MSS) for UDP agents is 1000 byte:

```
Agent/UDP set packetSize_ 1000 ; # max segment size
```

This OTcl instvar is bound to the C++ agent variable `size_`.

Applications can access UDP agents via the `sendmsg()` function in C++, or via the `send` or `sendmsg` methods in OTcl, as described in section 37.2.3.

The following is a simple example of how a UDP agent may be used in a program. In the example, the CBR traffic generator is started at time 1.0, at which time the generator begins to periodically call the UDP agent `sendmsg()` function.

```
set ns [new Simulator]
set n0 [$ns node]
set n1 [$ns node]
$ns duplex-link $n0 $n1 5Mb 2ms DropTail

set udp0 [new Agent/UDP]
$ns attach-agent $n0 $udp0
set cbr0 [new Application/Traffic/CBR]
$cbr0 attach-agent $udp0
$udp0 set packetSize_ 536 ; # set MSS to 536 bytes

set null0 [new Agent/Null]
$ns attach-agent $n1 $null0
$ns connect $udp0 $null0
$ns at 1.0 "$cbr0 start"
```

## 32.2 Commands at a glance

The following commands are used to setup UDP agents in simulation scripts:

```
set udp0 [new Agent/UDP]
```

This creates an instance of the UDP agent.

```
$ns_ attach-agent <node> <agent>
```

This is a common command used to attach any <agent> to a given <node>.

```
$traffic-gen attach-agent <agent>
```

This is a class Application/Traffic/<traffic type> method which connects the traffic generator to the given <agent>. For example, if we want to setup a CBR traffic flow for the udp agent, udp1, we give the following commands

```
set cbr1 [new Application/Traffic/CBR]
```

```
$cbr1 attach-agent $udp1
```

```
$ns_ connect <src-agent> <dst-agent>
```

This command sets up an end-to-end connection between two agents (at the transport layer).

```
$udp set packetSize_ <pktsize>
```

```
$udp set dst_addr_ <address>
```

```
$udp set dst_port_ <portnum>
```

```
$udp set class_ <class-type>
```

```
$udp set ttl_ <time-to-live>
```

```
..... etc
```

The above are different parameter values that may be set as shown above for udp agents. The default values can be found in *ns/tcl/lib/ns-default.tcl*.

For a typical example of setting up an UDP agent used in a simulation, see the above section 32.1.

# Chapter 33

## TCP Agents

This section describes the operation of the TCP agents in *ns*. There are two major types of TCP agents: one-way agents and a two-way agent. One-way agents are further subdivided into a set of TCP senders (which obey different congestion and error control techniques) and receivers (“sinks”). The two-way agent is symmetric in the sense that it represents both a sender and receiver. It is still under development.

The files described in this section are too numerous to enumerate here. Basically it covers most files matching the regular expression `~ns/tcp*.{cc,h}`.

The one-way TCP sending agents currently supported are:

- Agent/TCP - a “tahoe” TCP sender
- Agent/TCP/Reno - a “Reno” TCP sender
- Agent/TCP/Newreno - Reno with a modification
- Agent/TCP/Sack1 - TCP with selective repeat (follows RFC2018)
- Agent/TCP/Vegas - TCP Vegas
- Agent/TCP/Fack - Reno TCP with “forward acknowledgment”

The one-way TCP receiving agents currently supported are:

- Agent/TCPSink - TCP sink with one ACK per packet
- Agent/TCPSink/DelAck - TCP sink with configurable delay per ACK
- Agent/TCPSink/Sack1 - selective ACK sink (follows RFC2018)
- Agent/TCPSink/Sack1/DelAck - Sack1 with DelAck

The two-way experimental sender currently supports only a Reno form of TCP:

- Agent/TCP/FullTcp

The section comprises three parts: the first part is a simple overview and example of configuring the base TCP send/sink agents (the sink requires no configuration). The second part describes the internals of the base send agent, and last part is a description of the extensions for the other types of agents that have been included in the simulator.

## 33.1 One-Way TCP Senders

The simulator supports several versions of an abstracted TCP sender. These objects attempt to capture the essence of the TCP congestion and error control behaviors, but are not intended to be faithful replicas of real-world TCP implementations. They do not contain a dynamic window advertisement, they do segment number and ACK number computations entirely in packet units, there is no SYN/FIN connection establishment/teardown, and no data is ever transferred (e.g. no checksums or urgent data).

### 33.1.1 The Base TCP Sender (Tahoe TCP)

The “Tahoe” TCP agent `Agent/TCP` performs congestion control and round-trip-time estimation in a way similar to the version of TCP released with the 4.3BSD “Tahoe” UN\*X system release from UC Berkeley. The congestion window is increased by one packet per new ACK received during slow-start (when  $cwnd_ < ssthresh_$ ) and is increased by  $\frac{1}{cwnd_}$  for each new ACK received during congestion avoidance (when  $cwnd_ \geq ssthresh_$ ).

**Responses to Congestion** Tahoe TCP assumes a packet has been lost (due to congestion) when it observes `NUMDUPACKS` (defined in `tcp.h`, currently 3) duplicate ACKs, or when a retransmission timer expires. In either case, Tahoe TCP reacts by setting `ssthresh_` to half of the current window size (the minimum of `cwnd_` and `window_`) or 2, whichever is larger. It then initializes `cwnd_` back to the value of `windowInit_`. This will typically cause the TCP to enter slow-start.

**Round-Trip Time Estimation and RTO Timeout Selection** Four variables are used to estimate the round-trip time and set the retransmission timer: `rtt_`, `srtt_`, `rttvar_`, `tcpTick_`, and `backoff_`. TCP initializes `rttvar` to  $3/tcpTick_$  and `backoff` to 1. When any future retransmission timer is set, its timeout is set to the current time plus  $\max(bt(a + 4v + 1), 64)$  seconds, where  $b$  is the current backoff value,  $t$  is the value of `tcpTick`,  $a$  is the value of `srtt`, and  $v$  is the value of `rttvar`.

Round-trip time samples arrive with new ACKs. The RTT sample is computed as the difference between the current time and a “time echo” field in the ACK packet. When the first sample is taken, its value is used as the initial value for `srtt_`. Half the first sample is used as the initial value for `rttvar_`. For subsequent samples, the values are updated as follows:

$$srtt = \frac{7}{8} \times srtt + \frac{1}{8} \times sample$$
$$rttvar = \frac{3}{4} \times rttvar + \frac{1}{4} \times |sample - srtt|$$

### 33.1.2 Configuration

Running an TCP simulation requires creating and configuring the agent, attaching an application-level data source (a traffic generator), and starting the agent and the traffic generator.

### 33.1.3 Simple Configuration

#### Creating the Agent

```

set ns [new Simulator]                                ; # preamble initialization
set node1 [$ns node]                                  ; # agent to reside on this node
set node2 [$ns node]                                  ; # agent to reside on this node

set tcp1 [$ns create-connection TCP $node1 TCPSink $node2 42]
$tcp set window_ 50                                    ; # configure the TCP agent

set ftp1 [new Application/FTP]
$ftp1 attach-agent $tcp1

$ns at 0.0 "$ftp start"

```

This example illustrates the use of the simulator built-in function `create-connection`. The arguments to this function are: the source agent to create, the source node, the target agent to create, the target node, and the flow ID to be used on the connection. The function operates by creating the two agents, setting the flow ID fields in the agents, attaching the source and target agents to their respective nodes, and finally connecting the agents (i.e. setting appropriate source and destination addresses and ports). The return value of the function is the name of the source agent created.

**TCP Data Source** The TCP agent does not generate any application data on its own; instead, the simulation user can connect any traffic generation module to the TCP agent to generate data. Two applications are commonly used for TCP: FTP and Telnet. FTP represents a bulk data transfer of large size, and telnet chooses its transfer sizes randomly from `tcplib` (see the file `tcplib-telnet.cc`). Details on configuring these application source objects are in Section 37.4.

### 33.1.4 Other Configuration Parameters

In addition to the `window_` parameter listed above, the TCP agent supports additional configuration variables. Each of the variables described in this subsection is both a class variable and an instance variable. Changing the class variable changes the default value for all agents that are created subsequently. Changing the instance variable of a particular agent only affects the values used by that agent. For example,

```

Agent/TCP set window_ 100                                ; # Changes the class variable
$tcp set window_ 2.0                                     ; # Changes window_ for the $tcp object only

```

The default parameters for each TCP agent are:

```

Agent/TCP set window_ 20                                ; # max bound on window size
Agent/TCP set windowInit_ 1                             ; # initial/reset value of cwnd
Agent/TCP set windowOption_ 1                           ; # cong avoid algorithm (1: standard)
Agent/TCP set windowConstant_ 4                         ; # used only when windowOption != 1
Agent/TCP set windowThresh_ 0.002                       ; # used in computing averaged window
Agent/TCP set overhead_ 0                               ; # !=0 adds random time between sends
Agent/TCP set ecn_ 0                                    ; # TCP should react to ecn bit
Agent/TCP set packetSize_ 1000                          ; # packet size used by sender (bytes)
Agent/TCP set bugFix_ true                               ; # see explanation
Agent/TCP set slow_start_restart_ true                  ; # see explanation
Agent/TCP set tcpTick_ 0.1                              ; # timer granularity in sec (.1 is NONSTANDARD)
Agent/TCP set maxrto_ 64                                ; # bound on RTO (seconds)
Agent/TCP set dupacks_ 0                                ; # duplicate ACK counter
Agent/TCP set ack_ 0                                    ; # highest ACK received

```



```

Agent/TCP set cwnd_ 0 ; # congestion window (packets)
Agent/TCP set awnd_ 0 ; # averaged cwnd (experimental)
Agent/TCP set ssthresh_ 0 ; # slow-start threshold (packets)
Agent/TCP set rtt_ 0 ; # rtt sample
Agent/TCP set srtt_ 0 ; # smoothed (averaged) rtt
Agent/TCP set rttvar_ 0 ; # mean deviation of rtt samples
Agent/TCP set backoff_ 0 ; # current RTO backoff factor
Agent/TCP set maxseq_ 0 ; # max (packet) seq number sent

```

For many simulations, few of the configuration parameters are likely to require modification. The more commonly modified parameters include: `window_` and `packetSize_`. The first of these bounds the window TCP uses, and is considered to play the role of the receiver's advertised window in real-world TCP (although it remains constant). The packet size essentially functions like the MSS size in real-world TCP. Changes to these parameters can have a profound effect on the behavior of TCP. Generally, those TCPs with larger packet sizes, bigger windows, and smaller round trip times (a result of the topology and congestion) are more aggressive in acquiring network bandwidth.

### 33.1.5 Other One-Way TCP Senders

**Reno TCP** The Reno TCP agent is very similar to the Tahoe TCP agent, except it also includes *fast recovery*, where the current congestion window is “inflated” by the number of duplicate ACKs the TCP sender has received before receiving a new ACK. A “new ACK” refers to any ACK with a value higher than the highest seen so far. In addition, the Reno TCP agent does not return to slow-start during a fast retransmit. Rather, it reduces sets the congestion window to half the current window and resets `ssthresh_` to match this value.

**Newreno TCP** This agent is based on the Reno TCP agent, but which modifies the action taken when receiving new ACKs. In order to exit fast recovery, the sender must receive an ACK for the highest sequence number sent. Thus, new “partial ACKs” (those which represent new ACKs but do not represent an ACK for all outstanding data) do not deflate the window (and possibly lead to a stall, characteristic of Reno).

**Vegas TCP** This agent implements “Vegas” TCP ([4, 5]). It was contributed by Ted Kuo.

**Sack TCP** This agent implements selective repeat, based on selective ACKs provided by the receiver. It follows the ACK scheme described in [23], and was developed with Matt Mathis and Jamshid Mahdavi.

**Fack TCP** This agent implements “forward ACK” TCP, a modification of Sack TCP described in [22].

## 33.2 TCP Receivers (sinks)

The TCP senders described above represent one-way data senders. They must peer with a “TCP sink” object.

### 33.2.1 The Base TCP Sink

The base TCP sink object (`Agent/TCPSink`) is responsible for returning ACKs to a peer TCP source object. It generates one ACK per packet received. The size of the ACKs may be configured. The creation and configuration of the TCP sink object is generally performed automatically by a library call (see `create-connection` above).

#### configuration parameters

```
Agent/TCPSink set packetSize_ 40
```

### 33.2.2 Delayed-ACK TCP Sink

A delayed-ACK sink object (`Agent/Agent/TCPSink/DelAck`) is available for simulating a TCP receiver that ACKs less than once per packet received. This object contains a bound variable `interval_` which gives the number of seconds to wait between ACKs. The delayed ACK sink implements an aggressive ACK policy whereby only ACKs for in-order packets are delayed. Out-of-order packets cause immediate ACK generation.

#### configuration parameters

```
Agent/TCPSink/DelAck set interval_ 100ms
```

### 33.2.3 Sack TCP Sink

The selective-acknowledgment TCP sink (`Agent/TCPSink/Sack1`) implements SACK generation modeled after the description of SACK in RFC 2018. This object includes a bound variable `maxSackBlocks_` which gives the maximum number of blocks of information in an ACK available for holding SACK information. The default value for this variable is 3, in accordance with the expected use of SACK with RTTM (see RFC 2018, section 3). Delayed and selective ACKs together are implemented by an object of type `Agent/TCPSink/Sack1/DelAck`.

#### configuration parameters

```
Agent/TCPSink set maxSackBlocks_ 3
```

## 33.3 Two-Way TCP Agents (FullTcp)

The `Agent/TCP/FullTcp` object is a new addition to the suite of TCP agents supported in the simulator and is still under development. It is different from (and incompatible with) the other agents, but does use some of the same architecture. It differs from these agents in the following ways: following ways:

- connections may be established and torn down (SYN/FIN packets are exchanged)
- bidirectional data transfer is supported

- sequence numbers are in bytes rather than packets

The generation of SYN packets (and their ACKs) can be of critical importance in trying to model real-world behavior when using many very short data transfers. This version of TCP currently defaults to sending data on the 3rd segment of an initial 3-way handshake, a behavior somewhat different than common real-world TCP implementations. A “typical” TCP connection proceeds with an active opener sending a SYN, the passive opener responding with a SYN+ACK, the active opener responding with an ACK, and then some time later sending the first segment with data (corresponding to the first application write). Thus, this version of TCP sends data at a time somewhat earlier than typical implementations. This TCP can also be configured to send data on the initial SYN segment. Future changes to FullTCP may include a modification to send the first data segment later, and possibly to implement T/TCP functionality.

Currently FullTCP is only implemented with Reno congestion control, but ultimately it should be available with the full range of congestion control algorithms (e.g., Tahoe, SACK, Vegas, etc.).

### 33.3.1 Simple Configuration

Running an Full TCP simulation requires creating and configuring the agent, attaching an application-level data source (a traffic generator), and starting the agent and the traffic generator.

#### Creating the Agent

```
# set up connection (do not use "create-connection" method because
# we need a handle on the sink object)
set src [new Agent/TCP/FullTcp] ; # create agent
set sink [new Agent/TCP/FullTcp] ; # create agent
$ns_ attach-agent $node_(s1) $src ; # bind src to node
$ns_ attach-agent $node_(k1) $sink ; # bind sink to node
$src set fid_ 0 ; # set flow ID field
$sink set fid_ 0 ; # set flow ID field
$ns_ connect $src $sink ; # active connection src to sink

# set up TCP-level connections
$sink listen ; # will figure out who its peer is
$src set window_ 100;
```

The creation of the FullTcp agent is similar to the other agents, but the sink is placed in a listening state by the `listen` method. Because a handle to the receiving side is required in order to make this call, the `create-connection` call used above cannot be used.

**Configuration Parameters** The following configuration parameters are available through Tcl for the FullTcp agent:

```
Agent/TCP/FullTcp set segsperack_ 1 ; # segs received before generating ACK
Agent/TCP/FullTcp set segsize_ 536 ; # segment size (MSS size for bulk xfers)
Agent/TCP/FullTcp set tcprexmtthresh_ 3 ; # dupACKs thresh to trigger fast rexmt
Agent/TCP/FullTcp set iss_ 0 ; # initial send sequence number
Agent/TCP/FullTcp set nodelay_ false ; # disable sender-side Nagle algorithm
Agent/TCP/FullTcp set data_on_syn_ false ; # send data on initial SYN?
```

```

Agent/TCP/FullTcp set dupseg_fix_ true ; # avoid fast rxt due to dup segs+acks
Agent/TCP/FullTcp set dupack_reset_ false ; # reset dupACK ctr on !0 len data segs containing dup ACKs
Agent/TCP/FullTcp set interval_ 0.1 ; # as in TCP above, (100ms is non-std)

```

### 33.3.2 BayFullTcp

A different implementation of two-way TCP has been ported into ns from Kathy Nicholes/Van Jacobson's group. It is called BayFullTcp. The basic difference between BayFullTcp and FullTcp (the two-way tcp version already present in ns) are as follows:

- BayTcp supports a client-server application model while FullTcp makes no assumption about its application layer.
- The tcp-application interface is different for both;
- FullTcp supports partial ack (BayTcp doesn't).
- FullTcp supports different flavors of tcp (tahoe, reno etc) which is not the case for baytcp.
- Both implementations have different set of API's .

There might be other finer differences between the two as well. One of our future plans is to redefine the APIs to allow fulltcp to use baytcp's client-server model.

## 33.4 Architecture and Internals

The base TCP agent (class Agent/TCP) is constructed as a collection of routines for sending packets, processing ACKs, managing the send window, and handling timeouts. Generally, each of these routines may be over-ridden by a function with the same name in a derived class (this is how many of the TCP sender variants are implemented).

**The TCP header** The TCP header is defined by the `hdr_tcp` structure in the file `~ns/tcp.h`. The base agent only makes use of the following subset of the fields:

```

ts_ ; # current time packet was sent from source */
ts_echo_ ; # for ACKs: timestamp field from packet associated with this ACK */
seqno_ ; # sequence number for this data segment or ACK (Note: overloading!) */
reason_ ; # set by sender when (re)transmitting to trace reason for send */

```

**Functions for Sending Data** Note that generally the sending TCP never actually sends data (it only sets the packet size).

**send\_much(force, reason, maxburst)** - this function attempts to send as many packets as the current sent window allows. It also keeps track of how many packets it has sent, and limits to the total to *maxburst*.

The function `output(seqno, reason)` sends one packet with the given sequence number and updates the maximum sent sequence number variable (`maxseq_`) to hold the given sequence number if it is the greatest sent so far. This function also assigns the various fields in the TCP header (sequence number, timestamp, reason for transmission). This function also sets a retransmission timer if one is not already pending.

**Functions for Window Management** The usable send window at any time is given by the function **window()**. It returns the minimum of the congestion window and the variable `wnd_`, which represents the receiver's advertised window.

**opencwnd()** - this function opens the congestion window. It is invoked when a new ACK arrives. When in slow-start, the function merely increments `cwnd_` by each ACK received. When in congestion avoidance, the standard configuration increments `cwnd_` by its reciprocal. Other window growth options are supported during congestion avoidance, but they are experimental (and not documented; contact Sally Floyd for details).

**closecwnd(int how)** - this function reduces the congestion window. It may be invoked in several ways: when entering fast retransmit, due to a timer expiration, or due to a congestion notification (ECN bit set). Its argument `how` indicates how the congestion window should be reduced. The value **0** is used for retransmission timeouts and fast retransmit in Tahoe TCP. It typically causes the TCP to enter slow-start and reduce `ssthresh_` to half the current window. The value **1** is used by Reno TCP for implementing fast recovery (which avoids returning to slow-start). The value **2** is used for reducing the window due to an ECN indication. It resets the congestion window to its initial value (usually causing slow-start), but does not alter `ssthresh_`.

**Functions for Processing ACKs** **recv()** - this function is the main reception path for ACKs. Note that because only one direction of data flow is in use, this function should only ever be invoked with a pure ACK packet (i.e. no data). The function stores the timestamp from the ACK in `ts_peer_`, and checks for the presence of the ECN bit (reducing the send window if appropriate). If the ACK is a new ACK, it calls **newack()**, and otherwise checks to see if it is a duplicate of the last ACK seen. If so, it enters fast retransmit by closing the window, resetting the retransmission timer, and sending a packet by calling `send_much`.

**newack()** - this function processes a "new" ACK (one that contains an ACK number higher than any seen so far). The function sets a new retransmission timer by calling **newtimer()**, updates the RTT estimation by calling **rtt\_update**, and updates the highest and last ACK variables.

**Functions for Managing the Retransmission Timer** These functions serve two purposes: estimating the round-trip time and setting the actual retransmission timer. **rtt\_init** - this function initializes `srtt_` and `rtt_` to zero, sets `rttvar_` to `3/tcp_tick_`, and sets the backoff multiplier to one.

**rtt\_timeout** - this function gives the timeout value in seconds that should be used to schedule the next retransmission timer. It computes this based on the current estimates of the mean and deviation of the round-trip time. In addition, it implements Karn's exponential timer backoff for multiple consecutive retransmission timeouts.

**rtt\_update** - this function takes as argument the measured RTT and averages it in to the running mean and deviation estimators according to the description above. Note that `t_srtt_` and `t_rttvar_` are both stored in fixed-point (integers). They have 3 and 2 bits, respectively, to the right of the binary point.

**reset\_rtx\_timer** - This function is invoked during fast retransmit or during a timeout. It sets a retransmission timer by calling `set_rtx_timer` and if invoked by a timeout also calls `rtt_backoff`.

**rtt\_backoff** - this function backs off the retransmission timer (by doubling it).

**newtimer** - this function called only when a new ACK arrives. If the sender's left window edge is beyond the ACK, then `set_rtx_timer` is called, otherwise if a retransmission timer is pending it is cancelled.

## 33.5 Tracing TCP Dynamics

The behavior of TCP is often observed by constructing a sequence number-vs-time plot. Typically, a trace is performed by enabling tracing on a link over which the TCP packets will pass. Two trace methods are supported: the default one (used for tracing TCP agents), and an extension used only for FullTcp.

## 33.6 One-Way Trace TCP Trace Dynamics

TCP packets generated by one of the one-way TCP agents and destined for a TCP sink agent passing over a traced link (see section 25) will generate a trace file lines of the form:

```
+ 0.94176 2 3 tcp 1000 ----- 0 0.0 3.0 25 40
+ 0.94276 2 3 tcp 1000 ----- 0 0.0 3.0 26 41
d 0.94276 2 3 tcp 1000 ----- 0 0.0 3.0 26 41
+ 0.95072 2 0 ack 40 ----- 0 3.0 0.0 14 29
- 0.95072 2 0 ack 40 ----- 0 3.0 0.0 14 29
- 0.95176 2 3 tcp 1000 ----- 0 0.0 3.0 21 36
+ 0.95176 2 3 tcp 1000 ----- 0 0.0 3.0 27 42
```

The exact format of this trace file is given in section 25.4. When tracing TCP, packets of type `tcp` or `ack` are relevant. Their type, size, sequence number (ack number for ack packets), and arrival/depart/drop time are given by field positions 5, 6, 11, and 2, respectively. The `+` indicates a packet arrival, `d` a drop, and `-` a departure. A number of scripts process this file to produce graphical output or statistical summaries (see, for example, `~ns/test-suite.tcl`, the `finish` procedure).

## 33.7 One-Way Trace TCP Trace Dynamics

TCP packets generated by FullTcp and passing over a traced link contain additional information not displayed by default using the regular trace object. By enabling the flag `show_tcphdr_` on the trace object (see section `refsec:traceformat`), 3 additional header fields are written to the trace file: ack number, tcp-specific flags, and header length.

## 33.8 Commands at a glance

The following is a list of commands used to setup/manipulate TCP flows for simulations:

```
set tcp0 [new Agent/TCP]
```

This creates an instance of a TCP agent. There are several flavors of TCP-sender and TCP-receiver (or sink) agent currently implemented in ns. TCP-senders currently available are: `Agent/TCP`, `Agent/TCP/Reno`, `Agent/TCP/Newreno`, `Agent/TCP/Sack1`, `Agent/TCP/Vegas`, `Agent/TCP/Fack`.

TCP-receivers currently available are: `Agent/TCPSink`, `Agent/TCPSink/DelAck`, `Agent/TCPSink/Sack1`, `Agent/TCPSink/Sack1/DelAck`.

There is also a two-way implementation of tcp called `Agent/TCP/FullTcp`. For details on the different TCP flavors see earlier sections of this chapter.

Configuration parameters for TCP flows maybe set as follows:

```
$tcp set window_ <wnd-size>
```

For all possible configuration parameters available for TCP see section 33.1.4. The default configuration values can also be found in *ns/tcl/lib/ns-default.tcl*.

Following is an example of a simple TCP connection setup:

```
set tcp [new Agent/TCP]                                ; # create tcp agent
$ns_ attach-agent $node_(s1) $tcp                      ; # bind src to node
$tcp set fid_ 0                                         ; # set flow ID field
set ftp [new Application/FTP]                          ; # create ftp traffic
$ftp attach-agent $tcp                                 ; # bind ftp traffic to tcp agent
set sink [new Agent/TCPSink]                           ; # create tcpsink agent
$ns_ attach-agent $node_(k1) $sink                    ; # bind sink to node
$sink set fid_ 0                                         ; # set flow ID field
$ns_ connect $ftp $sink                                ; # active connection src to sink
$ns_ at $start-time "$ftp start"                      ; # start ftp flow
```

For an example of setting up a full-tcp connection see section 33.3.1.

# Chapter 34

## SCTP Agents

This chapter describes the SCTP agents developed for *ns* by the Protocol Engineering Lab at the University of Delaware. The SCTP agents are all two-way agents, which means they are symmetric in the sense that they represent both a sender and receiver. However, bi-directional data has not yet been implemented. Each instance of an SCTP agent is either a sender or receiver.

The SCTP agents are implemented in files matching the regular expression `~ns/sctp/sctp*.{cc, h}`.

The SCTP agents currently supported are:

- Agent/SCTP - RFC2960 + draft-ietf-tsvwg-sctpimpguide-09.txt + draft-ietf-tsvwg-usctp-01.txt
- Agent/SCTP/HbAfterRto - experimental extension (HEARTBEAT after RTO)
- Agent/SCTP/MultipleFastRtx - experimental extension (UD PEL's Multiple Fast Retransmit algorithm)
- Agent/SCTP/Timestamp - experimental extension (TIMESTAMP chunk)
- Agent/SCTP/MfrHbAfterRto - experimental extension that combines MultipleFastRtx and HbAfterRto
- Agent/SCTP/MfrTimestamp - experimental extension that combines MultipleFastRtx and Timestamp

Section 34.1 provides a simple overview of the base SCTP agent with details of configuration parameters and commands. Section 34.2 describes the SCTP extensions available. The details of the SCTP trace format used in packet trace files are explained in Section 34.3. Section 34.4 explains how to use legacy applications with SCTP and how to write SCTP-aware applications which exploit all SCTP's features. Section 34.5 provides example scripts for both singled homed and multihomed endpoints.

### 34.1 The Base SCTP Agent

The base SCTP agent `Agent /SCTP` supports the features in the following sections of RFC2960, including modifications up to draft-ietf-tsvwg-sctpimpguide-13.txt.

5.1 Normal Establishment of an Association (rudimentary handshake)

6.1 Transmission of DATA Chunks



- 6.2 Acknowledgment on Reception of DATA Chunks
- 6.3 Management Retransmission Timer
- 6.4 Multihomed SCTP Endpoints
- 6.5 Stream Identifier and Stream Sequence Number
- 6.6 Ordered and Unordered Delivery
- 6.7 Report Gaps in Received DATA TSNs
- 7.2 SCTP Slow-Start and Congestion Avoidance
- 8.1 Endpoint Failure Detection
- 8.2 Path Failure Detection
- 8.3 Path Heartbeat (without upper layer control)

This agent also supports the Partial Reliability extension as of draft-ietf-tsvwg-usc-01.txt.

**Association Establishment** The SCTP agent establishes an association using a four-way handshake, but the handshake is kept simple and does not strictly conform to RFC2960. The handshake does not exchange tags, and the INIT and COOKIE-ECHO chunks are not used to update the RTT. Instead, RTT estimation begins with the first DATA chunk.

**Association Shutdown** Currently, the SCTP agent does not perform a proper shutdown. The association is abruptly terminated when the simulated connection ends. A shutdown procedure may be added in a future release.

**Multihoming** The underlying infrastructure of ns-2 does not support multiple interfaces for a single node. To get around this limitation, our approach allows the general support for logically multihoming nodes that have a multihomed transport layer, such as SCTP. Each multihomed node is actually made up of more than one node. As shown in Figure 34.1, a logically multihomed node is made up of a single "core node" and multiple "interface nodes", one for each simulated interface. The core node is connected to each interface node via a uni-directional link towards the interface node, but traffic never traverses these links. These links are only in place for the core node to make routing decisions. An SCTP agent simultaneously resides on all these nodes (i.e., the core and interface nodes), but actual traffic only goes to/from the interface nodes. Whenever the SCTP agent needs to send data to a destination and does not know which outgoing interface to use, the agent first consults with the core node for a route lookup. Then, the SCTP agent performs the send from the appropriate interface node. Incoming data is received at one of the interface nodes directly and passed up to the SCTP agent. This solution is applicable to any transport protocol that requires multihoming functionality in ns-2. Note: the user must configure multihomed nodes using commands in Section 34.1.2 (an example is shown in Section 34.5.2).

**Packet Number vs TSN Numbering** While *ns* starts numbering packets at 0, the SCTP module starts numbering DATA chunk TSNs at 1 and assigns undefined TSN values (-1) to control chunks (ie, INIT, SACK, HEARTBEAT, etc). The four packets exchanged during the association establishment are counted in the packet enumeration, but do not show up in graphs. This information is important when doing things like specifying a drop list for the `ErrorModel` object. For example, packet 2 actually refers to the first SCTP packet with DATA chunk(s).

### 34.1.1 Configuration Parameters

SCTP supports several configuration variables which are TCL bindable. Each of the variables described in this subsection is both a class variable and an instance variable. Changing the class variable changes the default value for all agents that are

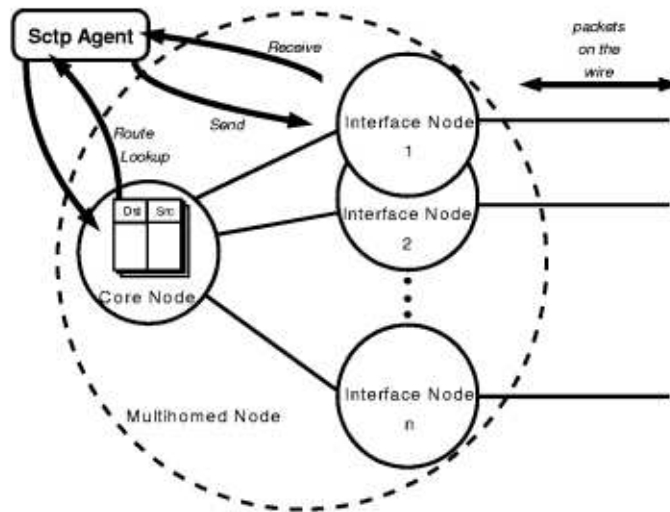


Figure 34.1: Example of a Multihomed Node

created subsequently. Changing the instance variable of a particular agent only affects the values used by that agent. For example,

```
Agent/SCTP set pathMaxRetrans_ 5 ; # Changes the class variable
$scp set pathMaxRetrans_ 5 ; # Changes pathMaxRetrans_for the $scp object only
```

The default parameters for each SCTP agent are:

```
Agent/SCTP set debugMask_ 0 ; # 32-bit mask for modular toggle debugging control (see explanation)
Agent/SCTP set debugFileIndex_ -1 ; # specifies debugging output file (see explanation)
Agent/SCTP set associationMaxRetrans_ 10 ; # RFC2960's Association.Max.Retrans
Agent/SCTP set pathMaxRetrans_ 5 ; # RFC2960's Path.Max.Retrans
Agent/SCTP set changePrimaryThresh_ -1 ; # change primary if error count exceeds thresh (default infinite)
Agent/SCTP set maxInitRetransmits_ 8 ; # RFC2960's Max.Init.Retransmits
Agent/SCTP set oneHeartbeatTimer_ 1 ; # toggle HB timer for each dest vs one for all dests
Agent/SCTP set heartbeatInterval_ 30 ; # RFC2960's HB.interval in seconds
Agent/SCTP set mtu_ 1500 ; # MTU in bytes including IP header
Agent/SCTP set initialRwnd_ 65536 ; # initial receiver window in bytes (set on receiver side)
Agent/SCTP set initialSsthresh_ 65536 ; # initial ssthresh value in bytes
Agent/SCTP set initialCwnd_ 2 ; # initial cwnd in multiple of (MTU - SCTP/IP headers)
Agent/SCTP set initialRto_ 3.0 ; # default initial RTO = 3 secs
Agent/SCTP set minRto_ 1.0 ; # default min RTO = 1 sec
Agent/SCTP set maxRto_ 60.0 ; # default max RTO = 60 secs
Agent/SCTP set fastRtxTrigger_ 4 ; # 4 missing reports trigger fast rtx
Agent/SCTP set numOutStreams_ 1 ; # number of outgoing streams
Agent/SCTP set numUnrelStreams_ 0 ; # number of partially reliable streams (all grouped starting at stream 0)
Agent/SCTP set reliability_ 0 ; # k-rtx value of all partially reliable streams
Agent/SCTP set unordered_ 0 ; # toggle all chunks are ordered/unordered
Agent/SCTP set ipHeaderSize_ 20 ; # IP header size
Agent/SCTP set dataChunkSize_ 1468 ; # includes data chunk header and restricted to 4 byte boundaries
```

```

Agent/SCTP set useDelayedSacks_ 1           ; # toggle on/off delayed sack algorithm (set on receiver side)
Agent/SCTP set sackDelay_ 0.200             ; # rfc2960 recommends 200 ms
Agent/SCTP set useMaxBurst_ 1                ; # toggle on/off max burst
Agent/SCTP set rtxToAlt_ 1                  ; # rtxs to which dest? 0 = same, 1 = alt, 2 = fast rtx to same + timeouts to alt
Agent/SCTP set dormantAction_ 0             ; # 0 = change dest, 1 = use primary, 2 = use last dest before dormant
Agent/SCTP set routeCalcDelay_ 0            ; # time to calculate a route (see explanation)
Agent/SCTP set routeCacheLifetime_ 1.2      ; # how long a route remains cached (see explanation)
Agent/SCTP set trace_all_ 0                 ; # toggle on/off print all variables on a trace event

```

The `debugMask_` parameter is a 32-bit mask to turn on/off debugging for particular functions. See `~ns/sctp/sctpDebug.h` for the mappings of the bitmask. A -1 may be used to clear all bits, and 0 is used to turn off all debugging. If `debug_` (the standard *ns* debug flag) is set to 1, then all the bits in `debugMask_` are set. Note: *ns* must be compiled with `-DDEBUG` for this option to work.

The `debugFileIndex_` parameter is an integer to specify the file used for debugging output by an SCTP agent. Each instance of an SCTP agent can independently output debugging info to a separate file. For example, the data sender can log debugging output to one file, while the receiver logs to another file. If `debugFileIndex_` is set to 0, the file used will be named `debug.SctpAgent.0`. If -1 is used, the debug output is sent to `stderr`. To avoid confusion, two SCTP agents should not send debug output to the same file. The default is -1. Note: *ns* must be compiled with `-DDEBUG` for this option to work.

The configuration parameters that deal with ordering and reliability options may be overridden by an SCTP-aware application (see Section 34.4).

The `routeCalcDelay_` and `routeCacheLifetime_` parameters are only used to optionally simulate overheads of reactive routing protocols in MANETs without actually simulating a MANET. (Do not use this feature if you are actually simulating a MANET!) The default setting for `routeCalcDelay_` is 0 seconds, which means that this feature is turned off. The default setting for `routeCacheLifetime_` is 1.2 seconds (ignored if this feature is turned off), which is purposely set slightly larger than the default min RTO to avoid a “cache miss” after a single timeout event.

### 34.1.2 Commands

SCTP provides certain commands that can be used within TCL scripts:

**trace** Tracks given variables. The variable (and associated information) is printed every time the value changes. Takes 1 argument:

`cwnd_` Used to trace the cwnds of all paths.

`rto_` Used to trace the RTOs of all paths.

`errorCount_` Used to trace the error counters of all paths.

`frCount_` Used to trace the number of times a fast retransmit is invoked.

`mfrCount_` Used to trace the number of times the multiple fast retransmit algorithm is invoked. This trace variable can only be used with the MultipleFastRtx extension agent. (See Section 34.2.2)

`timeoutCount_` Used to trace the total number of times a timeout has occurred on all paths.

`rcdCount_` Used to trace the total number of times a route calculation delay (see Section 34.1.1) has occurred on all paths.

Note: the actual value of these trace variables have no meaning. They are simply used to trace corresponding variables for potentially multihomed endpoints. For example, if a sender's peer endpoint has two destinations, the sender will maintain two cwnds. The `cwnd_` trace variable will trace both of these cwnds together.

**print** Provides the sampling method of tracing. This command simply prints a given variable (and associated information) per call. Takes 1 argument: one of the trace variables presented above.

**set-multihome-core** Sets the core node for multihomed endpoints. Takes 1 argument of type node. Mandatory for multihomed endpoints and must not be set more than once per endpoint.

**multihome-add-interface** Adds an interface to a multihomed endpoint. Takes 2 arguments of type node. Argument 1 is the core node of the multihomed endpoint. Argument 2 is the interface node to be added. Mandatory for multihomed endpoints. All interfaces must be added after `set-multihome-core` is called and before `multihome-attach-agent` is called.

**multihome-attach-agent** Attaches an SCTP agent to a multihomed endpoint. Takes 2 arguments. Argument 1 is the core node. Argument 2 is the SCTP agent. Mandatory for multihomed endpoints.

**set-primary-destination** Sets the interface node of the peer endpoint as the primary destination. Takes 1 argument of type node. Optional and may be set more than once per endpoint. If not used, a primary destination is chosen automatically.

**force-source** Sets the interface node that packets will be sent from. Takes 1 argument of type node. Optional and may be set more than once per endpoint. If not used, routing will automatically choose the source on a per packet basis.

## 34.2 Extensions

### 34.2.1 HbAfterRto SCTP

The HbAfterRto SCTP agent extends the current retransmission policy. In addition to SCTP's current policy of retransmitting to an alternate destination on a timeout, a heartbeat is sent immediately to the destination on which a timeout occurred. Extra heartbeats provide a mechanism for a sender to update an alternate destination's RTT estimate more frequently, thus resulting in a better RTT estimate on which to base the RTO value.

For example, suppose a packet is lost in transit to the primary destination, and later gets retransmitted to an alternate destination. Also suppose that the retransmission times out. The lost packet is retransmitted again to yet another alternate destination (if one exists; otherwise, the primary). More importantly, a heartbeat is also sent to the alternate destination which timed out. If the heartbeat is successfully acked, that destination acquires an additional RTT measurement to help reduce its recently doubled RTO [?].

### 34.2.2 MultipleFastRtx SCTP

The MultipleFastRtx SCTP agent attempts to minimize the number of timeouts which occur. Without the Multiple Fast Retransmit algorithm, SCTP may only Fast Retransmit a TSN once. If a Fast Retransmitted TSN is lost, a timeout is necessary

to retransmit the TSN again. The Multiple Fast Retransmit algorithm allows the same TSN to be Fast Retransmitted several times if needed. Without the Multiple Fast Retransmit algorithm, a large window of outstanding data may generate enough SACKs to incorrectly trigger more than one Fast Retransmit of the same TSN in a single RTT. To avoid these spurious Fast Retransmits, the Multiple Fast Retransmit algorithm introduces a *fastRtxRecover* state variable for each TSN Fast Retransmitted. This variable stores the highest outstanding TSN at the time a TSN is Fast Retransmitted. Then, only SACKs which newly ack TSNs beyond *fastRtxRecover* can increment the missing report for the Fast Retransmitted TSN. If the missing report threshold for the Fast Retransmitted TSN is reached again, the sender has enough evidence that this TSN was lost and can be Fast Retransmitted again [?].

### 34.2.3 Timestamp SCTP

The Timestamp SCTP agent introduces timestamps into each packet, thus allowing a sender to disambiguate original transmissions from retransmissions. By removing retransmission ambiguity, Karn's algorithm can be eliminated, and successful retransmissions on the alternate path can be used to update the RTT estimate and keep the RTO value more accurate. With timestamps, the sender has more samples for updating the RTT estimate of alternate destination(s) [?].

### 34.2.4 MfrHbAfterRto SCTP

The MfrHbAfterRto SCTP agent combines the functionality of the HbAfterRto and MultipleFastRtx SCTP agents.

### 34.2.5 MfrHbAfterRto SCTP

The MfrTimestamp SCTP agent combines the functionality of the Timestamp and MultipleFastRtx SCTP agents.

## 34.3 Tracing SCTP Dynamics

SCTP packets generated by one of the SCTP agents and destined for a peer SCTP agent over a traced link (see section 25) will generate a trace file with lines of the form:

```
+ 0.5 1 4 sctp 56 -----I 0 1.0 4.0 1 -1 4 65535 65535
- 0.5 1 4 sctp 56 -----I 0 1.0 4.0 1 -1 4 65535 65535
r 0.700896 1 4 sctp 56 -----I 0 1.0 4.0 1 -1 4 65535 65535
+ 0.700896 4 1 sctp 56 -----I 0 4.0 1.0 1 -1 5 65535 65535
- 0.700896 4 1 sctp 56 -----I 0 4.0 1.0 1 -1 5 65535 65535
r 0.901792 4 1 sctp 56 -----I 0 4.0 1.0 1 -1 5 65535 65535
+ 0.901792 1 4 sctp 36 -----I 0 1.0 4.0 1 -1 6 65535 65535
- 0.901792 1 4 sctp 36 -----I 0 1.0 4.0 1 -1 6 65535 65535
r 1.102368 1 4 sctp 36 -----I 0 1.0 4.0 1 -1 6 65535 65535
+ 1.102368 4 1 sctp 36 -----I 0 4.0 1.0 1 -1 7 65535 65535
- 1.102368 4 1 sctp 36 -----I 0 4.0 1.0 1 -1 7 65535 65535
r 1.302944 4 1 sctp 36 -----I 0 4.0 1.0 1 -1 7 65535 65535
+ 1.302944 1 4 sctp 1500 -----D 0 1.0 4.0 1 1 8 0 0
- 1.302944 1 4 sctp 1500 -----D 0 1.0 4.0 1 1 8 0 0
+ 1.302944 1 4 sctp 1500 -----D 0 1.0 4.0 1 2 9 0 1
- 1.326624 1 4 sctp 1500 -----D 0 1.0 4.0 1 2 9 0 1
```

```

r 1.526624 1 4 sctp 1500 -----D 0 1.0 4.0 1 1 8 0 0
r 1.550304 1 4 sctp 1500 -----D 0 1.0 4.0 1 2 9 0 1
+ 1.550304 4 1 sctp 48 -----S 0 4.0 1.0 1 2 11 65535 65535
- 1.550304 4 1 sctp 48 -----S 0 4.0 1.0 1 2 11 65535 65535
r 1.751072 4 1 sctp 48 -----S 0 4.0 1.0 1 2 11 65535 65535
...
+ 19.302944 4 1 sctp 56 -----H 0 2.0 5.0 1 -1 336 65535 65535
- 19.302944 4 1 sctp 56 -----H 0 2.0 5.0 1 -1 336 65535 65535
r 19.303264 4 1 sctp 56 -----H 0 4.0 1.0 1 -1 322 65535 65535
+ 19.303264 1 4 sctp 56 -----B 0 1.0 4.0 1 -1 337 65535 65535
- 19.327584 1 4 sctp 56 -----B 0 1.0 4.0 1 -1 337 65535 65535
r 19.52848 1 4 sctp 56 -----B 0 1.0 4.0 1 -1 337 65535 65535

```

When tracing SCTP, packets of type `sctp` are relevant. Their packet type, chunk type, packet size, TSN (CumAck point for SACK chunks), stream ID, SSN, and arrival/depart/drop time are given by field positions 5, 7 (flag position 8), 6, 12, 14, 15, and 2, respectively. If a packet has more than one chunk, a line is printed for each chunk. A future release should include a field to indicate which chunk of the packet a line refers to (e.g., 2:3 could identify the 2nd chunk of a packet which contains 3 chunks). Since control chunks do not use TSNs, stream IDs, or SSNs, the trace lines for these chunks use undefined numbers (-1 or 65535) for these fields. The `+` indicates a packet arrival, `d` a drop, and `-` a departure.

Flag position 8 of field 7 indicate the chunk type as follows. The `I` flag indicates an association initiation control chunk (INIT, INIT-ACK, COOKIE-ECHO, COOKIE-ACK). A future release should use `I` for the INIT and INIT-ACK chunks, and `C` for the COOKIE-ECHO and COOKIE-ACK chunks. The `D`, `S`, `H`, and `B` flags indicate a DATA chunk, a SACK, a HEARTBEAT chunk, and a HEARTBEAT-ACK chunk.

A number of scripts process this file to produce graphical output or statistical summaries (for example, see the `finish` procedure in `~ns/tcl/test/test-suite-sctp.tcl`).

## 34.4 SCTP Applications

SCTP supports legacy *ns* applications, but they obviously cannot completely exploit all SCTP's features. For these applications, the TCL-bound SCTP configuration parameters can be used to set reliability and ordering options. However, these options cannot be controlled per message using these parameters. Only SCTP-aware application can be written to do so. *ns* applications wishing to become SCTP-aware can use the `sendmsg` API as follows (see `~ns/apps/sctp_app1.{cc, h}` as an example).

1. Create and fill an instance of the `AppData_S` structure (see `~ns/sctp/sctp.h`). The `AppData_S` structure has the following fields:

`usNumStreams` is the number of outgoing streams to setup during negotiation. Although this field is passed with every `sendmsg` call, it is only used during association setup. Once the association is established, this field is ignored.

`usNumUnreliable` is the number of outgoing streams which are unreliable (now called partially reliable). The sender simply sets the lowest outgoing stream to unreliable/partially-reliable; the remaining ones are reliable. This field is also only used during association establishment.

`usStreamId` is the stream ID of a message.

`usReliability` is the reliability level (k-rtx value) of a message. This field is ignored if the message is sent on a reliable stream.

eUnordered is the unordered boolean flag for a message.

uiNumBytes is the number of bytes in a message.

2. Pass this object as the second parameter in SCTP's sendmsg:

```
sctpAgent->sendmsg(numBytes, (char *)appData);
```

## 34.5 Example Scripts

### 34.5.1 Singled Homed Example

```
Trace set show_sctphdr_ 1                                ; # this needs to be set for tracing SCTP packets

set ns [new Simulator]
set allchan [open all.tr w]
$ns trace-all $allchan

proc finish {
    exit 0
}

set n0 [$ns node]
set n1 [$ns node]
$ns duplex-link $n0 $n1 .5Mb 200ms DropTail

# NOTE: The debug files (in this example, they would be debug.SctpAgent.0
#       and debug.SctpAgent.1) contain a lot of useful info. They can be
#       used to trace every packet sent, received, and processed.
#
set sctp0 [new Agent/SCTP]
$ns attach-agent $n0 $sctp0
$sctp0 set debugMask_ 0x00303000                        ; # refer to sctpDebug.h for mask mappings
$sctp0 set debugFileIndex_ 0

set trace_ch [open trace.sctp w]
$sctp0 set trace_all_ 0                                  ; # do not trace all variables on one line
$sctp0 trace cwnd_                                       ; # trace cwnd for all destinations
$sctp0 attach $trace_ch

set sctp1 [new Agent/SCTP]
$ns attach-agent $n1 $sctp1
$sctp1 set debugMask_ -1                                ; # use -1 to turn on all debugging
$sctp1 set debugFileIndex_ 1

$ns connect $sctp0 $sctp1

set ftp0 [new Application/FTP]
$ftp0 attach-agent $sctp0

$ns at 0.5 "$ftp0 start"
$ns at 4.5 "$ftp0 stop"
```

```
$ns at 5.0 "finish"
```

```
$ns run
```

### 34.5.2 Multihomed Example

```
# This example demonstrates multihoming. Two SCTP endpoints, each  
# with 2 interfaces, are directly connected between each pair of  
# interfaces. In the middle of the association, a change primary  
# is done. Running nam helps to visualize the multihomed  
# architecture.
```

```
#  
#      host0_if0      O=====O      host1_if0  
#      /              \  
#      host0_core    O              O      host1_core  
#      \  
#      host0_if1      O=====O      host1_if1
```

```
Trace set show_sctphdr_ 1
```

```
set ns [new Simulator]  
set nf [open sctp.nam w]  
$ns namtrace-all $nf  
set allchan [open all.tr w]  
$ns trace-all $allchan
```

```
proc finish {  
    exec nam sctp.nam &  
    exit 0  
}
```

```
set host0_core [$ns node]  
set host0_if0 [$ns node]  
set host0_if1 [$ns node]  
$host0_core color Red  
$host0_if0 color Red  
$host0_if1 color Red  
$ns multihome-add-interface $host0_core $host0_if0  
$ns multihome-add-interface $host0_core $host0_if1
```

```
set host1_core [$ns node]  
set host1_if0 [$ns node]  
set host1_if1 [$ns node]  
$host1_core color Blue  
$host1_if0 color Blue  
$host1_if1 color Blue  
$ns multihome-add-interface $host1_core $host1_if0  
$ns multihome-add-interface $host1_core $host1_if1
```

```
$ns duplex-link $host0_if0 $host1_if0 .5Mb 200ms DropTail
```



```

$ns duplex-link $host0_if1 $host1_if1 .5Mb 200ms DropTail

set sctp0 [new Agent/SCTP]
$ns multihome-attach-agent $host0_core $sctp0

set trace_ch [open trace.sctp w]
$sctp0 set trace_all_ 1 ; # print all on a single trace event
$sctp0 trace rto_
$sctp0 trace errorCount_
$sctp0 attach $trace_ch

set sctp1 [new Agent/SCTP]
$ns multihome-attach-agent $host1_core $sctp1

$ns connect $sctp0 $sctp1

set ftp0 [new Application/FTP]
$ftp0 attach-agent $sctp0

$sctp0 set-primary-destination $host1_if0 ; # set primary before association starts

$ns at 7.5 "$sctp0 set-primary-destination $host1_if1" ; # change primary
$ns at 7.5 "$sctp0 print cwnd_" ; # print all dests' cwnds at time 7.5

$ns at 0.5 "$ftp0 start"
$ns at 9.5 "$ftp0 stop"
$ns at 10.0 "finish"

$ns run

```

# Chapter 35

## Agent/SRM

This chapter describes the internals of the SRM implementation in *ns*. The chapter is in three parts: the first part is an overview of a minimal SRM configuration, and a “complete” description of the configuration parameters of the base SRM agent. The second part describes the architecture, internals, and the code path of the base SRM agent. The last part of the chapter is a description of the extensions for other types of SRM agents that have been attempted to date.

The procedures and functions described in this chapter can be found in *~ns/tcl/mcast/srm.tcl*, *~ns/tcl/mcast/srm-adaptive.tcl*, *~ns/tcl/mcast/srm-nam.tcl*, *~ns/tcl/mcast/srm-debug.tcl*, and *~ns/srm.{cc, h}*.

### 35.1 Configuration

Running an SRM simulation requires creating and configuring the agent, attaching an application-level data source (a traffic generator), and starting the agent and the traffic generator.

#### 35.1.1 Trivial Configuration

##### Creating the Agent

```
set ns [new Simulator]                                ; # preamble initialization
$ns enableMcast
set node [$ns node]                                   ; # agent to reside on this node
set group [$ns allocaddr]                             ; # multicast group for this agent

set srm [new Agent/SRM]
$srmm set dst_ $group                                ; # configure the SRM agent
$ns attach-agent $node $srm

$srmm set fid_ 1                                       ; # optional configuration
$srmm log [open srmStats.tr w]                         ; # log statistics in this file
$srmm trace [open srmEvents.tr w]                     ; # trace events for this agent
```

The key steps in configuring a virgin SRM agent are to assign its multicast group, and attach it to a node.

Other useful configuration parameters are to assign a separate flow id to traffic originating from this agent, to open a log file for statistics, and a trace file for trace data<sup>1</sup>.

The file `tcl/mcast/srm-nam.tcl` contains definitions that overload the agent's send methods; this separates control traffic originating from the agent by type. Each type is allocated a separate flowID. The traffic is separated into session messages (flowid = 40), requests (flowid = 41), and repair messages (flowid = 42). The base flowid can be changed by setting global variable `ctrlFid` to one less than the desired flowid before sourcing `srm-nam.tcl`. To do this, the simulation script must source `srm-nam.tcl` before creating any SRM agents. This is useful for analysis of traffic traces, or for visualization in nam.

**Application Data Handling** The agent does not generate any application data on its own; instead, the simulation user can connect any traffic generation module to any SRM agent to generate data. The following code demonstrates how a traffic generation agent can be attached to an SRM agent:

```
set packetSize 210
set exp0 [new Application/Traffic/Exponential]           ; # configure traffic generator
$exp0 set packetSize_ $packetSize
$exp0 set burst_time_ 500ms
$exp0 set idle_time_ 500ms
$exp0 set rate_ 100k

$exp0 attach-agent $srm0                                ; # attach application to SRM agent
$srm0 set packetSize_ $packetSize                        ; # to generate repair packets of appropriate size
$srm0 set tg_ $exp0                                       ; # pointer to traffic generator object
$srm0 set app_fid_ 0                                       ; # fid value for packets generated by traffic generator
```

The user can attach any traffic generator to an SRM agent. The SRM agent will add the SRM headers, set the destination address to the multicast group, and deliver the packet to its target. The SRM header contains the type of the message, the identity of the sender, the sequence number of the message, and (for control messages), the round for which this message is being sent. Each data unit in SRM is identified as (sender's id, message sequence number).

The SRM agent does not generate its own data; it does not also keep track of the data sent, except to record the sequence numbers of messages received in the event that it has to do error recovery. Since the agent has no actual record of past data, it needs to know what packet size to use for each repair message. Hence, the instance variable `packetSize_` specifies the size of repair messages generated by the agent.

**Starting the Agent and Traffic Generator** The agent and the traffic generator must be started separately.

```
$srm start
$exp0 start
```

Alternatively, the traffic generator can be started from the SRM Agent:

```
$srm0 start-source
```

At start, the agent joins the multicast group, and starts generating session messages. The `start-source` triggers the traffic generator to start sending data.

---

<sup>1</sup>Note that the trace data can also be used to gather certain kinds of trace data. We will illustrate this later.

### 35.1.2 Other Configuration Parameters

In addition to the above parameters, the SRM agent supports additional configuration variables. Each of the variables described in this section is both an OTcl class variable and an OTcl object's instance variable. Changing the class variable changes the default value for all agents that are created subsequently. Changing the instance variable of a particular agent only affects the values used by that agent. For example,

```
Agent/SRM set D1_ 2.0 ; # Changes the class variable
$srn set D1_ 2.0 ; # Changes D1_ for the particular $srn object only
```

The default request and repair timer parameters [11] for each SRM agent are:

```
Agent/SRM set C1_ 2.0 ; # request parameters
Agent/SRM set C2_ 2.0
Agent/SRM set D1_ 1.0 ; # repair parameters
Agent/SRM set D2_ 1.0
```

It is thus possible to trivially obtain two flavors of SRM agents based on whether the agents use probabilistic or deterministic suppression by using the following definitions:

```
Class Agent/SRM/Deterministic -superclass Agent/SRM
Agent/SRM/Deterministic set C2_ 0.0
Agent/SRM/Deterministic set D2_ 0.0

Class Agent/SRM/Probabilistic -superclass Agent/SRM
Agent/SRM/Probabilistic set C1_ 0.0
Agent/SRM/Probabilistic set D1_ 0.0
```

In a later section (Section 35.7), we will discuss other ways of extending the SRM agent.

Timer related functions are handled by separate objects belonging to the class SRM. Timers are required for loss recovery and sending periodic session messages. There are loss recovery objects to send request and repair messages. The agent creates a separate request or repair object to handle each loss. In contrast, the agent only creates one session object to send periodic session messages. The default classes that express each of these functions are:

```
Agent/SRM set requestFunction_ "SRM/request"
Agent/SRM set repairFunction_ "SRM/repair"
Agent/SRM set sessionFunction_ "SRM/session"

Agent/SRM set requestBackoffLimit_ 5 ; # parameter to requestFunction_
Agent/SRM set sessionDelay_ 1.0 ; # parameter to sessionFunction_
```

The instance procedures `requestFunction{}`, `repairFunction{}`, and `sessionFunction{}` can be used to change the default function for individual agents. The last two lines are specific parameters used by the request and session objects. The following section (Section 35.2) describes the implementation of these objects in greater detail.

### 35.1.3 Statistics

Each agent tracks two sets of statistics: statistics to measure the response to data loss, and overall statistics for each request/repair. In addition, there are methods to access other information from the agent.

**Data Loss** The statistics to measure the response to data losses tracks the duplicate requests (and repairs), and the average request (and repair) delay. The algorithm used is documented in Floyd *et al*[11]. In this algorithm, each new request (or repair) starts a new request (or repair) period. During the request (or repair) period, the agent measures the number of first round duplicate requests (or repairs) until the round terminates either due to receiving a request (or repair), or due to the agent sending one. The following code illustrates how the user can simple retrieve the current values in an agent:

```
set statsList [$srm array get statistics_]
array set statsArray [$srm array get statistics_]
```

The first form returns a list of key-value pairs. The second form loads the list into the statsArray for further manipulation. The keys of the array are dup-req, ave-dup-req, req-delay, ave-req-delay, dup-rep, ave-dup-rep, rep-delay, and ave-rep-delay.

**Overall Statistics** In addition, each loss recovery and session object keeps track of times and statistics. In particular, each object records its startTime, serviceTime, distance, as are relevant to that object; startTime is the time that this object was created, serviceTime is the time for this object to complete its task, and the distance is the one-way time to reach the remote peer.

For request objects, startTime is the time a packet loss is detected, serviceTime is the time to finally receive that packet, and distance is the distance to the original sender of the packet. For repair objects, startTime is the time that a request for retransmission is received, serviceTime is the time send a repair, and the distance is the distance to the original requester. For both types of objects, the serviceTime is normalized by the distance. For the session object, startTime is the time that the agent joins the multicast group. serviceTime and distance are not relevant.

Each object also maintains statistics particular to that type of object. Request objects track the number of duplicate requests and repairs received, the number of requests sent, and the number of times this object had to backoff before finally receiving the data. Repair objects track the number of duplicate requests and repairs, as well as whether or not this object for this agent sent the repair. Session objects simply record the number of session messages sent.

The values of the timers and the statistics for each object are written to the log file every time an object completes the error recovery function it was tasked to do. The format of this trace file is:

```
<prefix> <id> <times> <stats>
where
<prefix> is      <time> n <node id> m <msg id> r <round>
                  <msg id> is expressed as <source id:sequence number>
<id> is          type <of object>
<times> is       list of key-value pairs of startTime, serviceTime, distance
<stats> is       list of key-value pairs of per object statistics
                  dupRQST, dupREPR, #sent, backoff      for request objects
                  dupRQST, dupREPR, #sent               for repair objects
                  #sent                                  for session objects
```

The following sample output illustrates the output file format (the lines have been folded to fit on the page):

```

3.6274 n 0 m <1:1> r 1 type repair serviceTime 0.500222 \
      startTime 3.5853553333333332 distance 0.0105 #sent 1 dupREPR 0 dupRQST 0
3.6417 n 1 m <1:1> r 2 type request serviceTime 2.66406 \
      startTime 3.5542666666666665 distance 0.0105 backoff 1 #sent 1 dupREPR 0 dupRQST 0
3.6876 n 2 m <1:1> r 2 type request serviceTime 1.33406 \
      startTime 3.5685333333333333 distance 0.021 backoff 1 #sent 0 dupREPR 0 dupRQST 0
3.7349 n 3 m <1:1> r 2 type request serviceTime 0.876812 \
      startTime 3.5828000000000002 distance 0.032 backoff 1 #sent 0 dupREPR 0 dupRQST 0
3.7793 n 5 m <1:1> r 2 type request serviceTime 0.669063 \
      startTime 3.5970666666666671 distance 0.042 backoff 1 #sent 0 dupREPR 0 dupRQST 0
3.7808 n 4 m <1:1> r 2 type request serviceTime 0.661192 \
      startTime 3.5970666666666671 distance 0.0425 backoff 1 #sent 0 dupREPR 0 dupRQST 0

```

**Miscellaneous Information** Finally, the user can use the following methods to gather additional information about the agent:

- `groupSize?{}` returns the agent's current estimate of the multicast group size.
- `distances?{}` returns a list of key-value pairs of distances; the key is the address of the agent, the value is the estimate of the distance to that agent. The first element is the address of this agent, and the distance of 0.
- `distance?{}` returns the distance to the particular agent specified as argument.

The default distance at the start of any simulation is 1.

<code>\$srm(i) groupSize?</code>	<i>; # returns \$srm(i)'s estimate of the group size</i>
<code>\$srm(i) distances?</code>	<i>; # returns list of &lt;address, distance&gt; tuples</i>
<code>\$srm(i) distance? 257</code>	<i>; # returns the distance to agent at address 257</i>

### 35.1.4 Tracing

Each object writes out trace information that can be used to track the progress of the object in its error recovery. Each trace entry is of the form:

`<prefix> <tag> <type of entry> <values>`

The prefix is as described in the previous section for statistics. The tag is **Q** for request objects, **P** for repair objects, and **S** for session objects. The following types of trace entries and parameters are written by each object:

Tag	Type of Object	Other values	Comments
Q	DETECT		
Q	INTERVALS	C1 <C1_> C2 <C2_> dist <distance> i <backoff_>	
Q	NTIMER	at <time>	Time the request timer will fire
Q	SENDNACK		
Q	NACK	IGNORE-BACKOFF <time>	Receive NACK, ignore other NACKs until <time>
Q	REPAIR	IGNORES <time>	Receive REPAIR, ignore NACKs until <time>
Q	DATA		Agent receives data instead of repair. Possibly indicates out of order arrival of data.
P	NACK	from <requester>	Receive NACK, initiate repair
P	INTERVALS	D1 <D1_> D2 <D2_> dist <distance>	
P	RTIMER	at <time>	Time the repair timer will fire
P	SENDREP		
P	REPAIR	IGNORES <time>	Receive REPAIR, ignore NACKs until <time>
P	DATA		Agent receives data instead of repair. Indicates premature request by an agent.
S	SESSION		logs session message sent

The following illustrates a typical trace for a single loss and recovery.

```

3.5543 n 1 m <1:1> r 0 Q DETECT
3.5543 n 1 m <1:1> r 1 Q INTERVALS C1 2.0 C2 0.0 d 0.0105 i 1
3.5543 n 1 m <1:1> r 1 Q NTIMER at 3.57527
3.5685 n 2 m <1:1> r 0 Q DETECT
3.5685 n 2 m <1:1> r 1 Q INTERVALS C1 2.0 C2 0.0 d 0.021 i 1
3.5685 n 2 m <1:1> r 1 Q NTIMER at 3.61053
3.5753 n 1 m <1:1> r 1 Q SENDNACK
3.5753 n 1 m <1:1> r 2 Q INTERVALS C1 2.0 C2 0.0 d 0.0105 i 2
3.5753 n 1 m <1:1> r 2 Q NTIMER at 3.61727
3.5753 n 1 m <1:1> r 2 Q NACK IGNORE-BACKOFF 3.59627
3.5828 n 3 m <1:1> r 0 Q DETECT
3.5828 n 3 m <1:1> r 1 Q INTERVALS C1 2.0 C2 0.0 d 0.032 i 1
3.5828 n 3 m <1:1> r 1 Q NTIMER at 3.6468
3.5854 n 0 m <1:1> r 0 P NACK from 257
3.5854 n 0 m <1:1> r 1 P INTERVALS D1 1.0 D2 0.0 d 0.0105
3.5854 n 0 m <1:1> r 1 P RTIMER at 3.59586
3.5886 n 2 m <1:1> r 2 Q INTERVALS C1 2.0 C2 0.0 d 0.021 i 2
3.5886 n 2 m <1:1> r 2 Q NTIMER at 3.67262
3.5886 n 2 m <1:1> r 2 Q NACK IGNORE-BACKOFF 3.63062
3.5959 n 0 m <1:1> r 1 P SENDREP
3.5959 n 0 m <1:1> r 1 P REPAIR IGNORES 3.62736
3.5971 n 4 m <1:1> r 0 Q DETECT
3.5971 n 4 m <1:1> r 1 Q INTERVALS C1 2.0 C2 0.0 d 0.0425 i 1
3.5971 n 4 m <1:1> r 1 Q NTIMER at 3.68207
3.5971 n 5 m <1:1> r 0 Q DETECT
3.5971 n 5 m <1:1> r 1 Q INTERVALS C1 2.0 C2 0.0 d 0.042 i 1
3.5971 n 5 m <1:1> r 1 Q NTIMER at 3.68107
3.6029 n 3 m <1:1> r 2 Q INTERVALS C1 2.0 C2 0.0 d 0.032 i 2

```

```

3.6029 n 3 m <1:1> r 2 Q NTIMER at 3.73089
3.6029 n 3 m <1:1> r 2 Q NACK IGNORE-BACKOFF 3.66689
3.6102 n 1 m <1:1> r 2 Q REPAIR IGNORES 3.64171
3.6172 n 4 m <1:1> r 2 Q INTERVALS C1 2.0 C2 0.0 d 0.0425 i 2
3.6172 n 4 m <1:1> r 2 Q NTIMER at 3.78715
3.6172 n 4 m <1:1> r 2 Q NACK IGNORE-BACKOFF 3.70215
3.6172 n 5 m <1:1> r 2 Q INTERVALS C1 2.0 C2 0.0 d 0.042 i 2
3.6172 n 5 m <1:1> r 2 Q NTIMER at 3.78515
3.6172 n 5 m <1:1> r 2 Q NACK IGNORE-BACKOFF 3.70115
3.6246 n 2 m <1:1> r 2 Q REPAIR IGNORES 3.68756
3.6389 n 3 m <1:1> r 2 Q REPAIR IGNORES 3.73492
3.6533 n 4 m <1:1> r 2 Q REPAIR IGNORES 3.78077
3.6533 n 5 m <1:1> r 2 Q REPAIR IGNORES 3.77927

```

The logging of request and repair traces is done by `SRM::evTrace{}`. However, the routine `SRM/Session::evTrace{}`, overrides the base class definition of `srn::evTrace{}`, and writes out nothing. Individual simulation scripts can override these methods for greater flexibility in logging options. One possible reason to override these methods might to reduce the amount of data generated; the new procedure could then generate compressed and processed output.

Notice that the trace file contains sufficient information and details to derive most of the statistics written out in the log file, or is stored in the statistics arrays.

## 35.2 Architecture and Internals

The SRM agent implementation splits the protocol functions into packet handling, loss recovery, and session message activity.

Packet handling consists of forwarding application data messages, sending and receipt of control messages. These activities are executed by C++ methods.

Error detection is done in C++ due to receipt of messages. However, the loss recovery is entirely done through instance procedures in OTcl.

The sending and processing of messages is accomplished in C++; the policy about when these messages should be sent is decided by instance procedures in OTcl.

We first describe the C++ processing due to receipt of messages (Section 35.3). Loss recovery and the sending of session messages involves timer based processing. The agent uses a separate `class SRM` to perform the timer based functions. For each loss, an agent may do either request or repair processing. Each agent will instantiate a separate loss recovery object for every loss, as is appropriate for the processing that it has to do. In the following section we describe the basic timer based functions and the loss recovery mechanisms (Section 35.5). Finally, each agent uses one timer based function for sending periodic session messages (Section 35.6).

## 35.3 Packet Handling: Processing received messages

The `recv()` method can receive four type of messages: data, request, repair, and session messages.

**Data Packets** The agent does not generate any data messages. The user has to specify an external agent to generate traffic. The `recv()` method must distinguish between locally originated data that must be sent to the multicast group, and data



received from multicast group that must be processed. Therefore, the application agent must set the packet's destination address to zero.

For locally originated data, the agent adds the appropriate SRM headers, sets the destination address to the multicast group, and forwards the packet to its target.

On receiving a data message from the group, `recv_data(sender, msgid)` will update its state marking message `(sender, msgid)` received, and possibly trigger requests if it detects losses. In addition, if the message was an older message received out of order, then there must be a pending request or repair that must be cleared. In that case, the compiled object invokes the OTcl instance procedure, `recv_data{sender, msgid}`<sup>2</sup>.

Currently, there is no provision for the receivers to actually receive any application data. The agent does not also store any of the user data. It only generates repair messages of the appropriate size, defined by the instance variable `packetSize_`. However, the agent assumes that any application data is placed in the data portion of the packet, pointed to by `packet->accessdata()`.

**Request Packets** On receiving a request, `recv_rqst(sender, msgid)` will check whether it needs to schedule requests for other missing data. If it has received this request before it was aware that the source had generated this data message (*i.e.*, the sequence number of the request is higher than the last known sequence number of data from this source), then the agent can infer that it is missing this, as well as data from the last known sequence number onwards; it schedules requests for all of the missing data and returns. On the other hand, if the sequence number of the request is less than the last known sequence number from the source, then the agent can be in one of three states: (1) it does not have this data, and has a request pending for it, (2) it has the data, and has seen an earlier request, upon which it has a repair pending for it, or (3) it has the data, and it should instantiate a repair. All of these error recovery mechanisms are done in OTcl; `recv_rqst()` invokes the instance procedure `recv_rqst{sender, msgid, requester}` for further processing.

**Repair Packets** On receiving a repair, `recv_repr(sender, msgid)` will check whether it needs to schedule requests for other missing data. If it has received this repair before it was aware that the source had generated this data message (*i.e.*, the sequence number of the repair is higher than the last known sequence number of data from this source), then the agent can infer that it is missing all data between the last known sequence number and that on the repair; it schedules requests for all of this data, marks this message as received, and returns. On the other hand, if the sequence number of the request is less than the last known sequence number from the source, then the agent can be in one of three states: (1) it does not have this data, and has a request pending for it, (2) it has the data, and has seen an earlier request, upon which it has a repair pending for it, or (3) it has the data, and probably scheduled a repair for it at some time; after error recovery, its hold down timer (equal to three times its distance to some requester) expired, at which time the pending object was cleared. In this last situation, the agent will simply ignore the repair, for lack of being able to do anything meaningful. All of these error recovery mechanisms are done in OTcl; `recv_repr()` invokes the instance procedure `recv_repr{sender, msgid}` to complete the loss recovery phase for the particular message.

**Session Packets** On receiving a session message, the agent updates its sequence numbers for all active sources, and computes its instantaneous distance to the sending agent if possible. The agent will ignore earlier session messages from a group member, if it has received a later one out of order.

Session message processing is done in `recv_sess()`. The format of the session message is: `(count of tuples in this message, list of tuples)`, where each tuple indicates the `(sender id, last sequence number from the source, time the last session message was received from this sender, time that that message was sent)`. The first tuple is the information about the local agent<sup>3</sup>.

---

<sup>2</sup>Technically, `recv_data()` invokes the instance procedure `recv_data (sender) (msgid)`, that then invokes `recv_data{}`. The indirection allows individual simulation scripts to override the `recv{}` as needed.

<sup>3</sup>Note that this implementation of session message handling is subtly different from that used in *wb* or described in [11]. In principle, an agent disseminates a list of the data it has actually received. Our implementation, on the other hand, only disseminates a count of the last message sequence number per source that the agent knows that the source has sent. This is a constraint when studying aspects of loss recovery during partition and healing. It is reasonable to expect that the maintainer of this code will fix this problem during one of his numerous intervals of copious spare time.

## 35.4 Loss Detection—The Class SRMinfo

A very small encapsulating class, entirely in C++, tracks a number of assorted state information. Each member of the group,  $n_i$ , uses one SRMinfo block for every other member of the group.

An SRMinfo object about group member  $n_j$  at  $n_i$ , contains information about the session messages received by  $n_i$  from  $n_j$ .  $n_i$  can use this information to compute its distance to  $n_j$ .

If  $n_j$  sends is active in sending data traffic, then the SRMinfo object will also contain information about the received data, including a bit vector indicating all packets received from  $n_j$ .

The agent keeps a list of SRMinfo objects, one per group member, in its member variable, `sip_`. Its method, `get_state(int sender)` will return the object corresponding to that sender, possibly creating that object, if it did not already exist. The class SRMinfo has two methods to access and set the bit vector, *i.e.*,

```
ifReceived(int id)  indicates whether the particular message from the appropriate sender, with id id was received
                    at  $n_i$ ,
setReceived(int id) to set the bit to indicate that the particular message from the appropriate sender, with id id was
                    received at  $n_i$ .
```

The session message variables to access timing information are public; no encapsulating methods are provided. These are:

```
int      lsess_;           /* # of last session msg received */
int      sendTime_;        /* Time sess. msg. # sent */
int      rcvTime_;        /* Time sess. msg. # received */
double   distance_;

/* Data messages */
int      ldata_;           /* # of last data msg sent */
```

## 35.5 Loss Recovery Objects

In the last section, we described the agent behavior when it receives a message. Timers are used to control when any particular control message is to be sent. The SRM agent uses a separate class SRM to do the timer based processing. In this section, we describe the basics of the class SRM, and the loss recovery objects. The following section will describe how the class SRM is used for sending periodic session messages. An SRM agent will instantiate one object to recover from one lost data packet. Agents that detect the loss will instantiate an object in the class SRM/request; agents that receive a request and have the required data will instantiate an object in the class SRM/repair.

**Request Mechanisms** SRM agents detect loss when they receive a message, and infer the loss based on the sequence number on the message received. Since packet reception is handled entirely by the compiled object, loss detection occurs in the C++ methods. Loss recovery, however, is handled entirely by instance procedures of the corresponding interpreted object in OTcl.

When any of the methods detects new losses, it invokes `Agent/SRM::request{}` with a list of the message sequence numbers that are missing. `request{}` will create a new `requestFunction_` object for each message that is missing. The agent stores the object handle in its array of `pending_` objects. The key to the array is the message identifier `<sender>:(msgid)`.

The default `requestFunction_` is class `SRM/request`. The constructor for the class `SRM/request` calls the base class constructor to initialize the simulator instance (`ns_`), the SRM agent (`agent_`), trace file (`trace_`), and the `times_` array. It then initializes its `statistics_` array with the pertinent elements.

A separate call to `set-params{}` sets the `sender_`, `msgid_`, `round_` instance variables for the request object. The object determines `C1_` and `C2_` by querying its `agent_`. It sets its distance to the sender (`times_(distance)`) and fixes other scheduling parameters: the backoff constant (`backoff_`), the current number of backoffs (`backoffCtr_`), and the limit (`backoffLimit_`) fixed by the agent. `set-params{}` writes the trace entry “Q DETECT”.

The final step in `request{}` is to schedule the timer to send the actual request at the appropriate moment. The instance procedure `SRM/request::schedule{}` uses `compute-delay{}` and its current backoff constant to determine the delay. The object schedules `send-request{}` to be executed after `delay_` seconds. The instance variable `eventID_` stores a handle to the scheduled event. The default `compute-delay{}` function returns a value uniformly distributed in the interval  $[C_1 d_s, (C_1 + C_2) d_s]$ , where  $d_s$  is twice `times_(distance)`. The `schedule{}` schedules an event to send a request after the computed delay. The routine writes a trace entry “Q NTIMER at (time)”.

When the scheduled timer fires, the routine `send-request{}` sends the appropriate message. It invokes “`$agent_ send request (args)`” to send the request. Note that `send{}` is an instproc-like, executed by the `command()` method of the compiled object. However, it is possible to overload the instproc-like with a specific instance procedure `send{}` for specific configurations. As an example, recall that the file `tcl/mcast/srm-nam.tcl` overloads the `send{}` command to set the flowid based on type of message that is sent. `send-request{}` updates the statistics, and writes the trace entry “Q SENDNACK”.

When the agent receives a control message for a packet for which a pending object exists, the agent will hand the message off to the object for processing.

When a request for a particular packet is received, the request object can be in one of two states: it is ignoring requests, considering them to be duplicates, or it will cancel its send event and re-schedule another one, after having backed off its timer. If ignoring requests it will update its statistics, and write the trace entry “Q NACK dup”. Otherwise, set a time based on its current estimate of the `delay_`, until which to ignore further requests. This interval is marked by the instance variable `ignore_`. If the object reschedules its timer, it will write the trace entry “Q NACK IGNORE-BACKOFF (ignore)”. Note that this re-scheduling relies on the fact that the agent has joined the multicast group, and will therefore receive a copy of every message it sends out.

When the request object receives a repair for the particular packet, it can be in one of two states: either it is still waiting for the repair, or it has already received an earlier repair. If it is the former, there will be an event pending to send a request, and `eventID_` will point to that event. The object will compute its `serviceTime`, cancel that event, and set a hold-down period during which it will ignore other requests. At the end of the hold-down period, the object will ask its agent to clear it. It will write the trace entry “Q REPAIR IGNORES (ignore)”. On the other hand, if this is a duplicate repair, the object will update its statistics, and write the trace entry “Q REPAIR dup”.

When the loss recovery phase is completed by the object, `Agent/SRM::clear{}` will remove the object from its array of `pending_` objects, and place it in its list of `done_` objects. Periodically, the agent will cleanup and delete the `done_` objects.

**Repair Mechanisms** The agent will initiate a repair if it receives a request for a packet, and it does not have a request object `pending_` for that packet. The default repair object belongs to the class `SRM/repair`. Barring minor differences, the sequence of events and the instance procedures in this class are identical to those for `SRM/request`. Rather than outline every single procedure, we only outline the differences from those described earlier for a request object.

The repair object uses the repair parameters, `D1_`, `D2_`. A repair object does not repeatedly reschedule its timers; therefore, it does not use any of the backoff variables such as that used by a request object. The repair object ignores all requests for the

same packet. The repair object does not use the `ignore_` variable that request objects use. The trace entries written by repair objects are marginally different; they are “P NACK from  $\langle$ requester $\rangle$ ”, “P RTIMER at  $\langle$ fi reTime $\rangle$ ”, “P SENDREP”, “P REPAIR IGNORES  $\langle$ holddown $\rangle$ ”.

Apart from these differences, the calling sequence for events in a repair object is similar to that of a request object.

**Mechanisms for Statistics** The agent, in concert with the request and repair objects, collect statistics about their response to data loss [11]. Each call to the agent `request{}` procedure marks a new period. At the start of a new period, `mark-period{}` computes the moving average of the number of duplicates in the last period. Whenever the agent receives a fi rst round request from another agent, and it had sent a request in that round, then it considers the request as a duplicate request, and increments the appropriate counters. A request object does not consider duplicate requests if it did not itself send a request in the fi rst round. If the agent has a repair object pending, then it does not consider the arrival of duplicate requests for that packet. The object methods `SRM/request::dup-request?{}` and `SRM/repair::dup-request?{}` encode these policies, and return 0 or 1 as required.

A request object also computes the elapsed time between when the loss is detected to when it receives the fi rst request. The agent computes a moving average of this elapsed time. The object computes the elapsed time (or delay) when it cancels its scheduled event for the fi rst round. The object invokes `Agent/SRM::update-ave` to compute the moving average of the delay.

The agent keeps similar statistics of the duplicate repairs, and the repair delay.

The agent stores the number of rounds taken for one loss recovery, to ensure that subsequent loss recovery phases for that packet that are not defi nitely not due to data loss do not account for these statistics. The agent stores the number of routes taken for a phase in the array `old_`. When a new loss recovery object is instantiated, the object will use the agent’s instance procedure `round?{}` to determine the number of rounds in a previous loss recovery phase for that packet.

## 35.6 Session Objects

Session objects, like the loss recovery objects (Section 35.5), are derived from the base class `SRM`. Unlike the loss recovery objects though, the agent only creates one session object for the lifetime of the agent. The constructor invokes the base class constructor as before; it then sets its instance variable `sessionDelay_`. The agent creates the session object when it `start{}`s. At that time, it also invokes `SRM/session::schedule`, to send a session message after `sessionDelay_` seconds.

When the object sends a session message, it will schedule to send the next one after some interval. It will also update its statistics. `send-session{}` writes out the trace entry “S SESSION”.

The class overrides the `evTrace{}` routine that writes out the trace entries. `SRM/session::evTrace` disable writing out the trace entry for session messages.

Two types of session message scheduling strategies are currently available: The function in the base class schedules sending session messages at fi xed intervals of `sessionDelay_` jittered around a small value to avoid synchronization among all the agents at all the nodes. class `SRM/session/logScaled` schedules sending messages at intervals of `sessionDelay` times  $\log_2(\text{groupSize}_)$  so that the frequency of session messages is inversely proportional to the size of the group.

The base class that sends messages at fi xed intervals is the default `sessionFunction_` for the agent.

## 35.7 Extending the Base Class Agent

In the earlier section on configuration parameters (Section 35.1.2), we had shown how to trivially extend the agent to get deterministic and probabilistic protocol behavior. In this section, we describe how to derive more complex extensions to the protocol for fixed and adaptive timer mechanisms.

### 35.7.1 Fixed Timers

The fixed timer mechanism are done in the derived class `Agent/SRM/Fixed`. The main difference with fixed timers is that the repair parameters are set to  $\log(\text{groupSize}_i)$ . Therefore, the repair procedure of a fixed timer agent will set  $D_1$  and  $D_2$  to be proportional to the group size before scheduling the repair object.

### 35.7.2 Adaptive Timers

Agents using adaptive timer mechanisms modify their request and repair parameters under three conditions (1) every time a new loss object is created; (2) when sending a message; and (3) when they receive a duplicate, if their relative distance to the loss is less than that of the agent that sends the duplicate. All three changes require extensions to the agent and the loss objects. The class `Agent/SRM/Adaptive` uses class `SRM/request/Adaptive` and class `SRM/repair/Adaptive` as the request and repair functions respectively. In addition, the last item requires extending the packet headers, to advertise their distances to the loss. The corresponding compiled class for the agent is the class `ASRMAgent`.

**Recompute for Each New Loss Object** Each time a new request object is created, `SRM/request/Adaptive::set-params` invokes `$agent_ recompute-request-params`. The agent method `recompute-request-params()` uses the statistics about duplicates and delay to modify  $C_1$  and  $C_2$  for the current and future requests.

Similarly, `SRM/request/Adaptive::set-params` for a new repair object invokes `$agent_ recompute-repair-params`. The agent method `recompute-repair-params()` uses the statistics objects to modify  $D_1$  and  $D_2$  for the current and future repairs.

**Sending a Message** If a loss object sends a request in its first round, then the agent, in the instance procedure `sending-request{}`, will lower  $C_1$ , and set its instance variable `closest_(requestor)` to 1.

Similarly, a loss object that sends a repair in its first round will invoke the agent's instance procedure, `sending-repair{}`, to lower  $D_1$  and set `closest_(repairor)` to 1.

**Advertising the Distance** Each agent must add additional information to each request/repair that it sends out. The base class `SRMAgent` invokes the virtual method `addExtendedHeaders()` for each SRM packet that it sends out. The method is invoked after adding the SRM packet headers, and before the packet is transmitted. The adaptive SRM agent overloads `addExtendedHeaders()` to specify its distances in the additional headers. When sending a request, that agent unequivocally knows the identity of the sender. As an example, the definition of `addExtendedHeaders()` for the adaptive SRM agent is:

```
void addExtendedHeaders(Packet* p) {
    SRMInfo* sp;
    hdr_srm* sh = (hdr_srm*) p->access(off_srm);
```

```

        hdr_asrm* seh = (hdr_asrm*) p->access(off_asrm_);
        switch (sh->type()) {
        case SRM_RQST:
            sp = get_state(sh->sender());
            seh->distance() = sp->distance_;
            break;

            ...
        }
    }
}

```

Similarly, the method `parseExtendedHeaders()` is invoked every time an SRM packet is received. It sets the agent member variable `pdistance_` to the distance advertised by the peer that sent the message. The member variable is bound to an instance variable of the same name, so that the peer distance can be accessed by the appropriate instance procedures. The corresponding `parseExtendedHeaders()` method for the Adaptive SRM agent is simply:

```

void parseExtendedHeaders(Packet* p) {
    hdr_asrm* seh = (hdr_asrm*) p->access(off_asrm_);
    pdistance_ = seh->distance();
}

```

Finally, the adaptive SRM agent's extended headers are defined as `struct hdr_asrm`. The header declaration is identical to declaring other packet headers in *ns*. Unlike most other packet headers, these are not automatically available in the packet. The interpreted constructor for the first adaptive agent will add the header to the packet format. For example, the start of the constructor for the `Agent/SRM/Adaptive` agent is:

```

Agent/SRM/Adaptive set done_ 0
Agent/SRM/Adaptive instproc init args {
    if ![${class set done_}] {
        set pm [[Simulator instance] set packetManager_]
        TclObject set off_asrm_ [$pm allochdr aSRM]
        $class set done_ 1
    }

    eval $self next $args
    ...
}

```

## 35.8 SRM objects

SRM objects are a subclass of agent objects that implement the SRM reliable multicast transport protocol. They inherit all of the generic agent functionalities. The methods for this object is described in the next section 35.9. Configuration parameters for this object are:

**packetSize\_** The data packet size that will be used for repair messages. The default value is 1024.

**requestFunction\_** The algorithm used to produce a retransmission request, e.g., setting request timers. The default value is `SRM/request`. Other possible request functions are `SRM/request/Adaptive`, used by the Adaptive SRM code.

**repairFunction\_** The algorithm used to produce a repair, e.g., compute repair timers. The default value is `SRM/repair`. Other possible request functions are `SRM/repair/Adaptive`, used by the Adaptive SRM code.

**sessionFunction\_** The algorithm used to generate session messages. Default is SRM/session

**sessionDelay\_** The basic interval of session messages. Slight random variation is added to this interval to avoid global synchronization of session messages. User may want to adjust this variable according to their specific simulation. Default value is 1.0.

**C1\_, C2\_** The parameters which control the request timer. Refer to [8] for detail. The default value is  $C1_ = C2_ = 2.0$ .

**D1\_, D2\_** The parameters which control the repair timer. Refer to [8] for detail. The default value is  $D1_ = D2_ = 1.0$ .

**requestBackoffLimit\_** The maximum number of exponential backoffs. Default value is 5.

State Variables are:

**stats\_** An array containing multiple statistics needed by adaptive SRM agent. Including: duplicate requests and repairs in current request/repair period, average number of duplicate requests and repairs, request and repair delay in current request/repair period, average request and repair delay.

**SRM/ADAPTIVE OBJECTS** SRM/Adaptive objects are a subclass of the SRM objects that implement the adaptive SRM reliable multicast transport protocol. They inherit all of the SRM object functionalities. State Variables are:

(Refer to the SRM paper by Sally et al [Fall, K., Floyd, S., and Henderson, T., Ns Simulator Tests for Reno FullTCP. URL <ftp://ftp.ee.lbl.gov/papers/fulltcp.ps>. July 1997.] for more detail.)

**pdistance\_** This variable is used to pass the distance estimate provided by the remote agent in a request or repair message.

**D1\_, D2\_** The same as that in SRM agents, except that they are initialized to  $\log_{10}(\text{group size})$  when generating the first repair.

**MinC1\_, MaxC1\_, MinC2\_, MaxC2\_** The minimum/maximum values of  $C1_$  and  $C2_$ . Default initial values are defined in [8]. These values define the dynamic range of  $C1_$  and  $C2_$ .

**MinD1\_, MaxD1\_, MinD2\_, MaxD2\_** The minimum/maximum values of  $D1_$  and  $D2_$ . Default initial values are defined in [8]. These values define the dynamic range of  $D1_$  and  $D2_$ .

**AveDups** Higher bound for average duplicates.

**AveDelay** Higher bound for average delay.

**eps AveDups** -dups determines the lower bound of the number of duplicates, when we should adjust parameters to decrease delay.

## 35.9 Commands at a glance

The following is a list of commands to create/manipulate srm agents in simulations:

```
set srm0 [new Agent/SRM]
```

This creates an instance of the SRM agent. In addition to the base class, two extensions of the srm agent have been implemented. They are Agent/SRM/Fixed and Agent/SRM/Adaptive. See section 35.7 for details about these extensions.

```
ns_ attach-agent <node> <srm-agent>
```

This attaches the srm agent instance to the given <node>.

```
set grp [Node allocaddr]
$srn set dst_ $grp
```

This assigns the srm agent to a multicast group represented by the mcast address <grp>.

Configuration parameters for srm agent may be set as follows:

```
$srn set fid_ <flow-id>
$srn set tg_ <traffic-generator-instance>
.. etc
```

For all possible parameters and their default values please lookup *ns/tcl/mcast/srm.tcl* and *ns/tcl/mcast/srm-adaptive.tcl*.

```
set exp [new Application/Traffic/Exponential]
$exp attach-agent $srn
```

This command attaches a traffic generator (an exponential one in this example), to the srm agent.

```
$srn start; $exp start
```

These commands start the srm agent and traffic generator. Note that the srm agent and traffic generator have to be started separately. Alternatively, the traffic generator may be started through the agent as follows:

```
$srn start-source.
```

See *ns/tcl/ex/srm.tcl* for a simple example of setting up a SRM agent.



# Chapter 36

## PLM

This chapter describes the ns implementation of the PLM protocol [19]. The code of the PLM protocol is written in both C++ and OTcl. The PLM Packet Pair generator is written in C++ and the PLM core machinery is written in OTcl. The chapter has simply three parts: the first part shows how to create and configure a PLM session; the second part describes the Packet Pair source generator; the third part describes the architecture and internals of the PLM protocol. In this last part, rather than giving a list of procedures and functions, we introduce the main procedures per functionality (instantiation of a PLM source, instantiation of a PLM receiver, reception of a packet, detection of a loss, etc.).

The procedures, functions, and variables described in this chapter can be found in: `~ns/plm/cbr-traffic-PP.cc`, `~ns/plm/loss-monitor-plm.cc`, `~ns/tcl/plm/plm.tcl`, `~ns/tcl/plm/plm-ns.tcl`, `~ns/tcl/plm/plm-topo.tcl`, `~ns/tcl/lib/ns-default.tcl`.

### 36.1 Configuration

#### Creating a simple scenario with one PLM flow (only one receiver)

This simple example can be run as is (several complex scenarios can be found in the file `~ns/tcl/ex/simple-plm.tcl`).

```
set packetSize 500 ; #Packet size (in bytes)
set plm_debug_flag 2 ; #Debugging output
set rates "50e3 50e3 50e3 50e3 50e3" ; #Rate of each layer
set rates_cum [calc_cum $rates] ; #Cumulated rate of the layers (mandatory)
set level [llength $rates] ; #Number of layers (mandatory)

set Queue_sched_ FQ ; #Scheduling of the queues
set PP_burst_length 2 ; #PP burst length (in packets)
set PP_estimation_length 3 ; #Minimum number of PP required to make an estimate

Class Scenario0 -superclass PLMTopology
Scenario0 instproc init args {
    eval $self next $args
    $self instvar ns node

    $self build_link 1 2 100ms 256Kb ; #Build a link
    set addr(1) [$self place_source 1 3] ; #Set a PLM source
    $self place_receiver 2 $addr(1) 5 1 ; #Set a PLM receiver
```

```

#set up the multicast routing
DM set PruneTimeout 1000                                ; #A large PruneTimeout value is required
set mproto DM
set mrthandle [$ns mrtproto $mproto {} ]
}

set ns [new Simulator -multicast on]                      ; #PLM needs multicast routing
$ns multicast
$ns namtrace-all [open out.nam w]                        ; #Nam output
set scn [new Scenario0 $ns]                               ; #Call of the scenario
$ns at 20 "exit 0"
$ns run

```

Several variables are introduced in this example. They all need to be set in the simulation script (there is no default value for these variables). In particular the two following lines are mandatory and must not be omitted:

```

set rates_cum [calc_cum $rates]
set level [llength $rates]

```

We describe now in detail each variable:

`packetSize` represents the size of the packets in bytes sent by the PLM source.

`plm_debug_flag` represents the verbose level of debugging output: from 0 no output to 3 full output. For `plm_debug_flag` set to 3 (full output), long lines output are generated which is not compatible with nam visualization.

`rates` is a list specifying the bandwidth of each layer (this is not the cumulated bandwidth!).

`rates_cum` is a list specifying the cumulated bandwidth of the layers: the first element of `rates_cum` is the bandwidth a layer 1, the second element of `rates_cum` is the sum of the bandwidth of layer 1 and layer 2, etc. The proc `calc_cum{}` computes the cumulated rates.

`level` is the number of layers.

`Queue_sched_` represents the scheduling of the queues. This is used by the PLMTopology instproc `build_link`. PLM requires FQ scheduling or a variation.

`PP_burst_length` represents the size of the Packet Pair bursts in packets.

`PP_estimation_length` represents the minimum number of Packet Pair required to compute an estimate (see section 36.3.3).

All the simulations for PLM should be setup using the PLMTopology environment (as in the example script where we define a PLMTopology superclass called Scenario0). The user interface is (all the instproc can be found in `~ns/tcl/plm/plm-topo.tcl`):

`build_link a b d bw` creates a duplex link between node a and b with a delay d and a bandwidth bw. If either node does not exist, `build_link` creates it.

`place_source n t` creates and places a PLM source at node n and starts it at time t. `place_source` returns addr which allows to attach receivers to this source.

`place_receiver n addr C nb` creates and places a PLM receiver at node `n` and attached it to the source which return the address `addr`. The check value for this PLM receiver is `C`. An optional parameter `nb` allows to get an instance of the PLM receiver called `PLMrcvr($nb)`. This instance is only useful to get some specific statistics about this receiver (mainly the number of packets received or lost).

## 36.2 The Packet Pair Source Generator

This section describes the Packet Pair source generator; the relevant files are: `~ns/plm/cbr-traffic-PP.cc`, `~ns/tcl/lib/ns-default.tcl`. The OTcl class name of the PP source is: `Application/Traffic/CBR_PP`. The Packet Pair (PP) source generator is in the file `~ns/plm/cbr-traffic-PP.cc`. This source generator is a variation of the CBR source generator in `~ns/cbr-traffic.cc`. We just describe the salient differences between the code of the CBR source and the code of the PP source. The default values in `~ns/tcl/lib/ns-default.tcl` for the PP source generator are the same than for the CBR source. We need for the PP source generator a new parameter `PBM_`:

```
Application/Traffic/CBR_PP set PBM_ 2 ; #Default value
```

The OTcl instvar bounded variable `PBM_` (same name in C++ and in OTcl) specifies the number of back-to-back packets to be sent. For `PBM_=1` we have a CBR source, for `PBM_=2` we have a Packet Pair source (a source which sends two packets back-to-back), etc. The mean rate of the PP source is `rate_`, but the packets are sent in burst of `PBM_` packets. Note that we also use the terminology Packet Pair source and Packet Pair burst for `PBM_ > 2`. We compute the `next_interval` as:

```
double CBR_PP_Traffic::next_interval(int& size)

/*(PP_ - 1) is the number of packets in the current burst.*/
    if (PP_ >= (PBM_ - 1))
        interval_ = PBM_*(double)(size_ << 3)/(double)rate_;
        PP_ = 0;

    else
        interval_ = 1e-100; //zero
        PP_ += 1 ;

...

```

The `timeout{}` method puts the `NEW_BURST` flag in the first packet of a burst. This is useful for the PLM protocol to identify the beginning of a PP burst.

```
void CBR_PP_Traffic::timeout()

...
if (PP_ == 0)
    agent_->sendmsg(size_, "NEW_BURST");
else
    agent_->sendmsg(size_);

...

```

## 36.3 Architecture of the PLM Protocol

The code of the PLM protocol is divided in three files: `~ns/tcl/plm/plm.tcl`, which contains the PLM protocol machinery without any specific interface with *ns*; `~ns/tcl/plm/plm-ns.tcl`, which contains the specific *ns* interface. However, we do not guarantee the strict validity of this *ns* interfacing; `~ns/tcl/plm/plm-topo.tcl`, which contains a user interface to build simulation scenarios with PLM flows.

In the following we do not discuss the various procedures per object (for instance all the instproc of the PLM class) but rather per functionality (for instance which instproc among the various classes are involved in the instantiation of a PLM receiver). For a given functionality, we do not describe in details all the code involved, but we give the principal steps.

### 36.3.1 Instantiation of a PLM Source

To create a PLM source, place it at node *n*, and start it at  $t_0$ , we call the PLMTopology instproc `place_source n t_0`. This instproc return `addr`, the address required to attach a receiver to this source. `place_source` calls the Simulator instproc `PLMbuild_source_set` that creates as many Application/Traffic/CBR\_PP instances as there are layers (in the following we call an instance of the class Application/Traffic/CBR\_PP a layer). Each layer corresponds to a different multicast group.

To speed up the simulations when the PLM sources start we use the following trick: At  $t = 0$ , `PLMbuild_source_set` restricts each layer to send only one packet (`maxpkts_` set to 1). That allows to build the multicast trees – one multicast tree per layer – without flooding the whole network. Indeed, each layer only sends one packet to build the corresponding multicast tree.

The multicast trees take at most the maximum RTT of the network to be established and must be established before  $t_0$ , the PLM source starting time. Therefore,  $t_0$  must be carefully chosen, otherwise the source sends a large number of useless packets. However, as we just need to start the PLM source after the multicast trees are established,  $t_0$  can be largely overestimated. At time  $t_0$ , we set `maxpkts_` to 268435456 for all the layers.

It is fundamental, in order to have persistent multicast trees, that the prune timeout is set to a large value. For instance, with DM routing:

```
DM set PruneTimeout 1000
```

Each layer of a same PLM source has the same flow id `fid_`. Consequently, each PLM source is considered as a single flow for a Fair Queueing scheduler. The PLM code manages automatically the `fid_` to prevent different sources to have the same `fid_`. The `fid_` starts at 1 for the first source and is increased by one for each new source. Be careful to avoid other flows (for instance concurrent TCP flows) to have the same `fid_` than the PLM sources. Additionally, If you consider `fid_` larger than 32, do not forget to increase the MAXFLOW in `~ns/fq.cc` (MAXFLOW must be set to the highest `fid_` considered in the simulation).

### 36.3.2 Instantiation of a PLM Receiver

All the PLM machinery is implemented at the receiver. In this section we describe the instantiation process of a receiver. To create, place at node *n*, attach to source *S*, and start at  $t_1$  a PLM receiver we call the PLMTopology instproc `build_receiver n addr t_1 C` where `addr` is the address returned by `place_source` when *S* was created, and *C* is the check value. The receiver created by `build_receiver` is an instance of the class PLM/ns, the *ns* interface to the PLM machinery. At the initialisation of the receiver, the PLM instproc `init` is called due to inheritance. `init` calls the PLM/ns instproc

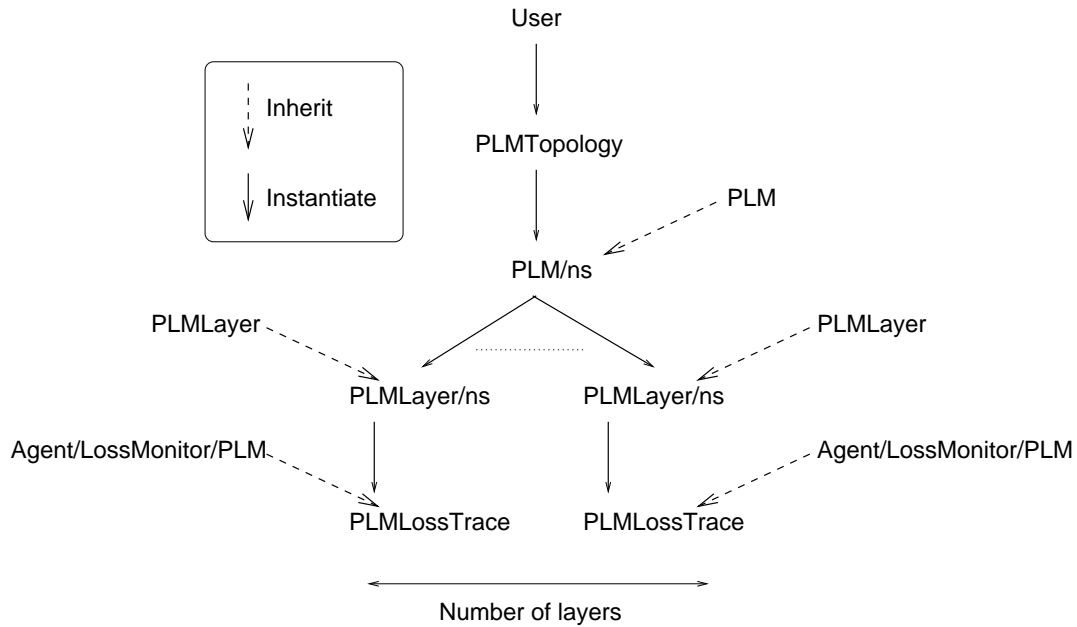


Figure 36.1: Inheritance and instantiation when we create a receiver.

create-layer and, by this way, creates as many instances of the class PLMLayer/ns (the ns interface to the PLMLayer class) as there are layers. Each instance of PLMLayer/ns creates an instance of the class PLMLossTrace which is responsible for monitoring the received and lost packets thanks to the fact that the class PLMLossTrace inherits from the class Agent/LossMonitor/PLM. Fig. 36.1 schematically describes the process of a PLM receiver instantiation. In the following we describe the behavior of a PLM receiver when it receives a packet and when it detects a loss.

### 36.3.3 Reception of a Packet

We create a new c++ class PLMLossMonitor (~ns/plm/loss-monitor-plm.cc) that inherits from LossMonitor. The OTcl class name of the c++ PLMLossMonitor class is Agent/LossMonitor/PLM.

```

class PLMLossMonitor : public LossMonitor
public:
    PLMLossMonitor();
    virtual void recv(Packet* pkt, Handler*);
protected:
    // PLM only
    int flag_PP_;
    double packet_time_PP_;
    int fid_PP_;
;

static class PLMLossMonitorClass : public TclClass
public:
    PLMLossMonitorClass() : TclClass("Agent/LossMonitor/PLM")
    TclObject* create(int, const char*const*)
        return (new PLMLossMonitor());

class_loss_mon_plm;

```

We add in `void PLMLossMonitor::recv(Packet* pkt, Handler*)` a Tcl call to the Agent/LossMonitor/PLM instproc `log-PP` each time a packet is received :

```
void LossMonitor::recv(Packet* pkt, Handler*)

...
if (expected_ >= 0)
...

Tcl::instance().evalf("%s log-PP", name());
```

The Agent/LossMonitor/PLM instproc `log-PP` is empty. In fact, we define the `log-PP` instproc for the class `PLMLossTrace`. `log-PP` computes an estimate of the available bandwidth based on a single PP burst (of length `PP_burst_length` in packets). Once `log-PP` has received the `PP_burst_length` packets of the burst, it computes the estimate and calls the PLM instproc `make_estimate` with the computed estimate as argument.

`make_estimate` puts the estimate based on a single PP (`PP_value`) in an array of estimate samples (`PP_estimate`). If `PP_value` is lower than the current subscription level (i.e. lower than the throughput achieved according to the current number of layers subscribed), `make_estimate` calls the PLM instproc `stability-drop` which simply drops layers until the current subscription level becomes lower than `PP_value`. `make_estimate` makes an estimate `PP_estimate_value` by taking the minimum `PP_value` received during the last `check_estimate` period (if there are at least `PP_estimation_length` single PP estimate received). Once `make_estimate` has a `PP_estimate_value` it calls the PLM instproc `choose_layer` which joins or drops layer(s) according to the current subscription level and to the `PP_estimate_value`. For details about the PLM instproc `make_estimate`, refer to its code in `~ns/tcl/plm/plm.tcl`.

### 36.3.4 Detection of a Loss

Each time a loss is detected by an instance of the class `PLMLossMonitor`, a call to the Agent/LossMonitor/PLM instproc `log-loss` is triggered. The Agent/LossMonitor/PLM instproc `log-loss` is empty. In fact, we define the `log-loss` instproc for the class `PLMLossTrace`. The `PLMLossTrace` instproc `log-loss` simply calls the PLM instproc `log-loss` which contains the PLM machinery in case of loss. In summary, `log-loss` only drops a layer when the loss rate exceeds 10% (this test is executed by the PLM instproc `exceed_loss_thresh`). After a layer drop `log-loss` precludes any other layer drop due to loss for 500ms. For details about the PLM instproc `log-loss`, refer to its code in `~ns/tcl/plm/plm.tcl`.

### 36.3.5 Joining or Leaving a Layer

To join a layer the PLM instproc `add-layer` is called. This instproc calls the `PLMLayer` instproc `join-group` which calls the `PLMLayer/ns` instproc `join-group`. To leave a layer the PLM instproc `drop-layer` is called. This instproc calls the `PLMLayer` instproc `leave-group` which calls the `PLMLayer/ns` instproc `leave-group`.

## 36.4 Commands at a Glance

Note: This section is a copy paste of the end of section 36.1. We add this section to preserve homogeneity with the ns manual.

All the simulations for PLM should be set using the PLMTopology environment (as in the example script where we define a PLMTopology superclass called Scenario0). The user interface is (all the instproc can be found in `~ns/tcl/plm/plm-topo.tcl`):

`build_link a b d bw` creates a duplex link between node `a` and `b` with a delay `d` and a bandwidth `bw`. If either node does not exist, `build_link` creates it.

`place_source n t` creates and places a PLM source at node `n` and starts it at time `t`. `place_source` returns `addr` which allows to attach receivers to this source.

`place_receiver n addr C nb` creates and places a PLM receiver at node `n` and attached it to the source which return the address `addr`. The check value for this PLM receiver is `C`. An optional parameter `nb` allows to get an instance of the PLM receiver called `PLMrcvr($nb)`. This instance is only useful to get some specific statistics about this receiver (mainly the number of packets received or lost).

# **Part VI**

## **Application**



## Chapter 37

# Applications and transport agent API

Applications sit on top of transport agents in *ns*. There are two basic types of applications: traffic generators and simulated applications. Figure 37.1 illustrates two examples of how applications are composed and attached to transport agents. Transport agents are described in Part V (Transport).

This chapter first describes the base class `Application`. Next, the transport API, through which applications request services from underlying transport agents, is described. Finally, the current implementations of traffic generators and sources are explained.

### 37.1 The class `Application`

`Application` is a C++ class defined as follows:

```
class Application : public TclObject {
public:
    Application();
    virtual void send(int nbytes);
    virtual void recv(int nbytes);
    virtual void resume();
protected:
    int command(int argc, const char*const* argv);
    virtual void start();
    virtual void stop();
    Agent *agent_;
    int enableRecv_;           // call OTcl recv or not
    int enableResume_;        // call OTcl resume or not
};
```

Although objects of class `Application` are not meant to be instantiated, we do not make it an abstract base class so that it is visible from OTcl level. The class provides basic prototypes for application behavior (`send()`, `recv()`, `resume()`, `start()`, `stop()`), a pointer to the transport agent to which it is connected, and flags that indicate whether a OTcl-level upcall should be made for `recv()` and `resume()` events.

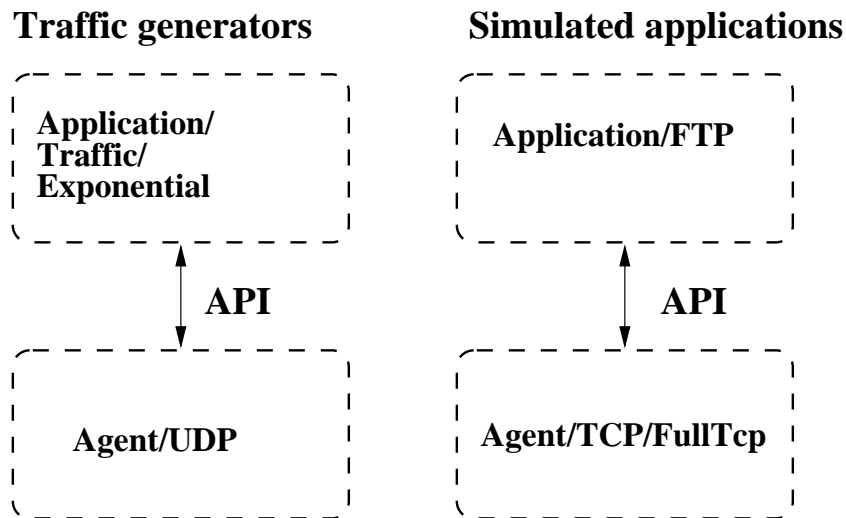


Figure 37.1: Example of Application Composition

## 37.2 The transport agent API

In real-world systems, applications typically access network services through an applications programming interface (API). The most popular of these APIs is known as “sockets.” In *ns*, we mimic the behavior of the sockets API through a set of well-defined API functions. These functions are then mapped to the appropriate internal agent functions (e.g., a call to `send(numBytes)` causes TCP to increment its “send buffer” by a corresponding number of bytes).

This section describes how agents and applications are hooked together and communicate with one another via the API.

### 37.2.1 Attaching transport agents to nodes

This step is typically done at OTcl level. Agent management was also briefly discussed in Section 5.2.

```

set src [new Agent/TCP/FullTcp]
set sink [new Agent/TCP/FullTcp]
$ns_ attach-agent $node_(s1) $src
$ns_ attach-agent $node_(k1) $sink
$ns_ connect $src $sink

```

The above code illustrates that in *ns*, agents are first attached to a node via `attach-agent`. Next, the `connect` instproc sets each agent’s destination target to the other. Note that, in *ns*, `connect()` has different semantics than in regular sockets. In *ns*, `connect()` simply establishes the destination address for an agent, but does not set up the connection. As a result, the overlying application does not need to know its peer’s address. For TCPs that exchange SYN segments, the first call to `send()` will trigger the SYN exchange.

To detach an agent from a node, the instproc `detach-agent` can be used; this resets the target for the agent to a null agent.

### 37.2.2 Attaching applications to agents

After applications are instantiated, they must be connected to a transport agent. The `attach-agent` method can be used to attach an application to an agent, as follows:

```
set ftp1 [new Application/FTP]
$ftp1 attach-agent $src
```

The following shortcut accomplishes the same result:

```
set ftp1 [$src attach-app FTP]
```

The `attach-agent` method, which is also used by `attach-app`, is implemented in C++. It sets the `agent_` pointer in class `Application` to point to the transport agent, and then it calls `attachApp()` in `agent.cc` to set the `app_` pointer to point back to the application. By maintaining this binding only in C++, OTcl-level instvars pointers are avoided and consistency between OTcl and C++ is guaranteed. The OTcl-level command `[$ftp1 agent]` can be used by applications to obtain the handler for the transport agent.

### 37.2.3 Using transport agents via system calls

Once transport agents have been configured and applications attached, applications can use transport services via the following system calls. These calls can be invoked at either OTcl or C++ level, thereby allowing applications to be coded in either C++ or OTcl. These functions have been implemented as virtual functions in the base class `Agent`, and can be redefined as needed by derived Agents.

- `send(int nbytes)`—Send `nbytes` of data to peer. For TCP agents, if `nbytes == -1`, this corresponds to an “infinite” send; i.e., the TCP agent will act as if its send buffer is continually replenished by the application.
- `sendmsg(int nbytes, const char* flags = 0)`—Identical to `send(int nbytes)`, except that it passes an additional string `flags`. Currently one flag value, “MSG\_EOF,” is defined; `MSG_EOF` specifies that this is the last batch of data that the application will submit, and serves as an implied close (so that TCP can send FIN with data).
- `close()`—Requests the agent to close the connection (only applicable for TCP).
- `listen()`—Requests the agent to listen for new connections (only applicable for Full TCP).
- `set_pkttype(int pkttype)`—This function sets the `type_` variable in the agent to `pkttype`. Packet types are defined in `packet.h`. This function is used to override the transport layer packet type for tracing purposes.

Note that certain calls are not applicable for certain agents; e.g., a call to `close()` a UDP connection results in a no-op. Additional calls can be implemented in specialized agents, provided that they are made `public` member functions.

### 37.2.4 Agent upcalls to applications

Since presently in *ns* there is no actual data being passed between applications, agents can instead announce to applications the occurrence of certain events at the transport layer through “upcalls.” For example, applications can be notified of the arrival of a number of bytes of data; this information may aid the application in modelling real-world application behavior more closely. Two basic “upcalls” have been implemented in base class `Application` and in the transport agents:

- `recv(int nbytes)`—Announces that `nbytes` of data have been received by the agent. For UDP agents, this signifies the arrival of a single packet. For TCP agents, this signifies the “delivery” of an amount of in-sequence data, which may be larger than that contained in a single packet (due to the possibility of network reordering).
- `resume()`—This indicates to the application that the transport agent has sent out all of the data submitted to it up to that point in time. For TCP, it does not indicate whether the data has been ACKed yet, only that it has been sent out for the first time.

The default behavior is as follows: Depending on whether the application has been implemented in C++ or OTcl, these C++ functions call a similarly named (`recv`, `resume`) function in the application, if such methods have been defined.

Although strictly not a callback to applications, certain Agents have implemented a callback from C++ to OTcl-level that has been used by applications such as HTTP simulators. This callback method, `done{}`, is used in TCP agents. In TCP, `done{}` is called when a TCP sender has received ACKs for all of its data and is now closed; it therefore can be used to simulate a blocked TCP connection. The `done{}` method was primarily used before this API was completed, but may still be useful for applications that do not want to use `resume()`.

To use `done{}` for FullTcp, for example, you can try:

```
set myagent [new Agent/TCP/FullTcp]
$myagent proc done
... code you want ...
```

If you want all the FullTCP's to have the same code you could also do:

```
Agent/TCP/FullTcp instproc done
... code you want ...
```

By default, `done{}` does nothing.

### 37.2.5 An example

Here is an example of how the API is used to implement a simple application (FTP) on top of a FullTCP connection.

```
set src [new Agent/TCP/FullTcp]
set sink [new Agent/TCP/FullTcp]
$ns_ attach-agent $node_(s1) $src
$ns_ attach-agent $node_(k1) $sink
$ns_ connect $src $sink

# set up TCP-level connections
$sink listen;
$src set window_ 100

set ftp1 [new Application/FTP]
$ftp1 attach-agent $src

$ns_ at 0.0 "$ftp1 start"
```

In the configuration script, the first five lines of code allocates two new FullTcp agents, attaches them to the correct nodes, and "connects" them together (assigns the correct destination addresses to each agent). The next two lines configure the TCP agents further, placing one of them in LISTEN mode. Next, ftp1 is defined as a new FTP Application, and the attach-agent method is called in C++ (app.cc).

The ftp1 application is started at time 0:

```
Application/FTP instproc start {} {
    [$self agent] send -1;    # Send indefinitely
}
```

Alternatively, the FTP application could have been implemented in C++ as follows:

```
void FTP::start()
{
    agent_>send(-1);    // Send indefinitely
}
```

Since the FTP application does not make use of callbacks, these functions are null in C++ and no OTcl callbacks are made.

### 37.3 The class TrafficGenerator

TrafficGenerator is an abstract C++ class defined as follows:

```
class TrafficGenerator : public Application {
public:
    TrafficGenerator();
    virtual double next_interval(int &) = 0;
    virtual void init() {}
    virtual double interval() { return 0; }
    virtual int on() { return 0; }
    virtual void timeout();
    virtual void recv() {}
    virtual void resume() {}
protected:
    virtual void start();
    virtual void stop();
    double nextPkttime_;
    int size_;
    int running_;
    TrafficTimer timer_;
};
```

The pure virtual function next\_interval() returns the time until the next packet is created and also sets the size in bytes of the next packet. The function start() calls init(void) and starts the timer. The function timeout() sends a packet and reschedules the next timeout. The function stop() cancels any pending transmissions. Callbacks are typically not used for traffic generators, so these functions (recv, resume) are null.

Currently, there are four C++ classes derived from the class TrafficGenerator:

1. **EXPOO\_Traffic**—generates traffic according to an Exponential On/Off distribution. Packets are sent at a fixed rate during on periods, and no packets are sent during off periods. Both on and off periods are taken from an exponential distribution. Packets are constant size.
2. **POO\_Traffic**—generates traffic according to a Pareto On/Off distribution. This is identical to the Exponential On/Off distribution, except the on and off periods are taken from a pareto distribution. These sources can be used to generate aggregate traffic that exhibits long range dependency.
3. **CBR\_Traffic**—generates traffic according to a deterministic rate. Packets are constant size. Optionally, some randomizing dither can be enabled on the interpacket departure intervals.
4. **TrafficTrace**—generates traffic according to a trace file. Each record in the trace file consists of 2 32-bit fields in network (big-endian) byte order. The first contains the time in microseconds until the next packet is generated. The second contains the length in bytes of the next packet.

These classes can be created from OTcl. The OTcl classes names and associated parameters are given below:

**Exponential On/Off** An Exponential On/Off object is embodied in the OTcl class `Application/Traffic/Exponential`. The member variables that parameterize this object are:

<code>packetSize_</code>	the constant size of the packets generated
<code>burst_time_</code>	the average “on” time for the generator
<code>idle_time_</code>	the average “off” time for the generator
<code>rate_</code>	the sending rate during “on” times

Hence a new Exponential On/Off traffic generator can be created and parameterized as follows:

```
set e [new Application/Traffic/Exponential]
$e set packetSize_ 210
$e set burst_time_ 500ms
$e set idle_time_ 500ms
$e set rate_ 100k
```

**NOTE:** The Exponential On/Off generator can be configured to behave as a **Poisson process** by setting the variable `burst_time_` to 0 and the variable `rate_` to a very large value. The C++ code guarantees that even if the burst time is zero, at least one packet is sent. Additionally, the next interarrival time is the sum of the assumed packet transmission time (governed by the variable `rate_`) and the random variate corresponding to `idle_time_`. Therefore, to make the first term in the sum very small, make the burst rate very large so that the transmission time is negligible compared to the typical idle times.

**Pareto On/Off** A Pareto On/Off object is embodied in the OTcl class `Application/Traffic/Pareto`. The member variables that parameterize this object are:

<code>packetSize_</code>	the constant size of the packets generated
<code>burst_time_</code>	the average "on" time for the generator
<code>idle_time_</code>	the average "off" time for the generator
<code>rate_</code>	the sending rate during "on" times
<code>shape_</code>	the "shape" parameter used by the pareto distribution

A new Pareto On/Off traffic generator can be created as follows:

```

set p [new Application/Traffic/Pareto]
$p set packetSize_ 210
$p set burst_time_ 500ms
$p set idle_time_ 500ms
$p set rate_ 200k
$p set shape_ 1.5

```

**CBR** A CBR object is embodied in the OTcl class `Application/Traffic/CBR`. The member variables that parameterize this object are:

<code>rate_</code>	the sending rate
<code>interval_</code>	(Optional) interval between packets
<code>packetSize_</code>	the constant size of the packets generated
<code>random_</code>	flag indicating whether or not to introduce random “noise” in the scheduled departure times (default is off)
<code>maxpkts_</code>	the maximum number of packets to send (default is $2^{28}$ )

Hence a new CBR traffic generator can be created and parameterized as follows:

```

set e [new Application/Traffic/CBR]
$e set packetSize_ 48
$e set rate_ 64Kb
$e set random_ 1

```

The setting of a CBR object’s `rate_` and `interval_` are mutually exclusive (the interval between packets is maintained as an interval variable in C++, and some example *ns* scripts specify an interval rather than a rate). In a simulation, either a rate or an interval (but not both) should be specified for a CBR object.

**Traffic Trace** A Traffic Trace object is instantiated by the OTcl class `Application/Traffic/Trace`. The associated class `Tracefile` is used to enable multiple Traffic/Trace objects to be associated with a single trace file. The Traffic/Trace class uses the method `attach-tracefile` to associate a Traffic/Trace object with a particular Tracefile object. The method `filename` of the Tracefile class associates a trace file with the Tracefile object. The following example shows how to create two `Application/Traffic/Trace` objects, each associated with the same trace file (called “example-trace” in this example). To avoid synchronization of the traffic generated, random starting places within the trace file are chosen for each Traffic/Trace object.

```

set tfile [new Tracefile]
$tfile filename example-trace

set t1 [new Application/Traffic/Trace]
$t1 attach-tracefile $tfile
set t2 [new Application/Traffic/Trace]
$t2 attach-tracefile $tfile

```

### 37.3.1 An example

The following code illustrates the basic steps to configure an Exponential traffic source over a UDP agent, for traffic flowing from node `s1` to node `k1`:

```

set src [new Agent/UDP]
set sink [new Agent/UDP]
$ns_ attach-agent $node_(s1) $src
$ns_ attach-agent $node_(k1) $sink
$ns_ connect $src $sink

set e [new Application/Traffic/Exponential]
$e attach-agent $src
$e set packetSize_ 210
$e set burst_time_ 500ms
$e set idle_time_ 500ms
$e set rate_ 100k

$ns_ at 0.0 "$e start"

```

## 37.4 Simulated applications: Telnet and FTP

There are currently two “simulate application” classes derived from Application: Application/FTP and Application/Telnet. These classes work by advancing the count of packets available to be sent by a TCP transport agent. The actual transmission of available packets is still controlled by TCP’s flow and congestion control algorithm.

**Application/FTP** Application/FTP, implemented in OTcl, simulates bulk data transfer. The following are methods of the Application/FTP class:

<code>attach-agent</code>	attaches an Application/FTP object to an agent.
<code>start</code>	start the Application/FTP by calling the TCP agent’s <code>send(-1)</code> function, which causes TCP to behave as if the application were continuously sending new data.
<code>stop</code>	stop sending.
<code>produce n</code>	set the counter of packets to be sent to <i>n</i> .
<code>producemore n</code>	increase the counter of packets to be sent by <i>n</i> .
<code>send n</code>	similar to <code>producemore</code> , but sends <i>n</i> bytes instead of packets.

**Application/Telnet** Application/Telnet objects generate packets in one of two ways. If the member variable `interval_` is non-zero, then inter-packet times are chosen from an exponential distribution with average equal to `interval_`. If `interval_` is zero, then inter-arrival times are chosen according to the `tcplib` distribution (see `tcplib-telnet.cc`). The `start` method starts the packet generation process.

## 37.5 Applications objects

An application object may be of two types, a traffic generator or a simulated application. Traffic generator objects generate traffic and can be of four types, namely, exponential, pareto, CBR and traffic trace.

**Application/Traffic/Exponential objects** Exponential traffic objects generate On/Off traffic. During “on” periods, packets are generated at a constant burst rate. During “off” periods, no traffic is generated. Burst times and idle times are taken from exponential distributions. Configuration parameters are:



**PacketSize\_** constant size of packets generated.

**burst\_time\_** average on time for generator.

**idle\_time\_** average off time for generator.

**rate\_** sending rate during on time.

**Application/Traffic/Pareto** Application/Traffic/Pareto objects generate On/Off traffic with burst times and idle times taken from pareto distributions. Configuration parameters are:

**PacketSize\_** constant size of packets generated.

**burst\_time\_** average on time for generator.

**idle\_time\_** average off time for generator.

**rate\_** sending rate during on time.

**shape\_** the shape parameter used by pareto distribution.

**Application/Traffic/CBR** CBR objects generate packets at a constant bit rate.

`$cbr start`

Causes the source to start generating packets.

`$cbr stop`

Causes the source to stop generating packets.

Configuration parameters are:

**PacketSize\_** constant size of packets generated.

**rate\_** sending rate.

**interval\_** (optional) interval between packets.

**random\_** whether or not to introduce random noise in the scheduled departure times. default is off.

**maxpkts\_** maximum number of packets to send.

**Application/Traffic/Trace** Application/Traffic/Trace objects are used to generate traffic from a trace file. `$trace attach-tracefile tfile`

Attach the Tracefile object tfile to this trace. The Tracefile object specifies the trace file from which the traffic data is to be read. Multiple Application/Traffic/Trace objects can be attached to the same Tracefile object. A random starting place within the Tracefile is chosen for each Application/Traffic/Trace object.

There are no configuration parameters for this object.

A simulated application object can be of two types, Telnet and FTP.

**Application/Telnet** TELNET objects produce individual packets with inter-arrival times as follows. If interval\_ is non-zero, then inter-arrival times are chosen from an exponential distribution with average interval\_. If interval\_ is zero, then inter-arrival times are chosen using the "tcplib" telnet distribution.

`$telnet start`

Causes the Application/Telnet object to start producing packets.

`$telnet stop`

Causes the Application/Telnet object to stop producing packets.

`$telnet attach <agent>`

Attaches a Telnet object to agent.

Configuration Parameters are:

**interval\_** The average inter-arrival time in seconds for packets generated by the Telnet object.

**Application FTP** FTP objects produce bulk data for a TCP object to send.

`$ftp start`

Causes the source to produce `maxpkts_` packets.

`$ftp produce <n>`

Causes the FTP object to produce `n` packets instantaneously.

`$ftp stop`

Causes the attached TCP object to stop sending data.

`$ftp attach agent`

Attaches a Application/FTP object to agent.

`$ftp producemore <count>`

Causes the Application/FTP object to produce `count` more packets.

Configuration Parameters are:

**maxpkts** The maximum number of packets generated by the source.

**TRACEFILE OBJECTS** Tracefile objects are used to specify the trace file that is to be used for generating traffic (see Application/Traffic/Trace objects described earlier in this section). `$tracefile` is an instance of the Tracefile Object. `$tracefile filename <trace-input>`

Set the filename from which the traffic trace data is to be read to `trace-input`.

There are no configuration parameters for this object. A trace file consists of any number of fixed length records. Each record consists of 2 32 bit fields. The first indicates the interval until the next packet is generated in microseconds. The second indicates the length of the next packet in bytes.

## 37.6 Commands at a glance

Following are some transport agent and application related commands commonly used in simulation scripts:

```
set tcpl [new Agent/TCP]
$ns_ attach-agent $node_(src) $tcpl
set sink1 [new Agent/TCPSink]
$ns_ attach-agent $node_(snk) $sink1
$ns_ connect $tcpl $sink1
```

This creates a source tcp agent and a destination sink agent. Both the transport agents are attached to their respective nodes. Finally an end-to-end connection is setup between the src and sink.

```
set ftp1 [new Application/FTP]
$ftp1 attach-agent $agent
```

or `set ftp1 [$agent attach-app FTP]` Both the above commands achieve the same. An application (ftp in this example) is created and attached to the source agent. An application can be of two types, a traffic generator or a simulated application. Types of Traffic generators currently present are: Application/Traffic/Exponential, Application/Traffic/Pareto, Application/Traffic/CBR, and Application/Traffic/Trace. See section 37.3 for details. Types of simulated applications currently implemented are: Application/FTP and Application/Telnet. See section 37.4 for details.

## Chapter 38

# Web cache as an application

All applications described above are “virtual” applications, in the sense that they do not actually transfer their own data in the simulator; all that matter is the *size* and the *time* when data are transferred. Sometimes we may want applications to transfer their own data in simulations. One such example is web caching, where we want HTTP servers to send HTTP headers to caches and clients. These headers contain page modification time information and other caching directives, which are important for some cache consistency algorithms.

In the following, we first describe general issues regarding transmitting application-level data in *ns*, then we discuss special issues, as well as APIs, related to transmitting application data using TCP as transport. We will then proceed to discuss the internal design of HTTP client, server, and proxy cache.

### 38.1 Using application-level data in *ns*

In order to transmit application-level data in *ns*, we provide a uniform structure to pass data among applications, and to pass data from applications to transport agents (Figure 38.1). It has three major components: a representation of a uniform application-level data unit (ADU), a common interface to pass data between applications, and two mechanisms to pass data between applications and transport agents.

#### 38.1.1 ADU

The functionality of an ADU is similar to that of a Packet. It needs to pack user data into an array, which is then included in the user data area of an *nspacket* by an Agent (this is not supported by current Agents. User must derive new agents to accept user data from applications, or use an wrapper like *TcpApp*. We’ll discuss this later).

Compared with Packet, ADU provides this functionality in a different way. In Packet, a common area is allocated for all packet headers; an offset is used to access different headers in this area. In ADU this is not applicable, because some ADU allocates their space dynamically according to the availability of user data. For example, if we want to deliver an OTcl script between applications, the size of the script is undetermined beforehand. Therefore, we choose a less efficient but more flexible method. Each ADU defines its own data members, and provides methods to serialize them (i.e., pack data into an array and extract them from an array). For example, in the abstract base class of all ADU, *AppData*, we have:

```
class AppData {
```

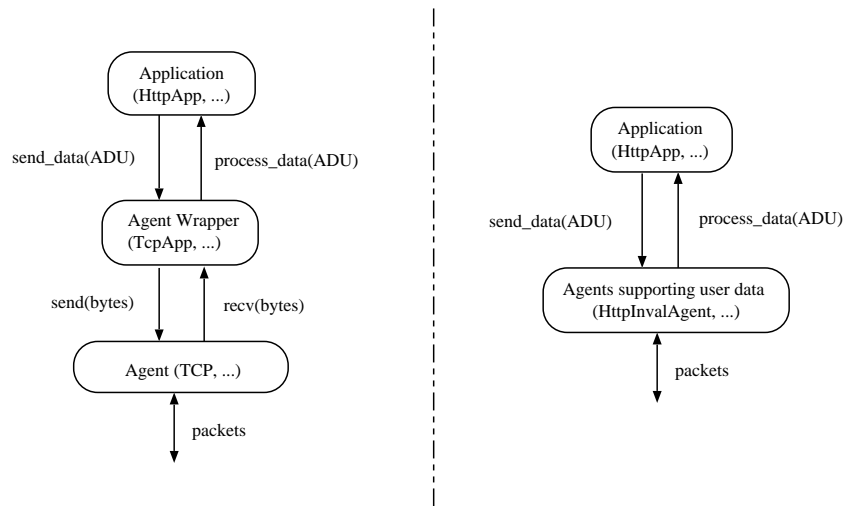


Figure 38.1: Examples of application-level data flow

```
private:
    AppDataType type_; // ADU type
public:
    struct hdr {
        AppDataType type_;
    };
public:
    AppData(char* b) {
        assert(b != NULL);
        type_ = ((hdr *)b)->type_;
    }
    virtual void pack(char* buf) const;
}
```

Here `pack(char* buf)` is used to write an `AppData` object into an array, and `AppData(char* b)` is used to build a new `AppData` from a “serialized” copy of the object in an array.

When deriving new ADU from the base class, users may add more data, but at the same time a new `pack(char* b)` and a new constructor should be provided to write and read those new data members from an array. For an example as how to derive from an ADU, look at `ns/webcache/http-aux.h`.

### 38.1.2 Passing data between applications

The base class of `Application`, `Process`, allows applications to pass data or request data between each other. It is defined as follows:

```
class Process {
public:
    Process() : target_(0) {}
    inline Process*& target() { return target_; }
```

```

        virtual void process_data(int size, char* data) = 0;
        virtual void send_data(int size, char* data) = 0;

protected:
        Process* target_;
};

```

Process enables Application to link together.

### 38.1.3 Transmitting user data over UDP

Currently there is no support in class Agent to transmit user data. There are two ways to transmit a serialized ADU through transport agents. First, for UDP agent (and all agents derived from there), we can derive from class UDP and add a new method `send(int nbytes, char *userdata)` to pass user data from Application to Agent. To pass data from an Agent to an Application is somewhat trickier: each agent has a pointer to its attached application, we dynamically cast this pointer to an `AppConnector` and then call `AppConnector::process_data()`.

As an example, we illustrate how class `HttpInvalAgent` is implemented. It is based on UDP, and is intended to deliver web cache invalidation messages (*ns/webcache/invalid-agent.h*). It is defined as:

```

class HttpInvalAgent : public Agent {
public:
        HttpInvalAgent();

        virtual void recv(Packet *, Handler *);
        virtual void send(int realsize, AppData* data);

protected:
        int off_inv_;
};

```

Here `recv(Packet*, Handler*)` overridden to extract user data, and a new `send(int, AppData*)` is provided to include user data in packets. An application (`HttpApp`) is attached to an `HttpInvalAgent` using `Agent::attachApp()` (a dynamic cast is needed). In `send()`, the following code is used to write user data from `AppData` to the user data area in a packet:

```

Packet *pkt = allocpkt(data->size());
hdr_inval *ih = (hdr_inval *)pkt->access(off_inv_);
ih->size() = data->size();
char *p = (char *)pkt->accessdata();
data->pack(p);

```

In `recv()`, the following code is used to read user data from packet and to deliver to the attached application:

```

hdr_inval *ih = (hdr_inval *)pkt->access(off_inv_);
((HttpApp*)app_)->process_data(ih->size(), (char *)pkt->accessdata());
Packet::free(pkt);

```

### 38.1.4 Transmitting user data over TCP

Transmitting user data using TCP is trickier than doing that over UDP, mainly because of TCP's reassembly queue is only available for FullTcp. We deal with this problem by abstracting a TCP connection as a FIFO pipe.

As indicated in section 37.2.4, transmission of application data can be implemented via agent upcalls. Assuming we are using TCP agents, all data are delivered in sequence, which means we can view the TCP connection as a FIFO pipe. We emulate user data transmission over TCP as follows. We first provide buffer for application data at the sender. Then we count the bytes received at the receiver. When the receiver has got all bytes of the current data transmission, it then gets the data directly from the sender. Class Application/TcpApp is used to implement this functionality.

A TcpApp object contains a pointer to a transport agent, presumably either a FullTcp or a SimpleTcp.<sup>1</sup> (Currently TcpApp doesn't support asymmetric TCP agents, i.e., sender is separated from receiver). It provides the following OTcl interfaces:

- `connect`: Connecting another TcpApp to this one. This connection is bi-directional, i.e., only one call to `connect` is needed, and data can be sent in either direction.
- `send`: It takes two arguments: `(nbytes, str)`. `nbytes` is the "nominal" size of application data. `str` is application data in string form.

In order to send application data in binary form, TcpApp provides a virtual C++ method `send(int nbytes, int dsize, const char *data)`. In fact, this is the method used to implement the OTcl method `send`. Because it's difficult to deal with binary data in Tcl, no OTcl interface is provided to handle binary data. `nbytes` is the number of bytes to be transmitted, `dsize` is the actual size of the array `data`.

TcpApp provides a C++ virtual method `process_data(int size, char*data)` to handle the received data. The default handling is to treat the data as a tcl script and evaluate the script. But it's easy to derive a class to provide other types of handling.

Here is an example of using Application/TcpApp. A similar example is `Test/TcpApp-2node` in `ns/tcl/test/test-suite-webcache.tcl`. First, we create FullTcp agents and connect them:

```
set tcp1 [new Agent/TCP/FullTcp]
set tcp2 [new Agent/TCP/FullTcp]
# Set TCP parameters here, e.g., window_, iss_, ...

$ns attach-agent $n1 $tcp1
$ns attach-agent $n2 $tcp2
$ns connect $tcp1 $tcp2
$tcp2 listen
```

Then we create TcpApps and connect them:

```
set app1 [new Application/TcpApp $tcp1]
set app2 [new Application/TcpApp $tcp2]
$app1 connect $app2
```

---

<sup>1</sup>A SimpleTcp agent is used solely for web caching simulations. It is actually an UDP agent. It has neither error recovery nor flow/congestion control. It doesn't do packet segmentation. Assuming a loss-free network and in-order packet delivery, SimpleTcp agent simplifies the trace files and hence aids the debugging of application protocols, which, in our case, is the web cache consistency protocol.

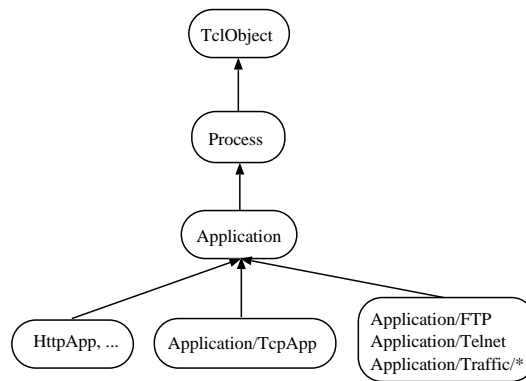


Figure 38.2: Hierarchy of classes related to application-level data handling

Now we let \$app1 be sender and \$app2 be receiver:

```
$ns at 1.0 "$app1 send 100 \"$app2 app-recv 100\""
```

Where app-recv is defined as:

```
Application/TcpApp instproc app-recv { size } {
    global ns
    puts "[$ns now] app2 receives data $size from app1"
}
```

### 38.1.5 Class hierarchy related to user data handling

We conclude this section by providing a hierarchy of classes involved in this section (Figure 38.2).

## 38.2 Overview of web cache classes

There are three major classes related to web cache, as it is in the real world: client (browser), server, and cache. Because they share a common feature, i.e., the HTTP protocol, they are derived from the same base class `Http` (Name of OTcl class, it's called `HttpApp` in C++). For the following reasons, it's not a real Application. First, an HTTP object (i.e., client/cache/server) may want to maintain multiple concurrent HTTP connections, but an Application contains only one `agent_`. Also, an HTTP object needs to transmit real data (e.g., HTTP header) and that's provided by `TcpApp` instead of any Agent. Therefore, we choose to use a standalone class derived from `Tclobj` for common features of all HTTP objects, which are managing HTTP connections and a set of pages. In the rest of the section, we'll discuss these functionalities of `Http`. In the next three sections, we'll in turn describe HTTP client, cache and server.

### 38.2.1 Managing HTTP connections

Every HTTP connection is embodied as a `TcpApp` object. `Http` maintains a hash of `TcpApp` objects, which are all of its active connections. It assumes that to any other `Http`, it has only one HTTP connection. It also allows dynamic establishment

and teardown of connections. Only OTcl interface is provided for establishing, tearing down a connection and sending data through a connection.

**OTcl methods** Following is a list of OTcl interfaces related to connection management in Http objects:

<code>id</code>	return the id of the Http object, which is the id of the node the object is attached to.
<code>get-cnc &lt;client&gt;</code>	return the TCP agent associated with \$client (Http object).
<code>is-connected &lt;server&gt;</code>	return 0 if not connected to \$server, 1 otherwise.
<code>send &lt;client&gt; &lt;bytes&gt; &lt;callback&gt;</code>	send \$bytes of data to \$client. When it's done, execute \$callback (a OTcl command).
<code>connect &lt;client&gt; &lt;TCP&gt;</code>	associate a TCP agent with \$client (Http object). That agent will be used to send packets to \$client.
<code>disconnect &lt;client&gt;</code>	delete the association of a TCP agent with \$client. Note that neither the TCP agent nor \$client is not deleted, only the association is deleted.

**Configuration parameter** By default, Http objects use Agent/SimpleTcp as transport agents (section 38.1.4). They can also use Agent/FullTcp agents, which allows Http objects to operate in a lossy network. Class variable `codeTRANSPORT_` is used for this purpose. E.g., `Http set TRANSPORT_FullTcp` tells all Http objects use FullTcp agents.

This configuration should be done *before* simulation starts, and it should not change during simulation, because FullTcp agents do not inter-operate with SimpleTcp agents.

### 38.2.2 Managing web pages

Http also provides OTcl interfaces to manage a set of pages. The real management of pages are handled by class `PagePool` and its subclasses. Because different HTTP objects have different requirements for page management, we allow different `PagePool` subclasses to be attached to different subclasses of `Http` class. Meanwhile, we export a common set of `PagePool` interfaces to OTcl through `Http`. For example, a browser may use a `PagePool` only to generate a request stream, so its `PagePool` only needs to contain a list of URLs. But a cache may want to store page size, last modification time of every page instead of a list of URLs. However, this separation is not clearcut in the current implementation.

Page URLs are represented in the form of: `<ServerName> : <SequenceNumber>` where the `ServerName` is the name of OTcl object, and every page in every server should have a unique `SequenceNumber`. Page contents are ignored. Instead, every page contains several *attributes*, which are represented in OTcl as a list of the following (`<name> <value>`) pairs: “`mod-time <val>`” (page modification time), “`size <val>`” (page size), and “`age <val>`”} The ordering of these pairs is not significant.

Following is a list of related OTcl methods.

<code>set-pagepool &lt;pagepool&gt;</code>	set page pool
<code>enter-page &lt;pageid&gt; &lt;attributes&gt;</code>	add a page with id \$pageid into pool. \$attributes is the attributes of \$pageid, as described above.
<code>get-page &lt;pageid&gt;</code>	return page attributes in the format described above.
<code>get-modtime &lt;pageid&gt;</code>	return the last modification time of the page \$pageid.
<code>exist-page &lt;pageid&gt;</code>	return 0 if \$pageid doesn't exist in this Http object, 1 otherwise.
<code>get-size &lt;pageid&gt;</code>	return the size of \$pageid.
<code>get-cachetime &lt;pageid&gt;</code>	return the time when page \$pageid is entered into the cache.



### 38.2.3 Debugging

HttpApp provides two debugging methods. `log` registers a file handle as the trace file for all HttpApp-specific traces. Its trace format is described in section 38.9. `evTrace` logs a particular event into trace file. It concatenates time and the id of the HttpApp to the given string, and writes it out. Details can be found in *ns/webcache/http.cc*.

## 38.3 Representing web pages

We represent web pages as the abstract class `Page`. It is defined as follows:

```
class Page {
public:
    Page(int size) : size_(size) {}
    int size() const { return size_; }
    int& id() { return id_; }
    virtual WebPageType type() const = 0;

protected:
    int size_;
    int id_;
};
```

It represents the basic properties of a web page: size and URL. Upon it we derive two classes of web pages: `ServerPage` and `ClientPage`. The former contains a list of page modification times, and is supposed to be used by servers. It was originally designed to work with a special web server trace; currently it is not widely used in *ns*. The latter, `ClientPage`, is the default web page for all page pools below.

A `ClientPage` has the following major properties (we omit some variables used by web cache with invalidation, which has too many details to be covered here):

- `HttpApp* server_` - Pointer to the original server of this page.
- `double age_` - Lifetime of the page.
- `int status_` - Status of the page. Its contents are explained below.

The status (32-bit) of a `ClientPage` is separated into two 16-bit parts. The first part (with mask `0x00FF`) is used to store page status, the second part (with mask `0xFF00`) is used to store expected page actions to be performed by cache. Available page status are (again, we omit those closely related to web cache invalidation):

`HTTP_VALID_PAGE`    Page is valid.

`HTTP_UNCACHEABLE`    Page is uncacheable. This option can be used to simulate CGI pages or dynamic server pages.

`ClientPage` has the following major C++ methods:

- `type()` - Returns the type of the page. Assuming pages of the same type should have identical operations, we let all `ClientPage` to be of type "HTML". If later on other types of web pages are needed, a class may be derived from `ClientPage` (or `Page`) with the desired type.

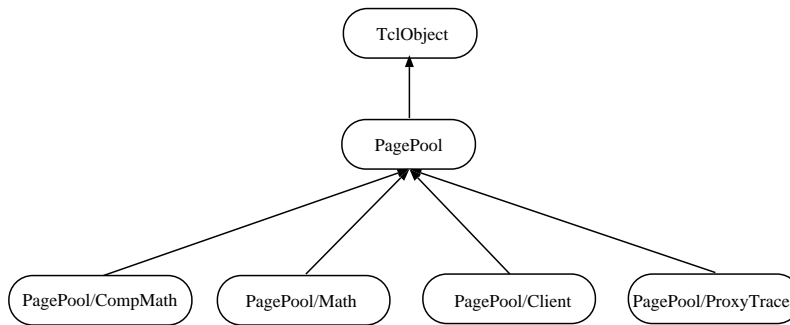


Figure 38.3: Class hierarchy of page pools

- `name(char *buf)` - Print the page's name into the given buffer. A page's name is in the format of: `<ServerName>:<PageID>`.
- `split_name(const char *name, PageID& id)` - Split a given page name into its two components. This is a static method.
- `mtime()` - Returns the last modification time of the page.
- `age()` - Returns the lifetime of the page.

## 38.4 Page pools

PagePool and its derived classes are used by servers to generate page information (name, size, modification time, lifetime, etc.), by caches to describe which pages are in storage, and by clients to generate a request stream. Figure 38.3 provides an overview of the class hierarchy here.

Among these, class PagePool/Client is mostly used by caches to store pages and other cache-related information; other three classes are used by servers and clients. In the following we describe these classes one by one.

### 38.4.1 PagePool/Math

This is the simplest type of page pool. It has only one page, whose size can be generated by a given random variable. Page modification sequence and request sequence are generated using two given random variables. It has the following OTcl methods:

<code>gen-pageid</code>	Returns the page ID which will be requested next. Because it has only one page, it always returns 0.
<code>gen-size</code>	Returns the size of the page. It can be generated by a given random variable.
<code>gen-modtime &lt;pageID&gt; &lt;mt&gt;</code>	Returns the next modification time of the page. <code>&lt;mt&gt;</code> gives the last modification time. It uses the lifetime random variable.
<code>ranvar-age &lt;rv&gt;</code>	Set the file lifetime random variable as <code>&lt;rv&gt;</code> .
<code>ranvar-size &lt;rv&gt;</code>	Set the file size random variable to be <code>&lt;rv&gt;</code> .

**NOTE:** There are two ways to generate a request sequence. With all page pools except PagePool/ProxyTrace, request sequence is generated with a random variable which describes the request interval, and the `gen-pageid` method of other page pools

gives the page ID of the next request. PagePool/ProxyTrace loads the request stream during initialization phase, so it does not need a random variable for request interval; see its description below.

An example of using PagePool/Math is at Section 38.8. That script is also available at *ns/tcl/ex/simple-webcache.tcl*.

### 38.4.2 PagePool/CompMath

It improves over PagePool/Math by introducing a compound page model. By a compound page we mean a page which consists of a main text page and a number of embedded objects, e.g., GIFs. We model a compound page as a main page and several component objects. The main page is always assigned with ID 0. All component pages have the same size; both the main page size and component object size is fixed, but adjustable through OTcl-bound variables *main\_size\_* and *comp\_size\_*, respectively. The number of component objects can be set using the OTcl-bound variable *num\_pages\_*.

PagePool/CompMath has the following major OTcl methods:

<code>gen-size &lt;pageID&gt;</code>	If <pageID> is 0, return <i>main_size_</i> , otherwise return <i>comp_size_</i> .
<code>ranvar-main-age &lt;rv&gt;</code>	Set random variable for main page lifetime. Another one, <code>ranvar-obj-age</code> , set that for component objects.
<code>gen-pageid</code>	Always returns 0, which is the main page ID.
<code>gen-modtime &lt;pageID&gt; &lt;mt&gt;</code>	Returns the next modification time of the given page <pageID>. If the given ID is 0, it uses the main page lifetime random variable; otherwise it uses the component object lifetime random variable.

An example of using PagePool/CompMath is available at *ns/tcl/ex/simple-webcache-comp.tcl*.

### 38.4.3 PagePool/ProxyTrace

The above two page pool synthesize request stream to a single web page by two random variables: one for request interval, another for requested page ID. Sometimes users may want more complicated request stream, which consists of multiple pages and exhibits spatial locality and temporal locality. There exists one proposal (SURGE [3]) which generates such request streams, we choose to provide an alternative solution: use real web proxy cache trace (or server trace).

The class PagePool/ProxyTrace uses real traces to drive simulation. Because there exist many web traces with different formats, they should be converted into a intermediate format before fed into this page pool. The converter is available at <http://mash.cs.berkeley.edu/dist/vint/webcache-trace-conv.tar.gz>. It accepts four trace formats: DEC proxy trace (1996), UCB Home-IP trace, NLANR proxy trace, and EPA web server trace. It converts a given trace into two files: *pglog* and *reqlog*. Each line in *pglog* has the following format:

```
[ <serverID> <URL_ID> <PageSize> <AccessCount> ]
```

Each line, except the last line, in *reqlog* has the following format:

```
[ <time> <clientID> <serverID> <URL_ID> ]
```

The last line in *reqlog* records the duration of the entire trace and the total number of unique URLs:

```
i <Duration> <Number_of_URL>
```

PagePool/ProxyTrace takes these two file as input, and use them to drive simulation. Because most existing web proxy traces do not contain complete page modification information, we choose to use a bimodal page modification model [7]. We allow user to select  $x\%$  of the pages to have one random page modification interval generator, and the rest of the pages to have another generator. In this way, it's possible to let  $x\%$  pages to be dynamic, i.e., modified frequently, and the rest static. Hot pages are evenly distributed among all pages. For example, assume 10% pages are dynamic, then if we sort pages into a list according to their popularity, then pages 0, 10, 20, ... are dynamic, rest are static. Because of this selection mechanism, we only allow bimodal ratio to change in the unit of 10%.

In order to distribute requests to different requestors in the simulator, PagePool/ProxyTrace maps the client ID in the traces to requestors in the simulator using a modulo operation.

PagePool/ProxyTrace has the following major OTcl methods:

get-poolsize	Returns the total number of pages.
get-duration	Returns the duration of the trace.
bimodal-ratio	Returns the bimodal ratio.
set-client-num <num>	Set the number of requestors in the simulation.
gen-request <ClientID>	Generate the next request for the given requestor.
gen-size <PageID>	Returns the size of the given page.
bimodal-ratio <ratio>	Set the dynamic pages to be <ratio>*10 percent. Note that this ratio changes in unit of 10%.
ranvar-dp <ranvar>	Set page modification interval generator for dynamic pages. Similarly, ranvar-sp <ranvar> sets the generator for static pages.
set-reqfile <file>	Set request stream file, as discussed above.
set-pgfile <file>	Set page information file, as discussed above.
gen-modtime <PageID> <LastModTime>	Generate next modification time for the given page.

An example of using PagePool/ProxyTrace is available at *ns/tcl/ex/simple-webcache-trace.tcl*.

#### 38.4.4 PagePool/Client

The class PagePool/Client helps caches to keep track of pages resident in cache, and to store various cache-related information about pages. It is mostly implemented in C++, because it is mainly used internally and little functionality is needed by users. It has the following major C++ methods:

- `get_page(const char* name)` - Returns a pointer to the page with the given name.
- `add_page(const char *name, int size, double mt, double et, double age)` - Add a page with given size, last modification time (mt), cache entry time (et), and page lifetime (age).
- `remove_page(const char* name)` - Remove a page from cache.

This page pool should support various cache replacement algorithms, however, it has not been implemented yet.

#### 38.4.5 PagePool/WebTraf

The class PagePool/WebTraf is a standalone Web traffic module that utilizes PagePool framework. However, this class has nothing to do with the HttpApp classes. Because we are only interested in using it to study Web traffic pattern here, and do

not want to be bothered with the burden of transmitting HTTP headers, etc. It has the following two major data structures. Details can be found in ns/webcache/webtraf.cc and ns/webcache/webtraf.h, the architecture WebTraf model is also loosely described in [10], Section 2.4, paragraph 3-4 and the appendix A.1.

- WebTrafSession - a class that models Web user session. It is defined as follows:

```
class WebTrafSession : public TimerHandler {
public:
    WebTrafSession(WebTrafPool *mgr, Node *src, int np, int id) : rvInterPage_(NULL),
        rvPageSize_(NULL), rvInterObj_(NULL), rvObjSize_(NULL), mgr_(mgr), src_(src),
        nPage_(np), curPage_(0), donePage_(0), id_(id), interPageOption_(1) {}
    virtual ~WebTrafSession();

    // Queried by individual pages/objects
    inline RandomVariable*& interPage() { return rvInterPage_; }
    inline RandomVariable*& pageSize() { return rvPageSize_; }
    inline RandomVariable*& interObj() { return rvInterObj_; }
    inline RandomVariable*& objSize() { return rvObjSize_; }

    void donePage(void* ClntData); // all the pages within this
                                // session have been sent
    void launchReq(void* ClntData, int obj, int size);
    inline int id() const { return id_; }
    inline WebTrafPool* mgr() { return mgr_; }
private:
    virtual void expire(Event *e = 0); // Launch request for a page
    virtual void handle(Event *e); // schedule the timer for next page

    RandomVariable *rvInterPage_, *rvPageSize_, *rvInterObj_, *rvObjSize_;
    WebTrafPool* mgr_;
    Node* src_; // One Web client (source of request) per session
    int nPage_; // number of pages per session
    int curPage_; // number of pages that have been sent
    int id_; // page ID
    int interPageOption_;
}
```

- WebPage - a class that models Web Page. It is defined as follows:

```
class WebPage : public TimerHandler {
public:
    WebPage(int id, WebTrafSession* sess, int nObj, Node* dst) :
        id_(id), sess_(sess), nObj_(nObj), curObj_(0),
        doneObj_(0), dst_(dst) {}
    virtual ~WebPage() {}
    inline void start() { // Call expire() and schedule the next one if needed
    void doneObject() { // All the objects within this page have been sent
    inline int id() const { return id_; }
    Node* dst() { return dst_; }
    inline int curObj() const { return curObj_; }
    inline int doneObj() const { return doneObj_; }
private:
    virtual void expire(Event* = 0) { // Launch request for an object
    virtual void handle(Event *e) { // schedule the timer for the next object
```

```

    int id_; // object ID
    WebTrafSession* sess_; // the session that requested this page
    int nObj_; // number of object in this page
    int curObj_ ; // number of object that have been sent
    Node* dst_; // server that this page has been requested from
}

```

Following is a list of related OTcl methods to the WebTraf class.

set-num-session <number-of-session>	set the total number of sessions in the WebTraf pool.
set-num-server <number-of-server>	set the total number of servers.
set-num-client <number-of-client>	set the total number clients.
set-interPageOption <option>	There are two ways to interpret <i>inter-page</i> time: One is the time between the start of two consecutive page downloads by the same user, and the other is the time between the end of previous page download and the start of the following page by the same user. \$option can be set to either 0 or 1 (default is 1). When \$option is set to 1, the second interpretation is used for "inter-page" time. The first interpretation is adopted when \$option is set to 0. Note the resulted traffic volume using the first interpretation is much higher than the second interpretation.
doneObj <webpage>	all the objects in \$webpage have been sent.
set-server <id> <node>	set \$node as server \$id.
set-client <id> <node>	set \$node as client \$id.
recycle <tcp> <sink>	Recycle a TCP source/sink pair.
create-session <session-index> <pages-per-sess> <launch-time> <inter-page-rv> <page-size-rv> <inter-obj-rv> <obj-size-rv>	Create a Web session. \$session-index is the session index. \$pages-per-sess is the total number of pages per session. \$launch-time is session starting time. \$inter-page-rv is the random variable that generates page inter-arrival time. \$page-size-rv is the random variable that generates number of objects per page. \$inter-obj-rv is the random variable that generates object inter-arrival time. \$obj-size-rv is the random variable that generates object size.

The example script is available at ns/tcl/ex/web-traffic.tcl (also see ns/tcl/ex/large-scale-web-traffic.tcl for use of a large-scale web traffic simulation)

## 38.5 Web client

Class Http/Client models behavior of a simple web browser. It generates a sequence of page requests, where request interval and page IDs are randomized. It's a pure OTcl class inherited from Http. Next we'll walk through its functionalities and usage.

**Creating a client** First of all, we create a client and connect it to a cache and a web server. Currently a client is only allowed to connect to a single cache, but it's allowed to connect to multiple servers. Note that this has to be called *AFTER* the simulation starts (i.e., after `$ns run` is called). This remains true for all of the following methods and code examples of Http and its derived classes, unless explicitly said.

```
# Assuming $server is a configured Http/Server.
set client [new Http/Client $ns $node]           ; # client resides on this node
$client connect $server                          ; # connecting client to server
```

**Configuring request generation** For every request, Http/Client uses PagePool to generate a random page ID, and use a random variable to generate intervals between two consecutive requests: <sup>2</sup>

```
$client set-page-generator $pgp                  ; # attach a configured PagePool
$client set-interval-generator $ranvar           ; # attach a random variable
```

Here we assume that PagePools of Http/Client share the same set of pages as PagePools of the server. Usually we simplify our simulation by letting all clients and servers share the same PagePool, i.e., they have the same set of pages. When there are multiple servers, or servers' PagePools are separated from those of clients', care must be taken to make sure that every client sees the same set of pages as the servers to which they are attached.

**Starting** After the above setup, starting requests is very simple:

```
$client start-session $cache $server             ; # assuming $cache is a configured Http/Cache
```

**OTcl interfaces** Following is a list of its OTcl methods (in addition to those inherited from Http). This is not a complete list. More details can be found in *ns/tcl/webcache/http-agent.tcl*.

send-request <server> <type> <pageid> <args>	send a request of page \$pageid and type \$type to \$server. The only request type allowed for a client is GET. \$args has a format identical to that of \$attributes described in <code>Http::enter-page</code> .
start-session <cache> <server>	start sending requests of a random page to \$server via \$cache.
start <cache> <server>	before sending requests, populate \$cache with all pages in the client's PagePool. This method is useful when assuming infinite-sized caches and we want to observe behaviors of cache consistency algorithms in steady state.
set-page-generator <pagepool>	attach a PagePool to generate random page IDs.
set-interval-generator <ranvar>	attach a random variable to generate random request intervals.

## 38.6 Web server

Class Http/Server models behavior of a HTTP server. Its configuration is very simple. All that a user needs to do is to create a server, attach a PagePool and wait:

```
set server [new Http/Server $ns $node]          ; # attach $server to $node
$server set-page-generator $pgp                 ; # attach a page pool
```

---

<sup>2</sup>Some PagePool, e.g., PagePool/Math, has only one page and therefore it always returns the same page. Some other PagePool, e.g., PagePool/Trace, has multiple pages and needs a random variable to pick out a random page.

An Http/Server object waits for incoming requests after simulation starts. Usually clients and caches initiate connection to an Http/Server. But it still has its own `connect` method, which allows an Http/Server object to actively connect to a certain cache (or client). Sometimes this is useful, as explained in `Test/TLC1::set-groups{ }` in `ns/tcl/test/test-suite-webcache.tcl`.

An Http/Server object accepts two types of requests: GET and IMS. GET request models normal client requests. For every GET request, it returns the attributes of the requested page. IMS request models If-Modified-Since used by TTL algorithms for cache consistency. For every IMS (If-Modified-Since) request, it compares the page modification time given in the request and that of the page in its PagePool. If the time indicated in the request is older, it sends back a response with very small size, otherwise it returns all of the page attributes with response size equal the real page size.

## 38.7 Web cache

Currently 6 types of web caches are implemented, including the base class Http/Cache. Its five derived subclasses implement 5 types of cache consistency algorithms: Plain old TTL, adaptive TTL, Omniscient TTL, Hierarchical multicast invalidation, and hierarchical multicast invalidation plus direct request.

In the following we'll only describe the base class Http/Cache, because all the subclasses involve discussion of cache consistency algorithms and it does not seem to be appropriate here.

### 38.7.1 Http/Cache

Class Http/Cache models behavior of a simple HTTP cache with infinite size. It doesn't contain removal algorithm, nor consistency algorithm. It is not intended to be used by itself. Rather, it is meant to be a base class for experimenting with various cache consistency algorithms and other cache algorithms.

**Creation and startup** Creating an Http/Cache requires the same set of parameters as Http/Client and Http/Server. After creation, a cache needs to connect to a certain server. Note that this creation can also be done dynamically, when a request comes in and the cache finds that it's not connected to the server. However, we do not model this behavior in current code. Following code is an example:

```
set cache [new HttpCache $ns $node]           ; # attach cache to $node
$cache connect $server                        ; # connect to $server
```

Like Http/Server, an Http/Cache object waits for requests (and packets from server) after it's initialized as above. When hierarchical caching is used, the following can be used to create the hierarchy:

```
$cache set-parent $parent                    ; # set parent cache
```

Currently all TTL and multicast invalidation caches support hierarchical caching. However, only the two multicast invalidation caches allow multiple cache hierarchies to inter-operate.

**OTcl methods** Although Http/Cache is a SplitObject, all of its methods are in OTcl. Most of them are used to process an incoming request. Their relations can be illustrated with the flowchart below, followed by explanations:



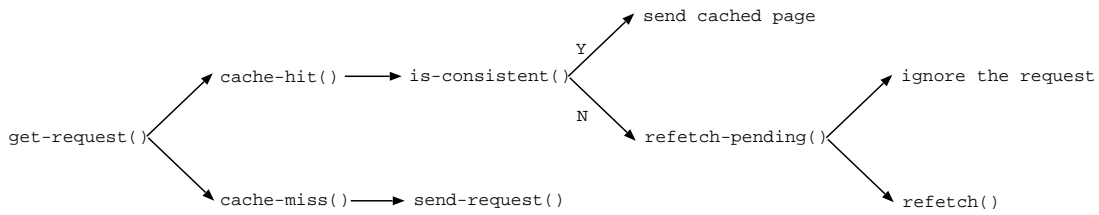


Figure 38.4: Handling of incoming request in Http/Cache

get-request <client> <type> <pageid>	The entry point of processing any request. It checks if the requested page \$pageid exists in the cache's page pool, then call either cache-hit or cache-miss.
cache-miss <client> <type> <pageid>	This cache doesn't have the page. Send a request to server (or parent cache) to refetch the page if it hasn't already done so. Register \$client in a list so that when the cache gets the page, it'll forward the page to all clients who have requested the page.
cache-hit <client> <type> <pageid>	Checks the validity of the cached page. If it's valid, send \$client the cached page, otherwise refetch the page.
is-consistent <client> <type> <pageid>	Returns 1 if \$pageid is valid. This is intended to be overridden by subclasses.
refetch <client> <type> <pageid>	Refetch an invalid page from server. This is intended to be overridden by subclasses.

## 38.8 Putting together: a simple example

We have seen all the pieces, now we present a script which provides a complete view of all pieces together. First, we build topology and other usual initializations:

```

set ns [new Simulator]

# Create topology/routing
set node(c) [$ns node]
set node(e) [$ns node]
set node(s) [$ns node]
$ns duplex-link $node(s) $node(e) 1.5Mb 50ms DropTail
$ns duplex-link $node(e) $node(c) 10Mb 2ms DropTail
$ns rtproto Session

```

Next we create the Http objects:

```

# HTTP logs
set log [open "http.log" w]

# Create page pool as a central page generator. Use PagePool/Math
set pgp [new PagePool/Math]
set tmp [new RandomVariable/Constant]
$tmp set val_ 1024
$pgp ranvar-size $tmp

```

; ## Page size generator  
; ## average page size

```

set tmp [new RandomVariable/Exponential]                ; ## Page age generator
$tmp set avg_ 5                                         ; ## average page age
$pgp ranvar-age $tmp

set server [new Http/Server $ns $node(s)]              ; ## Create a server and link it to the cen-
tral page pool
$server set-page-generator $pgp
$server log $log

set cache [new Http/Cache $ns $node(e)]                ; ## Create a cache
$cache log $log

set client [new Http/Client $ns $node(c)]              ; ## Create a client
set tmp [new RandomVariable/Exponential]              ; ## Poisson process as request sequence
$tmp set avg_ 5                                         ; ## average request interval
$client set-interval-generator $tmp
$client set-page-generator $pgp
$client log $log

set startTime 1                                         ; ## simulation start time
set finishTime 50                                       ; ## simulation end time
$ns at $startTime "start-connection"
$ns at $finishTime "finish"

```

Then we define a procedure which will be called after simulation starts. The procedure will setup connections among all Http objects.

```

proc start-connection {} {
    global ns server cache client
    $client connect $cache
    $cache connect $server
    $client start-session $cache $server
}

```

At the end, the usual closing:

```

proc finish {} {
    global ns log
    $ns flush-trace
    flush $log
    close $log
    exit 0
}
$ns run

```

This script is also available at `ns/tcl/ex/simple-webcache.tcl`. Examining its output `http.log`, one will find that the result of the absence cache consistency algorithm results in a lot of stale hits. This can be easily remedied by replacing “new Http/Cache” line with: `set cache [new Http/Cache/TTL $ns $node(e)]`. For more complicated cache consistency algorithm examples, see `ns/tcl/test/test-suite-webcache.tcl`.

## 38.9 Http trace format

The trace file of Http agents are constructed in a similar way as the SRM trace files. It consists of multiple entries, each of which occupies one line. The format of each entry is:

Time | ObjectID | Object Values

There are three types of objects: client (**C**), cache (**E**) and server (**S**). Following is a complete enumeration of all possible events and value types associated with these three types of objects.

<i>Object Type</i>	<i>Event Type</i>	<i>Values</i>
E	HIT	⟨Prefi x⟩
E	MISS	⟨Prefi x⟩ z ⟨RequestSize⟩
E	IMS	⟨Prefi x⟩ z ⟨Size⟩ t ⟨CacheEntryTime⟩
E	REF	p ⟨PageID⟩ s ⟨ServerID⟩ z ⟨Size⟩
E	UPD	p ⟨PageID⟩ m ⟨LastModifi edTime⟩ z ⟨PageSize⟩ s ⟨ServerID⟩
E	GUPD	z ⟨PageSize⟩
E	SINV	p ⟨PageID⟩ m ⟨LastModTime⟩ z ⟨PageSize⟩
E	GINV	p ⟨PageID⟩ m ⟨LastModTime⟩
E	SPF	p ⟨PageID⟩ c ⟨DestCache⟩
E	RPF	p ⟨PageID⟩ c ⟨SrcCache⟩
E	ENT	p ⟨PageID⟩ m ⟨LastModifi edTime⟩ z ⟨PageSize⟩ s ⟨ServerID⟩
C	GET	p ⟨PageID⟩ s ⟨PageServerID⟩ z ⟨RequestSize⟩
C	STA	p ⟨PageID⟩ s ⟨OrigServerID⟩ l ⟨StaleTime⟩
C	RCV	p ⟨PageID⟩ s ⟨PageServerID⟩ l ⟨ResponseTime⟩ z ⟨PageSize⟩
S	INV	p ⟨PageID⟩ m ⟨LastModifi edTime⟩ z ⟨Size⟩
S	UPD	p ⟨PageID⟩ m ⟨LastModifi edTime⟩ z ⟨Size⟩
S	SND	p ⟨PageID⟩ m ⟨LastModifi edTime⟩ z ⟨PageSize⟩ t ⟨Requesttype⟩
S	MOD	p ⟨PageID⟩ n ⟨NextModifyTime⟩

⟨Prefi x⟩ is the information common to all trace entries. It includes:

p ⟨PageID⟩ | c ⟨RequestClientID⟩ | s ⟨PageServerID⟩

*Short Explanation of event operations:*

<i>Object Type</i>	<i>Event Type</i>	<i>Explanation</i>
E	HIT	Cache hit. PageSererID is the id of the “owner” of the page.
E	MISS	Cache miss. In this case the cache will send a request to the server to fetch the page.
E	IMS	If-Modifi ed-Since. Used by TTL procotols to validate an expired page.
E	REF	Page refetch. Used by invalidation protocols to refetch an invalidated page.
E	UPD	Page update. Used by invalidation protocols to “push” updates from parent cache to children caches.
E	SINV	Send invalidation.
E	GINV	Get invalidation.
E	SPF	Send a pro forma
E	RPF	Receive a pro forma
E	ENT	Enter a page into local page cache.
C	GET	Client sends a request for a page.
C	STA	Client gets a stale hit. OrigModTime is the modifi cation time in the web server, CurrModTime is the local page’s modifi cation time.
C	RCV	Client receives a page.
S	SND	Server send a response.
S	UPD	Server pushes a page update to its “primary cache”. Used by invalidation protocol only.
S	INV	Server sends an invalidation message. Used by invalidation protocol only.
S	MOD	Server modifi ed a page. The page will be modifi ed next at <NextModifyTime>.

## 38.10 Commands at a glance

Following are the web cache related commands:

```
set server [new Http/Server <sim> <s-node>]
```

This creates an instance of an Http server at the specifi ed <s-node>. An instance of the simulator <sim> needs to be passed as an argument.

```
set client [new Http/Client <sim> <c-node>]
```

This creates an instance of a Http client at the given <c-node>.

```
set cache [new Http/Cache <sim> <e-node>]
```

This command creates a cache.

```
set pgp [new PagePool/<type-of-pagepool>]
```

This creates a pagepool of the type specifi ed. The different types of pagepool currently implemented are:

PagePool/Math, PagePool/CompMath, PagePool/ProxyTrace and PagePool/Client. See section 38.4 for details on Otcl interface for each type of Pagepool.

```
$server set-page-generator <pgp>
```

```
$server log <handle-to-log-file>
```

The above commands consist of server confi guration. First the server is attached to a central page pool <pgp>. Next it is attached to a log fi le.

```
client set-page-generator <pgp>
```

```
$client set-interval-generator <ranvar>
```

```
$client log <handle-to-log-file>
```

These consist configuration of the Http client. It is attached to a central page pool <pgp>. Next a random variable <ranvar> is attached to the client that is used by it (client) to generate intervals between two consecutive requests. Lastly the client is attached to a log file for logging its events.

```
$cache log <log-file>
```

This is part of cache configuration that allows the cache to log its events in a log-file.

```
$client connect <cache>
```

```
$cache connect <server>
```

Once the client, cache, and server are configured, they need to be connected as shown in above commands.

```
$client start-session <cache> <server>
```

This starts sending request for a random page from the client to the <server> via <cache>.

## Chapter 39

# Worm Model

In this chapter, we describe a scalable worm propagation model in *ns*, namely the detailed-network and abstract-network (DN-AN) model. It combines packet-level simulations with analytic worm spreading model. As shown in Figure 39.1, we model the Internet with two parts: detailed, and abstract part. A detailed-network could be an enterprise-network or the network run by an ISP. It simulates network connectivity and packet transmission. Users can evaluate worm detection algorithms in the detailed network. On the other hand, we abstract the rest of the Internet with a mathematical model, namely susceptible-infectious-removal (SIR) model (refer to [13] for detailed descriptions). Compared to the detailed network, we only track several state variables in the abstract world, such as the number of infected hosts. The interaction between DN and AN is through actual packet transmissions, that is, the probing traffic generated by compromised hosts in both parts.

For detailed description on DN-AN model, please refer to our draft paper. We implement the worm propagation model as applications. The source code can be found at `~ns//apps/worm.{cc,h}`. There is also a sample script to illustrate the DN-AN model under `~ns//tcl/ex/worm.tcl`.

### 39.1 Overview

We implement the worm propagation model with three classes: `class WormApp`, `DnhWormApp`, and `AnWormApp`. `class WormApp` and `class DnhWormApp` are used in the detailed network, representing invulnerable and vulnerable hosts respectively. `class AnWormApp` is the abstract network. Currently, our model only supports UDP-based worms.

An vulnerable host is compromised upon receiving a probing packet. Then, it chooses a target host (randomly or with certain preference to local neighbors) to scan. Probing packets have no effect on invulnerable hosts. When the abstract network receives probing packets, it updates its current states.

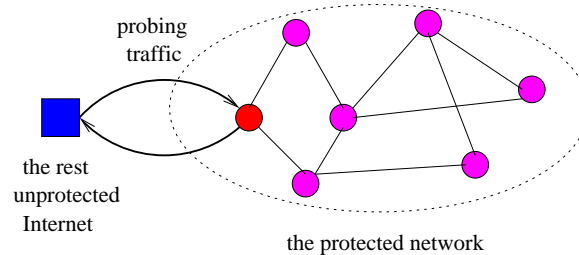


Figure 39.1: The DN-AN model.

## 39.2 Configuration

To set up simulation scenario, we first build the detailed network. We also need to create one extra node to represent the abstract network, and connect it to the detailed network.

For nodes in the detailed network, we first attach a `MessagePassing` agent to each node:

```
set a [new Agent/MessagePassing]
$n attach $a $probing_port
```

If the node represents a vulnerable host, we use `class DnhWormApp`

```
set w [new Application/Worm/Dnh]
$w attach-agent $a
```

Otherwise, we configure the node as invulnerable:

```
set w [new Application/Worm]
$w attach-agent $a
```

We configure the abstract network as:

```
set a [new Agent/MessagePassing]
$na attach $a $probing_port
set w [new Application/Worm/An]
$w attach-agent $a
```

In order for the abstract network to receive probing packets generated by nodes within the detailed networks, we need to use manual routing. There are some extra configuration for the abstract-network node:

```
set p [$na set dmux_]
$p defaulttarget $a
[$na entry] defaulttarget $p
```

## 39.3 Commands at a glance

Some common parameters can be configured through TCL script:

```
ScanRate      # the rate that a compromised host sends probing packets
ScanPort      # the vulnerable service port number
ScanPacketSize # the size of worm probing packets
```

By default, compromised hosts scan the Internet randomly. We can also simulate local-scanning worm by setting the local-scanning probability:

```
$w local-p 0.5
```

Following are some commands to configure parameters for the abstract network:

```
$w beta 0.1    # infection parameter
$w gamma 0     # removal parameter
$w addr-range 2000 200000 # the address space of the abstract network
$w dn-range 0 1999       # the address space of the detailed network
$w v_percent 0.01        # the percentage of vulnerable hosts in the abstract network
```



## Chapter 40

# PackMime-HTTP: Web Traffic Generation in NS-2

The PackMime Internet traffic model was developed by researchers in the Internet Traffic Research group at Bell Labs, based on recent Internet traffic traces. PackMime includes a model of HTTP traffic, called PackMime-HTTP. The traffic intensity generated by PackMime-HTTP is controlled by the rate parameter, which is the average number of new connections started each second. The PackMime-HTTP implementation in ns-2, developed at UNC-Chapel Hill, is capable of generating HTTP/1.0 and HTTP/1.1 (persistent, non-pipelined) connections.

PackMime-HTTP in ns-2 uses DelayBox (see Chapter 22), a module developed at UNC-Chapel Hill for delaying and/or dropping packets in a flow according to a distribution. DelayBox can be used with PackMime-HTTP to simulate a "cloud" of clients and servers that have different round-trip times, bottleneck link speeds, and amounts of packet loss (see 40.3 for more information).

The PackMime HTTP traffic model is described in detail in the following paper: J. Cao, W.S. Cleveland, Y. Gao, K. Jeffay, F.D. Smith, and M.C. Weigle, "Stochastic Models for Generating Synthetic HTTP Source Traffic", *Proceedings of IEEE INFOCOM*, Hong Kong, March 2004.

### 40.1 Implementation Details

PackMimeHTTP is an ns object that drives the generation of HTTP traffic. Each PackMimeHTTP object controls the operation of two types of Applications, a PackMimeHTTP server Application and a PackMimeHTTP client Application. Each of these Applications is connected to a TCP Agent (Full-TCP). **Note:** PackMime-HTTP only supports Full-TCP agents.

Each web server or web client cloud is represented by a single ns node that can produce and consume multiple HTTP connections at a time (Figure 40.1). For each HTTP connection, PackMimeHTTP creates (or allocates from the inactive pool, as described below) server and client Applications and their associated TCP Agents. After setting up and starting each connection, PackMimeHTTP sets a timer to expire when the next new connection should begin. The time between new connections is governed by the connection rate parameter supplied by the user. New connections are started according to the connection arrival times without regard to the completion of previous requests, but a new request between the same client and server pair (as with HTTP 1.1) begins only after the previous request-response pair has been completed.

PackMimeHTTP handles the re-use of Applications and Agents that have completed their data transfer. There are 5 pools used to maintain Applications and Agents – one pool for inactive TCP Agents and one pool each for active and inactive client and

# PackMime

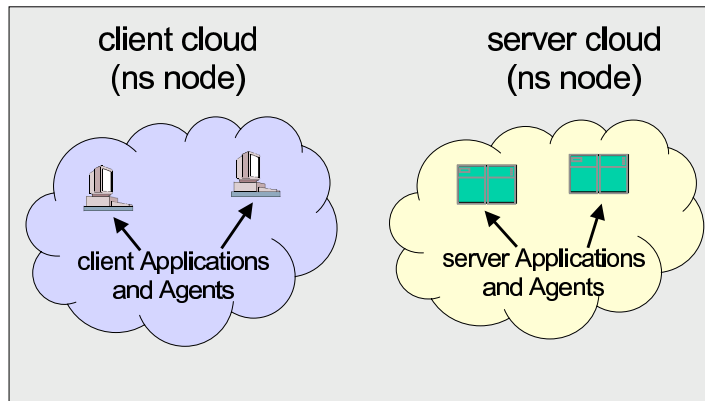


Figure 40.1: PackMimeHTTP Architecture. Each PackMimeHTTP object controls a server and a client cloud. Each cloud can represent multiple client or server Applications. Each Application represents either a single web server or a single web client.

server Applications. The pools for active Applications ensure that all active Applications are destroyed when the simulation is finished. Active TCP Agents do not need to be placed in a pool because each active Application contains a pointer to its associated TCP Agent. New objects are only created when there are no Agents or Applications available in the inactive pools.

## 40.1.1 PackMimeHTTP Client Application

Each PackMimeHTTP client controls the HTTP request sizes that are transferred. Each PackMimeHTTP client takes the following steps:

- sample the number of requests for this connection from the number-of-requests distribution (if the number of requests is 1, this is a non-persistent connection)
- sample the inter-request times from the inter-request-time distribution, if there will be more than 1 request
- sample the HTTP request sizes from the request-size distribution
- send the first HTTP request to the server
- listen for the HTTP response
- when the entire HTTP response has been received, the client sets a timer to expire when the next request should be made
- when the timer expires, the next HTTP request is sent, and the above process is repeated until all requests have been completed

### 40.1.2 PackMimeHTTP Server Application

Each web server controls the response sizes that are transferred. The server is started by when a new TCP connection is started. Each PackMimeHTTP client takes the following steps:

- listen for an HTTP request from the associated client
- when the entire request arrives, the server samples the server delay time from the server delay distribution
- set a timer to expire when the server delay has passed
- when the timer expires, the server samples the HTTP response sizes from the HTTP response size distribution
- this process is repeated until the requests are exhausted – the server is told how many requests will be sent in the connection
- send a FIN to close the connection

## 40.2 PackMimeHTTP Random Variables

This implementation of PackMimeHTTP provides several ns RandomVariable objects for specifying distributions of PackMimeHTTP connection variables. The implementations were taken from source code provided by Bell Labs and modified to fit into the ns RandomVariable framework. This allows PackMimeHTTP connection variables to be specified by any type of ns RandomVariable, which now include PackMimeHTTP-specific random variables. The PackMimeHTTP-specific random variable syntax for TCL scripts is as follows:

- `$ns [new RandomVariable/PackMimeHTTPFlowArrive <rate>]`, where `rate` is the specified PackMimeHTTP connection rate (number of new connections per second)
- `$ns [new RandomVariable/PackMimeHTTPFileSize <rate> <type>]`, where `type` is 0 for HTTP requests and 1 for HTTP responses
- `$ns [RandomVariable/PackMimeHTTPXmit <rate> <type>]`, where `type` is 0 for client-side delays and 1 for server-side delays. **Note:** This random variable is only used in conjunction with DelayBox. It returns 1/2 of the actual delay because it is meant to be used with 2 DelayBox nodes, each of which should delay the packets for 1/2 of the actual delay.

## 40.3 Use of DelayBox with PackMime-HTTP

PackMimeHTTP uses ns to model the TCP-level interaction between web clients and servers on the simulated link. To simulate network-level effects of HTTP transfer through the clouds, we implemented a new ns module called DelayBox (see 22). DelayBox is an ns analog to dummynet, often used in network testbeds to delay and drop packets. The delay times model the propagation and queuing delay incurred from the source to the edge of the cloud (or edge of the cloud to destination). Since all HTTP connections in PackMimeHTTP take place between only two ns nodes, there had to be an ns object to delay packets in each flow, rather than just having a static delay on the link between the two nodes. DelayBox also models bottleneck links and packet loss on an individual connection basis. Two DelayBox nodes are used as shown in Figure 40.3. One node is placed in front of the web client cloud ns node to handle client-side delays, loss, and bottleneck links. The other DelayBox node is placed in front of the web server cloud ns node to handle the server-side delays, loss, and bottleneck links.

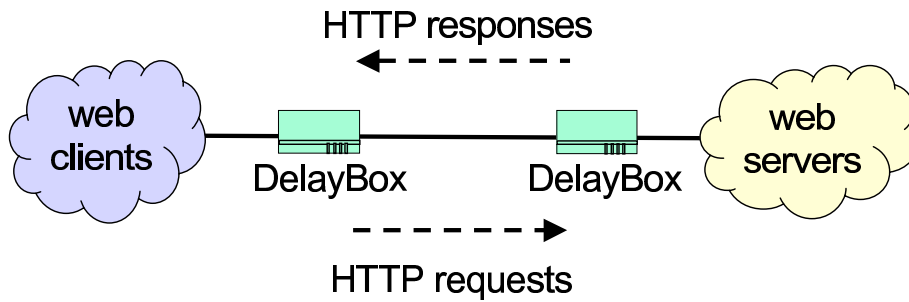


Figure 40.2: Example Topology Using PackMimeHTTP and DelayBox. The cloud of web clients is a single ns node, and the cloud of web servers is a single ns node. Each of the DelayBox nodes is a single ns node.

## 40.4 Example

More examples (including those that demonstrate the use of DelayBox with PackMime) are available in the `tcl/ex/packmime/` directory of the ns source code. The validation script `test-suite-packmime.tcl` is in `tcl/test/` and can be run with the command `test-all-packmime` from that directory.

```

# test-packmime.tcl

# useful constants
set CLIENT 0
set SERVER 1

remove-all-packet-headers;          # removes all packet headers
add-packet-header IP TCP;            # adds TCP/IP headers
set ns [new Simulator];              # instantiate the Simulator
$ns use-scheduler Heap;              # use the Heap scheduler

# SETUP TOPOLOGY
# create nodes
set n(0) [$ns node]

```

```

set n(1) [$ns node]
# create link
$ns duplex-link $n(0) $n(1) 10Mb 0ms DropTail

# SETUP PACKMIME
set rate 15
set pm [new PackMimeHTTP]
$pm set-client $n(0);           # name $n(0) as client
$pm set-server $n(1);          # name $n(1) as server
$pm set-rate $rate;             # new connections per second
$pm set-http-1.1;               # use HTTP/1.1

# SETUP PACKMIME RANDOM VARIABLES
global defaultRNG

# create RNGs (appropriate RNG seeds are assigned automatically)
set flowRNG [new RNG]
set reqsizeRNG [new RNG]
set rspsizeRNG [new RNG]

# create RandomVariables
set flow_arrive [new RandomVariable/PackMimeHTTPFlowArrive $rate]
set req_size [new RandomVariable/PackMimeHTTPFileSize $rate $CLIENT]
set rsp_size [new RandomVariable/PackMimeHTTPFileSize $rate $SERVER]

# assign RNGs to RandomVariables
$flow_arrive use-rng $flowRNG
$req_size use-rng $reqsizeRNG
$rsp_size use-rng $rspsizeRNG

# set PackMime variables
$pm set-flow_arrive $flow_arrive
$pm set-req_size $req_size
$pm set-rsp_size $rsp_size

# record HTTP statistics
$pm set-outfile "data-test-packmime.dat"

$ns at 0.0 "$pm start"
$ns at 300.0 "$pm stop"
$ns at 301.0 "exit 0"

$ns run

```

## 40.5 Commands at a Glance

The following commands on the PackMimeHTTP class can be accessed from OTcl:

```
[new PackMimeHTTP]
```

Creates a new PackMimeHTTP object.

`$packmime set-client <node>`  
Associates the node with the PackMimeHTTP client cloud

`$packmime set-server <node>`  
Associates the node with the PackMimeHTTP server cloud

`$packmime set-rate <float>`  
Set the average number of new connections started per second

`$packmime set-req_size <RandomVariable>`  
Set the HTTP request size distribution

`$packmime set-rsp_size <RandomVariable>`  
Set the HTTP response size distribution

`$packmime set-flow_arrive <RandomVariable>`  
Set the distribution of time between two consecutive connections starting

`$packmime set-server_delay <RandomVariable>`  
Set the web server delay for fetching pages

`$packmime start`  
Start generating connections

`$packmime stop`  
Stop generating new connections

`$packmime set-http-1.1`  
Use HTTP/1.1 distributions for persistent connections instead of HTTP/1.0.

`$packmime set-run <int>`  
Set the run number so that the RNGs used for the random variables will use the same substream (see Chapter 24 on RNG for more details).

`$packmime active-connections`  
Output the current number of active HTTP connections to standard error

`$packmime total-connections`  
Output the total number of completed HTTP connections to standard error

`$packmime get-pairs`  
Return the number of completed HTTP request-response pairs. See `tcl/ex/packmime/pm-end-pairs.tcl` for an example of using `get-pairs` to end the simulation after a certain number of pairs have completed.

`$packmime set-TCP <protocol>`  
Sets the TCP type (Reno, Newreno, or Sack) for all connections in the client and server clouds - Reno is the default

`$packmime set-outfile <filename>`  
Output the following fields (one line per HTTP request-response pair) to filename:

- time HTTP response completed
- HTTP request size (bytes)
- HTTP response size (bytes)

- HTTP response time (ms) – time between client sending HTTP request and client receiving complete HTTP response
- source node and port identifier
- number of active connections at the time this HTTP request-response pair completed

`$packmime set-debug <int>`

Set the debugging level:

- 1: Output the total number of connections created at the end of the simulation
- 2: Level 1 +  
output creation/management of TCP agents and applications  
output on start of new connection  
number of bytes sent by the client and expected response size  
number of bytes sent by server
- 3: Level 2 +  
output when TCP agents and applications are moved to the pool
- 4: Level 3 +  
output number of bytes received each time client or server receive a packet

## **Part VII**

# **Scale**



# Chapter 41

## Session-level Packet Distribution

This section describes the internals of the Session-level Packet Distribution implementation in *ns*. The section is in two parts: the first part is an overview of Session configuration (Section 41.1), and a “complete” description of the configuration parameters of a Session. The second part describes the architecture, internals, and the code path of the Session-level Packet distribution.

The procedures and functions described in this chapter can be found in `~ns/tcl/session/session.tcl`.

Session-level Packet Distribution is oriented towards performing multicast simulations over large topologies. The memory requirements for some topologies using session level simulations are:

2048 nodes, degree of connectivity = 8	$\approx$ 40 MB
2049–4096 nodes	$\approx$ 167 MB
4097–8194 nodes	$\approx$ 671 MB

Note however, that session level simulations ignore queueing delays. Therefore, the accuracy of simulations that use sources with a high data rate, or those that use multiple sources that get aggregated at points within the network is suspect.

### 41.1 Configuration

Configuration of a session level simulation consists of two parts, configuration of the session level details themselves (Section 41.1.1) and adding loss and error models to the session level abstraction to model specific behaviours (Section 41.1.2).

#### 41.1.1 Basic Configuration

The basic configuration consists of creating and configuring a multicast session. Each Session (*i.e.*, a multicast tree) must be configured strictly in this order: (1) create and configure the session source, (2) create the session helper and attach it to the session source, and finally, (3) have the session members join the session.

```
set ns [new SessionSim]                                ; # preamble initialization
set node [$ns node]
set group [$ns allocaddr]
```

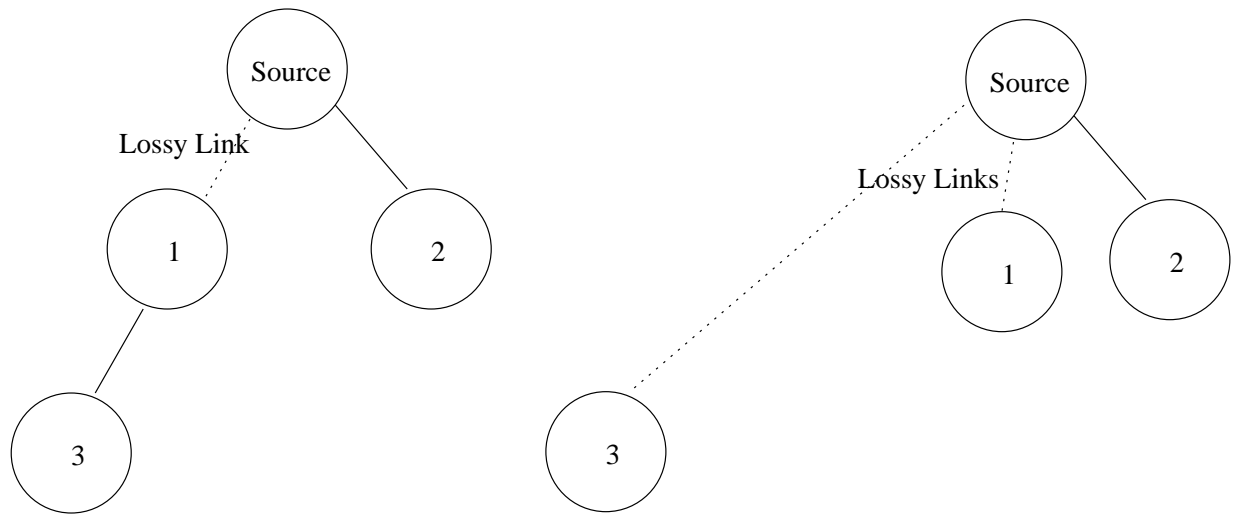


Figure 41.1: Comparison of Multicast Trees for Detailed vs. Session Routing

```

set udp [new Agent/UDP]                                ; # create and configure the source
$udp set dst_ $group
set src [new Application/Traffic/CBR]
$src attach-agent $udp
$ns attach-agent $node $udp

$ns create-session $node $udp                            ; # create attach session helper to src

set rcvr [new Agent/NULL]                                ; # configure the receiver
$ns attach-agent $node $rcvr
$ns at 0.0 "$node join-group $rcvr $group"              ; # joining the session

$ns at 0.1 "$src start"

```

A session level simulation scales by translating the topology into a virtual mesh topology. The steps involved in doing this are:

1. All of the classifiers and replicators are eliminated. Each node only stores instance variables to track its node id, and port ids.
2. Links do not consist of multiple components. Each link only stores instance variables to track the bandwidth and delay attributes.
3. The topology, consisting of links is translated into a virtual mesh.

Figure 41.1 shows the difference between a multicast tree in a detailed simulation and one in a session level simulation. Notice that the translation process results in a session level simulation ignoring queuing delays. For most simulations, *ns* already ignores processing delays at all of the nodes.

### 41.1.2 Inserting a Loss Module

When studying a protocol (*e.g.*, SRM error recovery mechanism), it might be useful to study protocol behavior over lossy links. However, since a session level simulation scales by abstracting out the internal topology, we need additional mechanisms to insert a loss module appropriately. This subsection describes how one can create these loss modules to model error scenarios.

**Creating a Loss Module** Before we can insert a loss module in between a source-receiver pair, we have to create the loss module. Basically, a loss module compares two values to decide whether to drop a packet. The first value is obtained every time when the loss module receives a packet from a random variable. The second value is fixed and configured when the loss module is created.

The following code gives an example to create a uniform 0.1 loss rate.

```
# creating the uniform distribution random variable
set loss_random_variable [new RandomVariable/Uniform]
$loss_random_variable set min_ 0                ;# set the range of the random variable
$loss_random_variable set max_ 100

set loss_module [new ErrorModel]                ;# create the error model
$loss_module drop-target [new Agent/Null]
$loss_module set rate_ 10                       ;# set error rate to 0.1 = 10 / (100 - 0)
$loss_module ranvar $loss_random_variable        ;# attach random var. to loss module
```

A catalogue of the random variable distributions was described earlier (Chapter 24). A more detailed discussion of error models was also described earlier in a different chapter (Chapter 13).

**Inserting a Loss Module** A loss module can only be inserted after the corresponding receiver has joined the group. The example code below illustrates how a simulation script can introduce a loss module.

```
set sessionhelper [$ns create-session $node $src] ;# keep a handle to the loss module
$ns at 0.1 "$sessionhelper insert-depended-loss $loss_module $rcvr"
```

## 41.2 Architecture

The purpose of Session-level packet distribution is to speed up simulations and reduce memory consumption while maintaining reasonable accuracy. The first bottleneck observed is the memory consumption by heavy-weight links and nodes. Therefore, in SessionSim (Simulator for Session-level packet distribution), we keep only minimal amount of states for links and nodes, and connect the higher level source and receiver applications with appropriate delay and loss modules. A particular source in a group sends its data packets to a replicator that is responsible for replicating the packets for all the receivers. Intermediate loss and delay modules between this replicator and the receivers will guarantee the appropriate end-to-end characteristics. To put it another way, a session level simulation abstracts out the topology, routing and queueing delays. Packets in SessionSim do not get routed. They only follow the established Session.

## 41.3 Internals

This section describes the internals of Session-level Packet Distribution. We first describe the OTcl primitives to configure a session level simulation (Section 41.3.1); we conclude with a brief note on how packet forwarding is achieved (Section 41.3.2).

### 41.3.1 Object Linkage

We describe three aspects of constructing a session level simulation in *ns*: the modified topology routines that permit the creation of abstract nodes and links, establish the session helper for each active source, add receivers to the session by inserting the appropriate loss and delay models when that receiver joins the appropriate group.

**Nodes and Links** The node contains only its node id and the port number for the next agent. A link only contains the values of its bandwidth and delay.

```
SessionNode instproc init {} {
    $self instvar id_ np_
    set id_ [Node getId]
    set np_ 0
}

SessionSim instproc simplex-link { n1 n2 bw delay type } {
    $self instvar bw_ delay_
    set sid [$n1 id]
    set did [$n2 id]

    set bw_($sid:$did) [expr [string trimright $bw Mb] * 1000000]
    set delay_($sid:$did) [expr [string trimright $delay ms] * 0.001]
}
```

**Session Helper** Each active source in a session requires a “session helper”. The session helper in *ns* is realised through a replicator. This session helper is created when the user issues a `create-session{}` to identify the source agent. The simulator itself keeps a reference to the session helper in its instance variable array, `session_`, indexed by the source and destination address of the source.

Note that the destination of source agent must be set before calling `create-session{}`.

```
SessionSim instproc create-session { node agent } {
    $self instvar session_

    set nid [$node id]
    set dst [$agent set dst_]
    set session_($nid:$dst) [new Classifier/Replicator/Demuxer]
    $agent target $session_($nid:$dst) ;# attach the replicator to the source
    return $session_($nid:$dst) ;# keep the replicator in the SessionSim instance variable array session_
}
```

**Delay and Loss Modules** Each receiver in a group requires a delay module that reflects its delay with respect to the particular source. When the receiver joins a group, `join-group{}` identifies all session helpers in `session_`. If the destination index matches the group address the receiver are joining, then the following actions are performed.

1. A new slot of the session helper is created and assigned to the receiver.
2. The routine computes the accumulated bandwidth and delay between the source and receiver using the SessionSim instance procedures `get-bw{}` and `get-delay{}`.
3. A constant random variable is created; it will generate random delivery times using the accumulative delay as an estimate of the average delay.
4. A new delay module is created with the end-to-end bandwidth characteristics, and the random variable generator provides the delay estimates.
5. The delay module is inserted into the session helper and interposed between the helper and the receiver.

See Section 41.1.2 for similarly inserting a loss module for a receiver.

```
SessionSim instproc join-group { agent group } {
    $self instvar session_

    foreach index [array names session_] {
        set pair [split $index :]
        if {[lindex $pair 1] == $group} {
            # Note: must insert the chain of loss, delay,
            # and destination agent in this order:

            $session_($index) insert $agent           ;# insert destination agent into session replicator

            set src [lindex $pair 0]                  ;# find accum. b/w and delay
            set dst [[ $agent set node_] id]
            set accu_bw [$self get-bw $dst $src]
            set delay [$self get-delay $dst $src]

            set random_variable [new RandomVariable/Constant] ;# set delay variable
            $random_variable set avg_ $delay

            set delay_module [new DelayModule]          ;# configure the delay module
            $delay_module bandwidth $accu_bw
            $delay_module ranvar $random_variable

            $session_($index) insert-module $delay_module $agent ;# insert the delay module
        }
    }
}
```

### 41.3.2 Packet Forwarding

Packet forwarding activities are executed in C++. A source application generates a packet and forwards to its target which must be a replicator (session helper). The replicator copies the packet and forwards to targets in the active slots which are either delay modules or loss modules. If loss modules, a decision is made whether to drop the packet. If yes, the packet is

Figure 41.2: Architectural Realization of a Session Level Simulation Session

forwarded to the loss modules drop target. If not, the loss module forwards it to its target which must be a delay module. The delay module will forward the packet with a delay to its target which must be a receiver application.

## 41.4 Commands at a glance

Following is a list of session-level related commands:

```
set ns [new SessionSim]
```

This command creates an instance of the sessionmode simulator.

```
$ns_ create-session <node> <agent>
```

This command creates and attaches a session-helper, which is basically a replicator, for the source <agent> created at the <node>.

## Chapter 42

# Asim: approximate analytical simulation

This chapter describes a fast approximate network simulator, Asim. Asim solves the steady state of the network using approximate fixed points. The overall structure is shown in Figure 42.1. The user feeds a regular ns script and turns on the asim flag. Asim would do a fast approximate simulation of the network scenario and would present to the user the drop probabilities of the routers, the delays and the approximate aggregate throughput of the links and the flows.

In particular, we the following links/traffic are supported:

- Drop Tail Queues
- RED Queues
- Bulk TCP flows with FTP traffic
- Short lived TCP flows

The data structures of Asim are populated by a module within the Tcl space of ns from the user supplied script. Upon executing Asim, the results can be accessed using Tcl routines. To use the Asim within a script the user has to use

Simulator set useasim\_1

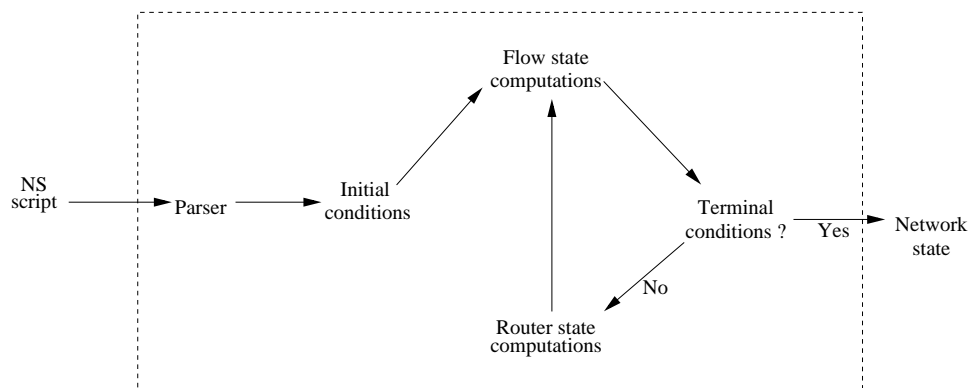


Figure 42.1: The structure of Asim

By default, this flag is set to 0

A simple script is given below

```
proc addsrc { s } {
    global ns
    set t [$ns set src_]
    lappend t $s
    $ns set src_ $t
}

proc adddst { src } {
    global ns
    set t [$ns set dst_]
    lappend t $src
    $ns set dst_ $t
}

proc finish {} {

    global ns fmon

    set drops [$fmon set pdrops_]
    set pkts [$fmon set parrivals_]
    set notDropped [$fmon set pdepartures_]

    set overflow_prob [expr 1.0 * $drops / $pkts]
    puts [format "tdrops $drops tpkts $pkts o_prob. %7.4f" $overflow_prob]

    exit 0
}

set N_ 100000
set arrival 0
set available $N_
set endTime_ 200

set ns [new Simulator]
$ns set useasim_ 1
$ns at $endTime_ "finish"

set src_ ""
set dst_ ""

$ns set src_ $src_
$ns set dst_ $dst_

set n(0) [$ns node]
```



```

set n(1) [$ns node]
set link(0:1) [$ns duplex-link $n(0) $n(1) 1Mbps 50ms RED]

for {set i 0} { $i < 4} {incr i} {

set ltcp($i) [new Agent/TCP]
set ltcpsink($i) [new Agent/TCPSink]
$ns attach-agent $n(0) $ltcp($i)
$ns attach-agent $n(1) $ltcpsink($i)
$ns connect $ltcp($i) $ltcpsink($i)

set lftp($i) [new Application/FTP]
$lftp($i) attach-agent $ltcp($i)
$ns at 0 "$lftp($i) start"

}

# Short term flows
addsrc 1
adddst 0

set pool [new PagePool/WebTraf]

# Set up server and client nodes
$pool set-num-client [llength [$ns set src_]]
$pool set-num-server [llength [$ns set dst_]]
global n
set i 0
foreach s [$ns set src_] {
    $pool set-client $i $n($s)
    incr i
}
set i 0
foreach s [$ns set dst_] {
    $pool set-server $i $n($s)
    incr i
}

# Number of Pages per Session
set numPage 100000

$pool set-num-session 1

set interPage [new RandomVariable/Exponential]
$interPage set avg_ 0.5

set pageSize [new RandomVariable/Constant]
$pageSize set val_ 1

set interObj [new RandomVariable/Exponential]
$interObj set avg_ 1

set objSize [new RandomVariable/Constant]
$objSize set val_ 20

```

```

# This is needed
$pool use-asim

$pool create-session 0 $numPage 0 $interPage $pageSize $interObj $objSize

# Dumps internal data structures to this dumpfile
$ns asim-dump dumpfile

# Calls asim-run
$ns asim-run

# Access asim statistics
set l [$ns link $n(0) $n(1)]
puts "after asim run, link bw = [$ns asim-getLinkTput $l] packets"
puts "after asim run, flow bw = [$ns asim-getFlowTput $ltcp(0)] packets"

```

## **Part VIII**

# **Emulation**

# Chapter 43

## Emulation

This chapter describes the *emulation* facility of *ns*. Emulation refers to the ability to introduce the simulator into a live network. Special objects within the simulator are capable of introducing live traffic into the simulator and injecting traffic from the simulator into the live network.

### Emulator caveats:

- While the interfaces described below are not expected to change drastically, this facility is still under development and should be considered experimental and subject to change.
- The facility described here has been developed under FreeBSD 2.2.5, and use on other systems has not been tested by the author.
- Because of the currently limited portability of emulation, it is only compiled into *nse* (build it with “make nse”), not standard *ns*.

### 43.1 Introduction

The emulation facility can be subdivided into two modes:

1. **opaque mode** – live data treated as opaque data packets
2. **protocol mode** – live data may be interpreted/generated by simulator

In opaque mode, the simulator treats network data as uninterpreted packets. In particular, real-world protocol fields are not directly manipulated by the simulator. In opaque mode, live data packets may be dropped, delayed, re-ordered, or duplicated, but because no protocol processing is performed, protocol-specific traffic manipulation scenarios (e.g. “drop the TCP segment containing a retransmission of sequence number 23045”) may not be performed. In protocol mode, the simulator is able to interpret and/or generate live network traffic containing arbitrary field assignments. **To date (Mar 1998), only Opaque Mode is currently implemented.**

The interface between the simulator and live network is provided by a collection of objects including *tap agents* and *network objects*. Tap agents embed live network data into simulated packets and vice-versa. Network objects are installed in tap agents and provide an entrypoint for the sending and receipt of live data. Both objects are described in the following sections.

When using the emulation mode, a special version of the system scheduler is used: the `RealTime` scheduler. This scheduler uses the same underlying structure as the standard calendar-queue based scheduler, but ties the execution of events to real-time. It is described below.

## 43.2 Real-Time Scheduler

The real-time scheduler implements a soft real-time scheduler which ties event execution within the simulator to real time. Provided sufficient CPU horsepower is available to keep up with arriving packets, the simulator virtual time should closely track real-time. If the simulator becomes too slow to keep up with elapsing real time, a warning is continually produced if the skew exceeds a pre-specified constant “slop factor” (currently 10ms).

The main dispatch loop is found in the routine `RealTimeScheduler::run()`, in the file `scheduler.cc`. It follows essentially the following algorithm:

- While simulator is not halted
  - get current real time (“now”)
  - dispatch all pending simulator events prior to now
  - fetch next (future) event if there is one
  - delay until the next simulator event is ready or a Tcl event occurs
  - if a tcl event occurred, re-insert next event in simulator event queue and continue
  - otherwise, dispatch simulator event, continue
  - if there was no future even, check for Tcl events and continue

The real-time scheduler should always be used with the emulation facility. Failure to do so may easily result in the simulator running faster than real-time. In such cases, traffic passing through the simulated network will not be delayed by the proper amount of time. Enabling the real-time scheduler requires the following specification at the beginning of a simulation script:

```
set ns [new Simulator]
$ns use-scheduler RealTime
```

## 43.3 Tap Agents

The class `TapAgent` is a simple class derived from the base `Agent` class. As such, it is able to generate simulator packets containing arbitrarily-assigned values within the `ns` common header. The tap agent handles the setting of the common header packet size field and the type field. It uses the packet type `PT_LIVE` for packets injected into the simulator. Each tap agent can have at most one associated network object, although more than one tap agent may be instantiated on a single simulator node.

**Configuration** Tap agents are able to send and receive packets to/from an associated `Network` object. Assuming a network object `$netobj` refers to a network object, a tap agent is configured using the `network` method:

```
set a0 [new Agent/Tap]
```

```

$a0 network $netobj
$a0 set fid_ 26
$a0 set prio_ 2
$ns connect $a0 $a1

```

Note that the configuration of the flow ID and priority are handled through the Agent base class. The purpose of setting the flow id field in the common header is to label packets belonging to particular flows of live data. Such packets can be differentially treated with respect to drops, reorderings, etc. The `connect` method instructs agent `$a0` to send its live traffic to the `$a1` agent via the current route through the simulated topology.

## 43.4 Network Objects

Network objects provide access to a live network. There are several forms of network objects, depending on the protocol layer specified for access to the underlying network, in addition to the facilities provided by the host operating system. Use of some network objects requires special access privileges where noted. Generally, network objects provide an entrypoint into the live network at a particular protocol layer (e.g. link, raw IP, UDP, etc) and with a particular access mode (read-only, write-only, or read-write). Some network objects provide specialized facilities such as filtering or promiscuous access (i.e. the pcap/bpf network object) or group membership (i.e. UDP/IP multicast). The C++ class `Network` is provided as a base class from which specific network objects are derived. Three network objects are currently supported: pcap/bpf, raw IP, and UDP/IP. Each are described below.

### 43.4.1 Pcap/BPF Network Objects

These objects provide an extended interface to the LBNL packet capture library (libpcap). (See <ftp://ftp.ee.lbl.gov/libpcap.tar.gz> for more info). This library provides the ability to capture link-layer frames in a promiscuous fashion from network interface drivers (i.e. a copy is made for those programs making use of libpcap). It also provides the ability to read and write packet trace files in the “tcpdump” format. The extended interface provided by `ns` also allows for writing frames out to the network interface driver, provided the driver itself allows this action. Use of the library to capture or create live traffic may be protected; one generally requires at least read access to the system’s packet filter facility which may need to be arranged through a system administrator.

The packet capture library works on several UNIX-based platforms. It is optimized for use with the Berkeley Packet Filter (BPF) [25], and provides a filter compiler for the BPF pseudomachine machine code. On most systems supporting it, a kernel-resident BPF implementation processes the filter code, and applies the resulting pattern matching instructions to received frames. Those frames matching the patterns are received through the BPF machinery; those not matching the pattern are otherwise unaffected. BPF also supports sending link-layer frames. This is generally not suggested, as an entire properly-formatted frame must be created prior to handing it off to BPF. This may be problematic with respect to assigning proper link-layer headers for next-hop destinations. It is generally preferable to use the raw IP network object for sending IP packets, as the system’s routing function will be used to determine proper link-layer encapsulating headers.

**Configuration** Pcap network objects may be configured as either associated with a live network or with a trace file. If associated with a live network, the particular network interface to be used may be specified, as well as an optional promiscuous flag. As with all network objects, they may be opened for reading or writing. Here is an example:

```

set me [exec hostname]
set pfl [new Network/Pcap/Live]
$pfl set promisc_ true

```

```

set intf [$pfl open readonly]
puts "pfl configured on interface $intf"
set filt "(ip src host foobar) and (not ether broadcast)"
set nbytes [$pfl filter $filt]
puts "filter compiled to $nbytes bytes"
puts "drops: [$pfl pdrops], pkts: [$pfl pkts]"

```

This example first determines the name of the local system which will be used in constructing a BPF/libpcap filter predicate. The new `Network/Pcap/Live` call creates an instance of the pcap network object for capturing live traffic. The `promisc_` flag tells the packet filter whether it should configure the underlying interface in promiscuous mode (if it is supported). The `open` call activates the packet filter, and may be specified as `readonly`, `writeonly`, or `readwrite`. It returns the name of the network interface the filter is associated with. The `open` call takes an optional extra parameter (not illustrated) indicating the name of the interface to use in cases where a particular interface should be used on a multi-homed host. The `filter` method is used to create a BPF-compatible packet filter program which is loaded into the underlying BPF machinery. The `filter` method returns the number of bytes used by the filter predicate. The `pdrops` and `pkts` methods are available for statistics collection. They report the number of packets dropped by the filter due to buffer exhaustion and the total number of packets that arrived at the filter, respectively (*not* the number of packets accepted by the filter).

### 43.4.2 IP Network Objects

These objects provide raw access to the IP protocol, and allow the complete specification of IP packets (including header). The implementation makes use of a *raw socket*. In most UNIX systems, access to such sockets requires super-user privileges. In addition, the interface to raw sockets is somewhat less standard than other types of sockets. The class `Network/IP` provides raw IP functionality plus a base class from which other network objects implementing higher-layer protocols are derived.

**Configuration** The configuration of a raw IP network object is comparatively simple. The object is not associated with any particular physical network interface; the system's IP routing capability will be used to emit the specified datagram out whichever interface is required to reach the destination address contained in the header. Here is an example of configuring an IP object:

```

set ipnet [new Network/IP]
$ipnet open writeonly
...
$ipnet close

```

The IP network object supports only the `open` and `close` methods.

### 43.4.3 IP/UDP Network Objects

These objects provide access to the system's UDP implementation along with support for IP multicast group membership operations. **IN PROGRESS**

## 43.5 An Example

The following code illustrates a small but complete simulation script for setting up an emulation test using BPF and IP network objects. It was run on a multi-homed machine, and the simulator essentially provides routing capability by reading frames from one interface, passing them through the simulated network, and writing them out via the raw IP network object:

```
set me "10.0.1.1"
set ns [new Simulator]

$ns use-scheduler RealTime

#
# we want the test machine to have ip forwarding disabled, so
# check this (this is how to do so under FreeBSD at least)
#

set ipforw [exec sysctl -n net.inet.ip.forwarding]
if $ipforw
    puts "can not run with ip forwarding enabled"
    exit 1

#
# allocate a BPF type network object and a raw-IP object
#
set bpf0 [new Network/Pcap/Live]
set bpf1 [new Network/Pcap/Live]
$bpf0 set promisc_ true
$bpf1 set promisc_ true

set ipnet [new Network/IP]

set nd0 [$bpf0 open readonly fxp0]
set nd1 [$bpf1 open readonly fxp1]
$ipnet open writeonly

#
# try to filter out weird stuff like netbios pkts, arp requests, dns,
# also, don't catch stuff to/from myself or broadcasted
#
set notme "(not ip host $me)"
set notbcast "(not ether broadcast)"
set ftp "and port ftp-data"
set f0len [$bpf0 filter "(ip dst host bit) and $notme and $notbcast"]
set f1len [$bpf1 filter "(ip src host bit) and $notme and $notbcast"]

puts "filter lengths: $f0len (bpf0), $f1len (bpf1)"
puts "dev $nd0 has address [$bpf0 linkaddr]"
puts "dev $nd1 has address [$bpf1 linkaddr]"

set a0 [new Agent/Tap]
set a1 [new Agent/Tap]
set a2 [new Agent/Tap]
```



```

puts "install nets into taps..."
$a0 network $bpf0
$a1 network $bpf1
$a2 network $ipnet

set node0 [$ns node]
set node1 [$ns node]
set node2 [$ns node]

$ns simplex-link $node0 $node2 10Mb 10ms DropTail
$ns simplex-link $node1 $node2 10Mb 10ms DropTail

$ns attach-agent $node0 $a0
$ns attach-agent $node1 $a1
$ns attach-agent $node2 $a2

$ns connect $a0 $a2
$ns connect $a1 $a2

puts "okey"
$ns run

```

## 43.6 Commands at a glance

Following is a list of emulation related commands:

`$ns_ use-scheduler RealTime`

This command sets up the real-time scheduler. Note that a real-time scheduler should be used with any emulation facility. Otherwise it may result the simulated network running faster than real-time.

`set netob [new Network/<network-object-type>]`

This command creates an instance of a network object. Network objects are used to access a live network. Currently the types of network objects available are Network/Pcap/Live, Network/IP and Network/IP/UDP. See section 43.4 for details on network objects.

## **Part IX**

# **Visualization with Nam - The Network Animator**

# Chapter 44

## Nam

### 44.1 Introduction

Nam is a Tcl/Tk based animation tool for viewing network simulation traces and real world packet tracedata. The design theory behind nam was to create an animator that is able to read large animation data sets and be extensible enough so that it could be used indifferent network visualization situations. Under this constraint nam was designed to read simple animation event commands from a large trace file. In order to handle large animation data sets a minimum amount of information is kept in memory. Event commands are kept in the file and reread from the file whenever necessary.

The first step to use nam is to produce the trace file. The trace file contains topology information, e.g., nodes, links, as well as packet traces. The detailed format is described in the section 45.1. Usually, the trace file is generated by ns. During an ns simulation, user can produce topology configurations, layout information, and packet traces using tracing events in ns. However any application can generate a nam trace file.

When the trace file is generated, it is ready to be animated by nam. Upon startup, nam will read the tracefile, create topology, pop up a window, do layout if necessary, and then pause at time 0. Through its user interface, nam provides control over many aspects of animation. These functionalities will be described in detail in the USER INTERFACE section.

There are bugs in nam however each successive has become much more stable than the previous one. Please mail ns-users@isi.edu if you encounter any bugs, or have suggestions for additional desired functionality.

### 44.2 Nam Command Line Options

```
nam [ -g <geometry> ] [ -t <graphInput> ] [ -i <interval> ] [ -j <startup time> ]  
    [ -k <initial socket port number> ] [ -N <application name> ] [ -c <cache size> ]  
    [ -f <configuration file> ] [ -r initial animation rate ]  
    [ -a ] [ -p ] [ -S ]  
    [ <tracefile(s)> ]
```

#### Command Line Options

-g Specify geometry of the window upon startup.

-t Instruct nam to use tkgraph, and specify input file nam for tkgraph.

-i [Information for this option may not be accurate] Specify rate (real) milliseconds as the screenupdate rate. The default rate is 100.

-N Specify the application name of this nam instance. This application name may later be used in peer synchronization.

-c The maximum size of the cache used to store 'active' objects when doing animating in reverse.

-f Name of the initialization files to be loaded during startup. In this file, user can define functions which will be called in the initialization.

-a Create a separate instance of nam.

-p Print out nam trace file format.

-S Enable synchronous X behavior so it is easier for graphics debugging. For UNIX system running X only.

<tracefile> is the name of the file containing the trace data to be animated. If <tracefile> cannot be read, nam will try to open <tracefile>.

## 44.3 User Interface

Starting up nam will first create the nam console window. You can have multiple animations running under the same nam instance. At the top of all nam windows is a menu bar. For the nam console there are 'File' and 'Help' menus. Under the 'File' there is a 'New' command for creating a new topology using the nam editor (under construction), an 'Open' command which allows you to open existing tracefiles, a 'WinList' command that popup a window with the names of all currently opened tracefiles, and a 'Quit' command which exits nam. The 'Help' menu contains a very limited popup help screen and a command to show version and copyright information.

Once a tracefile has been loaded into nam (either by using the 'Open' menu command or by specifying the tracefile on the command line) an animation window will appear. It has a 'Save layout' command which will save the current network layout to a file and a 'Print' command which will print the current network layout.

The 'Views' menu has 4 buttons:

- New view button: Creates a new view of the same animation. User can scroll and zoom on the newview. All views will be animated synchronously.
- Show monitors checkbox: If checked, will show a pane at the lower half of window, where monitors will be displayed.
- Show autolayout checkbox: If checked, will show a pane at the lower half of window, which contains input boxes and a button for automatic layout adjustments. This box will not be enabled when using link orientain layouts.
- Show annotation checkbox: If checked, will show a listbox at the lower half of window, which will be used to list annotations in the ascending order of time.

Below the menu bar, there is a control bar containing 6 buttons, a label, and a small scrollbar(scale). They can be clicked in any order. We will explain them from left to right.

- Button 1 (◀) - Rewind. When clicked, animation time will go back at the rate of 25 times the current screen update rate.
- Button 2 (<) - Backward play. When clicked, animation will be played backward with time decreasing.
- Button 3 (square) - Stop. When clicked, animation will pause.
- Button 4 (>) - Forward play. When clicked, animation will be played forward with time increasing.
- Button 5 (▶) - Fast Forward. When clicked, animation time will go forward at the rate of 25 times the current screen update rate.

- Button 6 (Chevron logo) - Close current animation window.

Time label - Show the current animation time (i.e., simulation time as in the trace file). Rate Slider - Controls the screen update rate (animation granularity). The current rate is displayed in the label above the slider.

Below the first control bar, there is Main Display, which contains a tool bar and a main view pane with two panning scroll bars. All new views created by menu command 'Views/New view' will have these three components. The tool bar contains two zoom buttons. The button with an up arrow zooms in, the button with a down arrow zooms out. The two scroll bars are used to pan the main animation view.

Clicking the left button on any of the objects in the main view pane will pop up a information window. For packet and agent objects, there is a 'monitor' button in the popup window. Clicking that button will bring out the monitor pane (if it is not already there), and add a monitor to the object. For link objects, there will be a 'Graph' button. Clicking on that button will bring up another popup window, where users can select between drawing a bandwidth utilization graph or drawing a link loss graph of one simplex edge of the duplex link.

Below the user interface objects we have discussed so far, there may or may not be a Monitor pane, depending on whether the checkbox 'Views/Show monitors' is set. (The default is unset). All monitors will be shown in this pane. A monitor looks like a big button in the pane. Currently only packets and agents may have monitors.

A packet monitor shows the size, id, and sent time. When the packet reaches its destination, the monitor will still be there, but will say that the packet is invisible. An agent monitor shows the name of the agent, and if there are any variable traces associated with this agent, they will be shown there as well.

Below the monitor pane (or in its place if the monitor pane isn't there), there is a Time Slider. It looks like a scaled ruler, with a tag 'TIME' which can be dragged along the ruler. It is used to set the current animation time. As you drag the 'TIME' tag, current animation time will be displayed in the time label in the control bar above. The left edge of the slider represents the earliest event time in the trace file and the right edge represents the last event time. Clicking left button on the ruler (not on the tag) has the same effect as Rewind or Fast Forward, depending on the clicking position.

The Automatic Layout Pane may be visible or hidden. If visible, it is below the time slider. It has three inputboxes and one relayout button. The labeled input boxes let user adjust two automatic layout constants, and the number of iterations during next layout. When user press ENTER in any of the input boxes, or click the 'relayout' button, that number of iterations will be performed. Refer to the AUTOMATIC LAYOUT section for details of usage.

The bottom component of the nam window is a Annotation Listbox, where annotations are displayed. An annotation is a (time, string) pair, which describes a event occurring at that time. Refer to ns(1) for functions to generate annotations. Double-clicking on an annotation in the listbox will bring nam to the time when that annotation is recorded. When the pointer is within the listbox, clicking the right button will stop the animation and bring up a popup menu with 3 options: Add, Delete, Info. 'Add' will bring up a dialog box with a text input to add a new annotation entry which has the current animation time. The user can type an annotation string in the dialog box. 'Delete' will delete the annotation entry pointed by the pointer. 'Info' will bring out a pane which shows both the annotation time and the annotation string.

## 44.4 Keyboard Commands

Most of the buttons have keyboard equivalents. Note they only function when the mouse cursor is inside the nam window.

- <return> - Typing a <return> will pause nam if it's not already paused. If nam is paused, <return> will step the animation one simulated clock tick. (If your keyboard autorepeats, holding down <return> is a good way to slow-step through the animation.)

- 'p' or 'P' - Pause but not step if paused.
- 'c' or 'C' - Continue after a pause.
- 'b' or 'B' - Decrease animation time for one screen update interval.
- 'r' or 'R' - Rewind.
- 'f' or 'F' - Fast Forward.
- 'n' or 'N' - Jump to time of next event in the trace file.
- 'x' or 'X' - Undo the last rate change.
- 'u' or 'U' - Undo the last time slider drag.
- '>' or '.' Increase the granularity (speed up) by 5%.
- '<' or ',' Decrease the granularity (slow down) by 5%.
- <space bar> - Toggle the pause state of nam.
- 'q', 'Q' or <control-c> - Quit.

## 44.5 Generating External Animations from Nam

Nam animations can be saved and converted to animated gifs or MPEG movies.

To save the frames of your movie, first start nam with your trace and set it up where you want it to start and adjust other parameters (step rate, size, etc.) Select 'Record Animation' from the File menu to start saving frames. Each animation step will be saved in a X-window dump file called "nam%d.xwd" where %d is the frame number. These files can then be assembled into animated GIFs or MPEGs with the appropriate post-processing tools.

The following shell script (sh, not csh) converts these files into an animated gif:

```
for i in *.xwd; do
  xwdtoppm <$i |
  pmtogif -interlace -transparent '#e5e5e5' >'basename $i .xwd'.gif;
done
gifmerge -l0 -2 -229,229,229 *.gif >movie.gif
```

Please note that the programs xwdtoppm, pmtogif, and gifmerge are *not* part of ns. You can get the first two from <http://download.sourceforge.net/netpbm/> and gifmerge from <http://www.the-labs.com/GIFMerge/>.

## 44.6 Network Layout

In nam, a topology is specified by alternating node objects with edge objects. But to display the topology in a comprehensible way, a layout mechanism is needed. Currently nam provides three layout methods. First, user may specify layout by the link's orientation. A link orientation is the angle between the edge and a horizontal line, in the interval  $[0, 2\pi]$ . During layout, nam will honor the given link orientation. Generally, it will first choose a reference node, then place other nodes using link orientations and link length. The link length is determined by link delay and connecting node sizes. This works well for small and manually generated topologies.

Second, when dealing with randomly generated topologies, we may want to do layout automatically. An automatic graph layout algorithm has been adapted and implemented. The basic idea of the algorithm is to model the graph as balls (nodes) connected by springs (links). Balls will repulse each other, while springs pull them together. This system will (hopefully) converge after some number of iterations. In practice, after a small number of iterations (tens or hundreds), most small to medium sized graphs will converge to a visually comprehensible structure. Larger graphs may take a combination of automatic layout and hand placement to achieve an acceptable layout.

There are 3 parameters to tune the automatic layout process: Ca Attractive force constant, which controls springs's force between balls. Cr Repulsive force constant, which controls the repulsive force between balls. Number of iterations How many times to run the autolayout procedure.

For small topologies with tens of nodes, using the default parameters (perhaps with 20 to 30 more iterations) will suffice to produce a nice layout. But for larger topology, careful parameter tuning is necessary. Following is a empirical method to layout a 100 node random transit stub topology generated by Georgia Tech's ITM internet topology modeler. First, set Ca and Cr to 0.2, do about 30 iterations, then set Cr to 1.0, Ca to about 0.01, then do about 10 iterations, then set Ca to 0.5, Cr to 1.0, do about 6 iterations.

Third, there is a x,y coordinate style layout. This was developed for use in displaying a wireless topologies in which permanent links don't exist. Using this style, nodes events are given x and y coordinate values indicating where those nodes should be placed in a cartesian world.

## 44.7 Animation Objects

Nam does animation using the following building blocks which are defined below:

**Node** Nodes are created from 'n' trace event in trace file. It represents a source, host, or router. Nam will skip over any duplicate definitions for the same node. A node may have three shapes, (circle, square, and hexagon), but once created it cannot change its shape. Nodes can change its color during animation. Nodes can be labeled.

**Link** Links are created between nodes to form a network topology. Internally nam links are consist of 2 simplex links. The trace event 'l' creates two simplex links and does other necessary setup. Therefore, for a users perspective all links are duplex links. Links can be labeled and also can change color during the animation. Links can be labeled as well.

**Queue** Queues need to be constructed in nam between two nodes. A nam queue is associated to only one edge of a duplex link. Queues are visualized as stacked packets. Packets are stacked along a line, the angle between the line and the horizontal line can be specified in the queue trace event.

**Packet** Packets are visualized as a block with an arrow. The direction of the arrow shows the flow direction of the packet. Queued packets are shown as little squares. A packet may be dropped from a queue or a link. Dropped packets are shown as falling rotating squares, and disappear at the end of the screen. Unfortunately, due to nam's design dropped packets are not visible during backward animation.

**Agent** Agents are used to separate protocol states from nodes. They are always associated with nodes. An agent has a name, which is a unique identifier of the agent. It is shown as a square with its name inside, and is drawn next to its associated node.

## Chapter 45

# Nam Trace

Nam is a Tcl/Tk based animation tool that is used to visualize the ns simulations and real world packet trace data. The first step to use nam is to produce a nam trace file. The nam trace file should contain topology information like nodes, links, queues, node connectivity etc as well as packet trace information. In this chapter we shall describe the nam trace format and simple ns commands/APIs that can be used to produce topology configurations and control animation in nam.

The underlying design constraints for nam were that it is able to handle large amounts of trace data and that its animation primitives be adaptable so that it may be used in different types of network visualization. As a result, internally nam reads information from a file and keeps only a minimum amount of animation event information in memory. Its animation event has a fairly simple and consistent structure so that it can many different visualization situations.

### 45.1 Nam Trace Format

The C++ class Trace used for ns tracing is used for nam tracing as well. Description of this class may be found under section 25.3. The method Trace::format() defines nam format used in nam trace files which are used by nam for visualization of ns simulations. Trace class method Trace::format() is described in section 25.4 of chapter 25. If the macro NAM\_TRACE has been defined (by default it is defined in trace.h), then the following code is executed as part of the Trace::format() function:

```
if (namChan_ != 0)
    sprintf(nwrk_,
            "%c -t \"TIME_FORMAT\" -s %d -d %d -p %s -e %d -c %d
-i %d -a %d -x %s.%s %s.%s %d %s %s",
            tt,
            Scheduler::instance().clock(),
            s,
            d,
            name,
            th->size(),
            iph->flowid(),
            th->uid(),
            iph->flowid(),
            src_nodeaddr,
            src_portaddr,
            dst_nodeaddr,
```



```
dst_portaddr,
seqno, flags, sname);
```

A nam trace file has a basic format to it. Each line is a nam event. The first character on the line defines the type of event and is followed by several flags to set options on that event. Each event is terminated by a newline character.

```
<event-type> -t <time> <more flags>...
```

Depending on the event type, there are different flags following the time flag.

There are 2 sections in that file, static initial configuration events and animation events. All events with -t \* in them are configuration events and should be at the beginning of the file. One thing to note is that nam can also be fed the trace file from a stream which enables it to be used with realtime applications. See the section **Using Streams with Realtime Applications** for more information.

Following we describe nam trace file format for different classes events and animation objects.

### 45.1.1 Initialization Events

The first section of a trace file must contain initialization information. All initialization events will have the flag -t \*. This tells nam that this event needs to be parsed before any animation has started.

**Version** The following line define the nam version as required to visualize the given trace:

```
V -t <time> -v <version> -a <attr>
```

Normally there is only one version string in a given tracefile, and it is usually the first line of the file. An example is the following:

```
V -t * -v 1.0a5 -a 0
```

The flag -v 1.0a5 tells nam that this script requires a version of nam > 1.0a5. For more information on this event look at the file tcl/stats.tcl under the procedure nam\_analysis.

**Wireless** If you want to use wireless nodes in nam you need the wireless initialization event.

```
W -t * -x 600 -y 600
```

This gives nam the size of the layout for the wireless world. The -x value is the width and -y is height. For more information look at the file animator.tcl in the procedure infer-network-model.

**Hierarchy** Hierarchical address information is defined by:

```
A -t <time> -n <levels> -o <address-space size> -c <mcastshift> -a <mcastmask> -h
<nth level> -m <mask in nth level> -s <shift in nth level>
```

This trace gives the details of hierarchy, if hierarchical addressing is being used for simulation. Flag -n <levels> indicate the total number of hierarchical tiers, which is 1 for flat addressing, 2 for a 2-level hierarchy etc. Flag -o <address-space size> denotes the total number of bits used for addressing. Flag -h <nth level> specifies the level of the address hierarchy. Flag -m <mask> and -s <shift> describes the address mask and the bit shift of a given level in the address hierarchy, respectively. Here is an example of a trace for topology with 3 level hierarchy:

```

A -t * -n 3 -p 0 -o 0xffffffff -c 31 -a 1
A -t * -h 1 -m 1023 -s 22
A -t * -h 2 -m 2047 -s 11
A -t * -h 3 -m 2047 -s 0

```

Look at tcl/netModel.tcl under the nam\_addressing procedure for more information.

**Color Table Entry** A table of color entries can be built using:

```
c -t <time> -i <color id> -n <color name>
```

Nam allows one to associate color names with integers. This is very useful in coloring packets. The flow id of a packet is used to color the packet using the corresponding color table entry color. Notice the color name should be one of the names listed in color database in X11 (/usr/X11/lib/rgb.txt).

In addition to the above node and link layout events are also included in the initialization section.

## 45.1.2 Nodes

The nam trace format defining node state is:

```
n -t <time> -a <src-addr> -s <src-id> -S <state> -v <shape> -c <color> -i <l-color> -o <color>
```

"n" denotes the node state.

Flags "-t" indicates time and "-a" and "-s" denotes the node address and id.

"-S" gives the node state transition. The possible state transition values are:

- UP, DOWN indicates node recovery and failure.
- COLOR indicates node color change. If COLOR is given, a following -c <color> is expected which gives the new color value. Also, flag -o is expected so that backtracing can restore the old color of a node.
- DLABEL indicates addition of label to node. If DLABEL is given, a following -l <old-label> -L <new-label> is expected that gives the old-label, if any (for backtracing) and current label. Shape gives the node shape. The color of a node label can be specified via the -i flag.

"-v" is the shape of the node. The possible values are:

- circle
- box
- hexagon

As an example, the line

```
n -t * -a 4 -s 4 -S UP -v circle -c tan -i tan
```

defines a node with address and id of 4 that has the shape of a circle, and color of tan and label-color (-i) of tan.

## 45.1.3 Links

The nam trace for link states is given by:

```
l -t <time> -s <src> -d <dst> -S <state> -c <color> -o orientation -r <bw> -D <delay>
```

where <state> and <color> indicate the same attributes (and the same format) as described above in the node state traces.

Flag -o gives the link orientation (angle between link and the horizon). Valid orientation values are:

- up
- down
- right
- left
- up-right
- down-right
- up-left
- down-left
- angle between 0 and 2pi

Flags `-r` and `-D` give the bandwidth (in Mb) and delay (in ms), respectively. An example of a link trace is:

```
l -t * -s 0 -d 1 -S UP -r 1500000 -D 0.01 -c black -o right
```

#### 45.1.4 Queues

The nam trace queue states is given by:

```
q -t <time> -s <src> -d <dst> -a <attr>
```

Queues are visualized in nam as a straight line along which packets (small squares) are packed. In queue trace events, the flag `-a` specifies the orientation of the line of the queue (angle between the queue line and the horizontal line, counter-clockwise). For example, the following line specifies a queue that grows vertically upwards with respect to the screen (here 0.5 means the angle of the queue line is  $\pi/2$ ):

```
q -t * -s 0 -d 1 -a 0.5
```

#### 45.1.5 Packets

When a trace line describes a packet, the event type may be + (enqueue), - (dequeue), r (receive), d (drop), or h (hop).

```
<type> -t <time> -e <extent> -s <source id> -d <destination id> -c <conv> -i <id>
```

`<type>` is one of:

**h** Hop: The packet started to be transmitted on the link from `<source id>` to `<destination id>` and is forwarded to the next hop.

**r** Receive: The packet finished transmission and started to be received at the destination.

**d** Drop: The packet was dropped from the queue or link from `<source id>` to `<destination id>`. Drop here doesn't distinguish between dropping from queue or link. This is decided by the drop time.

**+** Enter queue: The packet entered the queue from `<source id>` to `<destination id>`.

- Leave queue: The packet left the queue from <source id> to <destination id>.

The other flags have the following meaning:

- t <time> is the time the event occurred.
- s <source id> is the originating node.
- d <destination id> is the destination node.
- p <pkt-type> is the descriptive name of the type of packet seen. See section 25.5 for the different types of packets supported in *ns*.
- e <extent> is the size (in bytes) of the packet.
- c <conv> is the conversation id or flow-id of that session.
- i <id> is the packet id in the conversation.
- a <attr> is the packet attribute, which is currently used as color id.
- x <src-na.pa> <dst-sa.na> <seq> <flags> <sname> is taken from ns-traces and it gives the source and destination node and port addresses, sequence number, flags (if any) and the type of message. For example `-x {0.1 -2147483648.0 -1 ----- SRM_SESS}` denotes an SRM message being sent from node 0 (port 1).

Additional flags may be added for some protocols.

- P <packet type> gives an ASCII string specifying a comma separated list of packet types. Some values are:

**TCP** A tcp data packet.

**ACK** Generic acknowledgement.

**NACK** Generic negative acknowledgement.

**SRM** SRM data packet.

- n <sequence number> gives the packet sequence number.

### 45.1.6 Node Marking

Node marks are colored concentric circles, boxes, or hexagons around nodes. They are created by:

```
m -t <time> -n <mark name> -s <node> -c <color> -h <shape> -o <color>
```

and can be deleted by:

```
m -t <time> -n <mark name> -s <node> -X
```

Note that once created, a node mark cannot change its shape. The possible choices for shapes are, circle, box, and hexagon. They are defined as lowercase strings exactly as above. A nam trace showing a node mark is:

```
m -t 4 -s 0 -n m1 -c blue -h circle
```

indicating node 0 is marked with a blue circle at time 4.0. The name of the mark is m1.

### 45.1.7 Agent Tracing

Agent trace events are used to visualize protocol state. They are always associated with nodes. An agent event has a name, which is a *unique* identifier of the agent. An agent is shown as a square with its name inside, and a line link the square to its associated node

Agent events are constructed using the following format:

```
a -t <time> -n <agent name> -s <src>
```

Because in *ns*, agents may be detached from nodes, an agent may be deleted in nam with:

```
a -t <time> -n <agent name> -s <src> -X
```

For example, the following nam trace line creates an agent named `srm(5)` associated with node 5 at time 0:

```
a -t 0.000000000000000000 -s 5 -n srm(5)
```

### 45.1.8 Variable Tracing

To visualize state variables associated with a protocol agent, we use feature trace events. Currently we allow a feature to display a simple variable, i.e., a variable with a single value. Notice that the value is simple treated as a string (without white space). Every feature is required to be associated with an agent. Then, it can be added or modified at any time after its agent is created. The trace line to create a feature is:

```
f -t <time> -s <src> -a <agentname> -T <type> -n <varname> -v <value> -o <prev value>
```

Flag <type> is

**v** for a simple variable

**l** for a list

**s** for a stopped timer

**u** for an up-counting timer

**d** for a down-counting timer.

However, only **v** is implemented in *ns*. Flag **-v <value>** gives the new value of the variable. Variable values are simple ASCII strings obeying the TCL string quoting conventions. List values obey the TCL list conventions. Timer values are ASCII numeric. Flag **-o <prev value>** gives the previous value of the variable. This is used in backward animation. Here is an example of a simple feature event:

```
f -t 0.000000000000000000 -s 1 -n Cl_ -a srm(1) -v 2.25000 -T v
```

Features may be deleted using:

```
f -t <time> -a <agent name> -n <var name> -o <prev value> -X
```

### 45.1.9 Executing Tcl Procedures and External Code from within Nam

There is a special event that can be put in a nam tracefile which allows us to run different tcl procedures. This event is represented by event type **v**.

```
v -t <time> -e <nam tcl procedure call>
```

This event is very generic, in that it may execute several different procedures at a given time, as long as it is in one line (no more than 256 characters). There may be white spaces in the string which are passed as arguments to the tcl procedure. Unlike other events, the order of flags and the tcl procedure call is important.

Here are some examples of this event in use:

### Setting playback speed

Normally the user chooses a playback rate via the rate slider in the animation window. A trace file may set the playback rate via the `set_rate_ext` command:

```
v -t <time> -e set_rate_ext <time> <update-peers?>
```

For example:

```
v -t 2.5 -e set_rate_ext 20ms 1
```

For compatibility with earlier versions of nam, the `set_rate` command is also supported. Instead of specifying the step size directly, you use  $10 \times \log_{10} \langle \text{time-in-seconds} \rangle$ . For example:

```
v -t 2.5 -16.9897000433602 1
```

In order to have more readable files, `set_rate_ext` is preferred.

### Annotation

This procedure is used for displaying text annotation at specific times:

```
v -t <time> -e sim_annotation <time> <unique sequence id> <text to display>
```

For example:

```
v -t 4 -e sim_annotation 4 3 node 0 added one mark
```

This line calls a special tcl function `sim_annotation` in nam, which inserts the given string `node 0 added one mark` into nam's annotation pane. Look at `Animator instproc sim_annotation` in `tcl/annotation.tcl` for the implementation details.

### Node Exec Button

In nam releases, version 1.0a10 and later there is support for running external userdefinable scripts or programs from nam by clicking on a node button.

```
v -t 1.0 -e node_tclscript <node_id> <button label> <tcl script>
```

This line when read in a tracefile will add an extra button to node objects that will execute a tcl script when clicked.

For example:

```
v -t 1.0 -e node_tclscript 2 "Echo Hello" {puts [exec echo hello]}
```

The above line adds a button to node 2's info window with the label "Echo Hello" and when this button is pressed the shell command "echo hello" will be run and it's output will be returned to nam and then output to the terminal via the tcl procedure puts.

The functions that implement the different nam trace formats described above may be found in the following files: *ns/trace.cc*, *ns/trace.h*, *ns/tcl/lib/ns-namsupp.tcl*.

## 45.1.10 Using Streams for Realtime Applications

In addition to reading from files nam can also read from a stream such as STDIN. Here is a little tutorial on how to send a nam trace stream to nam to make it operate with real-time data. First some background on how nam works internally. Basically, it thinks it is reading from a nam tracefile. The file has a format to it. Each line is a nam event. The first character on the line defines the type of event and is followed by several flags to set options on that event. Each event is terminated by a newline character. A nam tracefile has 2 sections, static configuration events and animation events. All events with -t \* in them are configuration events and should be sent to nam in one burst. Lines beginning with a # are comment lines. Currently comments should only be placed in the animation section of the file after the first animation event.

First of all you need to pipe your data to nam's stdin and add the '-' flag to the nam command.

For example:

```
% cat wireless.nam | nam -
```

nam will read the information from stdin

Following is a short wireless animation example. The first part of the script has line with -t \* which tells nam that these are initial configuration information.

```
V -t * -v 1.0a5 -a 0
W -t * -x 600 -y 600
```

The first 2 lines are used in the nam initialization. They need to be the first 2 lines sent to nam from your program. V is the minimum nam version needed to correctly run this script. W means this script contains wireless nodes which will be within the canvas size of width x and height y.

```
n -t * -s 0 -x 0.0 -y 0.0 -z 20 -v circle -c black -w
n -t * -s 1 -x 0.0 -y 200.0 -z 20 -v box -c black -w
```

Next is the network layout information. The first line n creates a wireless (-w) circular (-v circle) node with id 0 (-s 0) at position 0.0,0.0 (-x 0.0 -y 0.0). Its size (-z) is 20 and its color (-c) is black. The second is a wireless square (-v box) node with id 1 (-s 1) at 0.0,200.0. Each node has to have a unique id number given by the -s flag.

```
A -t * -n 1 -p 0 -o 0xffffffff -c 31 -a 1
A -t * -h 1 -m 2147483647 -s 0
```

The A event line has to do with setting up hierarchical addressing in nam. It is necessary in wireless nam because packets are treated as broadcast to every node.

Now we are done with the configuration part of the nam file. Next are the animation events. In order for nam to operate in a close to real-time mode it needs to constantly receive updates. As it is playing it will keep reading lines from the nam trace and playing them back. The sequence of events must be in chronological order. For example the following lines change the color of node 1 from black to green back to black and then to black again.

```
n -t 0.0 -s 1 -S COLOR -c green -o black
n -t 0.01 -s 1 -S COLOR -c black -o green
n -t 0.02 -s 1 -S COLOR -c black -o black
```

Notice that the "-t <time>" flags are always increasing. You cannot issue one event at -t 0.2 and then another later on at -t 0.1. Nam has an internal counter of time and it executes an event once its time counter passes that event time. It will execute events in chronological order but only if they are given to it in chronological order. So the following WILL NOT work.

```
n -t 0.0 -s 1 -S COLOR -c black -o black
n -t 0.02 -s 1 -S COLOR -c green -o black
n -t 0.01 -s 1 -S COLOR -c black -o green
```

Since nam has its own internal representation of time which is different than current real world time you have to try and synchronize them. There is no explicit and totally accurate way to do this but you can have a rough synchronization of time by having your application periodically send nam events even if nothing has happened. We have created a dummy or "no-op" event (T) for this purpose.

```
T -t 0.5
```

As described above, you MUST feed events to nam in non-decreasing timestamp order. Successive events at the same time are OK. Before animating to a given time, nam needs to know that it's got all the events for that time, and so it actually has to read an event AFTER that time. Therefore if you're driving nam from an external process in real-time it will refuse to animate time t until it sees an event at time t+i (for any i > 0). To make nam play out events smoothly, you may therefore need to generate dummy events with intermediate timestamps (so that nam knows it can advance). Events of type "T" are dummy events, so this stream would produce jerky animation:

```
n -t 1.0 -s 1 -S COLOR -c green -o black
n -t 2.0 -s 1 -S COLOR -c black -o green
n -t 3.0 -s 1 -S COLOR -c black -o black
```

while this would be animated much smoother:

```
T -t 0.0
T -t 0.5
n -t 1.0 -s 1 -S COLOR -c green -o black
T -t 1.5
n -t 2.0 -s 1 -S COLOR -c black -o green
T -t 2.5
n -t 3.0 -s 1 -S COLOR -c black -o black
T -t 3.5
T -t 4.0
...
```



If nam ever gets to the end of an event stream it will block and the program will appear as if it has frozen. The screen will not be updated until it can read another event so you must be sure to feed events to nam faster than or as fast as it can read them. This technique works pretty well and allows nam to look as if it is running in real-time although in reality there will be a slight delay which is usually acceptable.

One other thing to remember is that your application should send these events based on it's representation of time from when the application started. Also, when sending events to nam they should be unbuffered so you will want to `fflush(stdout);` after each event.

Another event which you can keep sending to nam would be an note which will show at the bottom of the nam window.

```
v -t 0.08 -e sim_annotation 0.08 1 Time is 0.08
v -t 0.09 -e sim_annotation 0.09 2 Time is 0.09
v -t 0.10 -e sim_annotation 0.08 3 Time is 0.10
```

The 'v' event means that you will execute a tcl script at time -t, everything after -e is the script to execute. `sim_annotation` writes a note at the bottom of the screen. The numbers after it are the time to write and a unique note id. Notice how I incremented the note id with each successive note. The remaining is what is written to the screen. For example "Time is 0.08" followed by "Time is 0.09", etc...

That is the basic idea behind making nam work in a real-time fashion. Following are two examples on how to generate wireless packet animations when using nam. To make a wireless broadcast which is shown as quickly expanding circle you can use the following.

```
+ -t 0.16 -s 0 -d -1 -p AODV -e 100 -c 2 -a 0 -i 0 -k MAC
- -t 0.16 -s 0 -d -1 -p AODV -e 100 -c 2 -a 0 -i 0 -k MAC
h -t 0.16 -s 0 -d -1 -p AODV -e 100 -c 2 -a 0 -i 0 -k MAC
```

'+' event puts the packet onto the transmission queue '-' event remove the packet from the queue and makes it ready to broadcast 'h' send the packet to the next hop which actually causes the animation Here are the meanings of the flags -t time -s transmitting node id -d destination node id (-1 indicates broadcast to world) -e size of transmission -c ultimate destination of the packet

To show a packet being send from one particular node to another use the following

```
r -t 0.255 -s 1 -d -1 -p MAC -e 512 -c 0 -a 0 -i 0 -k MAC
+ -t 0.255 -s 1 -d 0 -p AODV -e 512 -c 0 -a 0 -i 0 -k MAC
- -t 0.255 -s 1 -d 0 -p AODV -e 512 -c 0 -a 0 -i 0 -k MAC
h -t 0.255 -s 1 -d 0 -p AODV -e 512 -c 0 -a 0 -i 0 -k MAC
r -t 0.255 -s 0 -d 1 -p AODV -e 512 -c 0 -a 0 -i 0 -k MAC
```

First the packet is received ('r') from the wireless broadcast to node 1. It is then added to the outgoing queue ('+') on node 1. Next, it is removed from the queue('-') and ready to be sent to node 0. Then the packet is sent to the next hop ('h') node 0. This will trigger an animation of a line the length of the packet size moving from node 1 to node 0. Finally it is received ('r') by node 0 from node 1.

For more nam events you can look at the nam section in the ns manual

Also, you can save a copy of the trace from a realtime source using the unix 'tee' command. For example:

```
% cat wireless.nam | tee saved_tracefile.nam | nam -
```

Sometimes it is a bug in nam and sometimes it is a problem with the way your tracefile is formatted. You may expect nam to do something that it won't do. Part of the philosophy in the design of nam is that the detail of an animation is handled by the tracefile which makes nam more flexible but pushes some of the animation complexity on to the programmer generating the tracefile.

### 45.1.11 Nam Trace File Format Lookup Table

This is a listing of all possible nam trace event codes and the flags associated with them. It was taken from the source code in the file parser.cc. You can generate your own table by running `nam -p`.

# : comment – this line is ignored

T : Dummy event to be used in time synchronization

-t <time> time

n : node

-t	<time>	time
-s	<int>	node id
-v	<shape>	shape (circle, box, hexagon)
-c	<color>	color
-z	<double>	size of node
-a	<int>	address
-x	<double>	x location
-y	<double>	y location
-i	<color>	label color
-b	<string>	label
-l	<string>	label
-o	<color>	previous color
-S	<string>	state (UP, DOWN, COLOR)
-L	<string>	previous label
-p	<string>	label location
-P	<string>	previous label location
-i	<color>	inside label color
-I	<color>	previous inside label color
-e	<color>	label color
-E	<color>	previous label color
-u	<string>	x velocity
-U	<string>	x velocity
-V	<string>	y velocity
-T	<double>	node stop time
-w	<flag>	wireless node

# l : link

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-r	<double>	transmission rate
-D	<double>	delay
-h	<double>	length
-O	<orientation>	orientation
-b	<string>	label
-c	<color>	color
-o	<color>	previous color
-S	<string>	state (UP, DOWN)
-l	<string l>	label
-L	<string>	previous label
-e	<color>	label color
-E	<color>	previous label color

# + : enqueue packet

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	packet color attribute id
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type

# - : dequeue packet

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type

h : hop

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type
-R	<double>	wireless broadcast radius
-D	<double>	wireless broadcast duration

r : receive

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type
-R	<double>	wireless broadcast radius
-D	<double>	wireless broadcast duration

d : drop line

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type

E : session enqueue

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type

D : session dequeue

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type

P : session drop

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-e	<int>	extent
-a	<int>	attribute
-i	<int>	id
-l	<int>	energy
-c	<string>	conversation
-x	<comment>	comment
-p	<string>	packet type
-k	<string>	packet type

a : agent

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-x	<flag>	remove agent
-n	<string>	agent name

f : feature

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-x	<flag>	remove feature
-T	<char>	type
-n	<string>	name
-a	<string>	agent
-v	<string>	value
-o	<string>	previous value

G : group

-t	<time>	time
-n	<string>	name
-i	<int>	node id
-a	<int>	group id
-x	<flag>	remove from group

L : lan link

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-o	<orientation>	orientation
-O	<orientation>	orientation

m : mark node

-t	<time>	time
-n	<string>	name
-s	<int>	node id
-c	<string>	color
-h	<string>	shape (circle, square, hexagon)
-X	<flag>	remove mark

R : routing event

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-g	<int>	multicast group
-p	<packet source>	packet source id or *
-n	<flag>	negative cache
-x	<flag>	this route timed out
-T	<double>	timeout
-m	<string>	mode (iif or oif)

v : execute tcl expression

-t	<time>	time
-e	<tcl expression>	tcl script

V : set trace file version

-t	<time>	time
-v	<string>	time
-a	<int>	time

N : use nam graph

W : wireless range

-t	<time>	time
-x	<int>	X
-y	<int>	Y

g : energy status – for future use

-t	<time>	time
----	--------	------

A : hierarchical address space configuration – initialization only

-t	<time>	time
-n	<int>	hierarchy
-p	<int>	port shift
-o	<hex>	port mask
-c	<int>	multicast shift
-a	<int>	multicast mask
-h	<int>	hierarchy
-m	<int>	node shift
-s	<int>	node mask

c : color table configuration – initialization only

-t	<time>	time
-i	<int>	id
-n	<string>	color

q : create packet queue – initialization only

-t	<time>	time
-s	<int>	source id
-d	<int>	destination id
-a	<orientation>	orientaion

```

X: layout lan
-t      <time>      time
-n      <string>     name
-r      <double>     rate
-D      <double>     delay
-o      <orientation> orientation
-O      <orientation> orientation

```

## 45.2 Ns commands for creating and controlling nam animations

This section describes different APIs in *ns* that may be used to manipulate nam animations for objects like nodes, links, queues and agents. The implementation of most of these APIs is contained in *ns/tcl/lib/ns-nam supp.tcl*. Demonstration of nam APIs may be found in *ns/tcl/ex/nam-example.tcl*.

### 45.2.1 Node

Nodes are created from the "n" trace event in trace file. Each node represents a host or a router. Nam terminates if there are duplicate definitions of the same node. Attributes specific to node are color, shape, label, label-color, position of label and adding/deleting mark on the node. Each node can have 3 shapes: circle (default), square, or hexagon. But once created, the shape of a node cannot be changed during the simulation. Different node may have different colors, and its color may be changed during animation. The following OTcl procedures are used to set node attributes, they are methods of the class Node:

```

$node color [color]      ;# sets color of node
$node shape [shape]      ;# sets shape of node
$node label [label]      ;# sets label on node
$node label-color [lcolor] ;# sets color of label
$node label-at [ldirection] ;# sets position of label
$node add-mark [name] [color] [shape] ;# adds a mark to node
$node delete-mark [name] ;# deletes mark from node

```

### 45.2.2 Link/Queue

Links are created between nodes to form a network topology. *namlinks* are internally simplex, but it is invisible to the users. The trace event "l" creates two simplex links and other necessary setups, hence it looks to users identical to a duplex link. Link may have many colors and it can change its color during animation. Queues are constructed in nam between two nodes. Unlike link, nam queue is associated to a simplex link. The trace event "q" only creates a queue for a simplex link. In nam, queues are visualized as stacked packets. Packets are stacked along a line, and the angle between the line and the horizontal line can be specified in the trace event "q". Commands to setup different animation attributes of a link are as follows:

```

$ns duplex-link-op <attribute> <value>

```

The <attribute> may be one of the following: orient, color, queuePos, label. Orient or the link orientation defines the angle between the link and horizontal. The optional orientation values may be defined in degrees or by text like right (0), right-up (45), right-down (-45), left (180), left-up (135), left-down (-135), up (90), down (-90). The queuePos or position of queue is defined as the angle of the queue line with horizontal. Examples for each attribute are given as following :



```

$ns duplex-link-op orient right      ;# orientation is set as right. The order
                                     ;# in which links are created in nam
                                     ;# depends on calling order of this function.

$ns duplex-link-op color "green"
$ns duplex-link-op queuePos 0.5
$ns duplex-link-op label "A"

```

### 45.2.3 Agent and Features

Agents are used to separate protocol states from nodes. They are always associated with nodes. An agent has a name, which is a unique identifier of the agent. It is shown as a square with its name inside, and a line link the square to its associated node. The following are commands that support agent tracing:

```

$ns add-agent-trace <agent> <name> <optional:tracefile>
$ns delete-agent-trace <agent>
$ns monitor-agent-trace <agent>

```

Once the above command is used to create an agent in nam trace, the `tracevar` method of the `nsagent` can be used to create feature traces of a given variable in the agent. For example, the following code segment creates traces of the variable `C1_` in an SRM agent `$srm(0)`:

```

$ns attach-agent $n($i) $srm(0)
$ns add-agent-trace $srm($i) srm(0)
$ns monitor-agent-trace $srm(0) ;# turn nam monitor on from the start
$srm(0) tracevar C1_

```

### 45.2.4 Some Generic Commands

`$ns color <color-id>` defines color index for nam. Once specified, `color-id` can be used in place of the color name in nam traces.

`$ns trace-annotate <annotation>` inserts an annotation in nam. Notice that if `<annotation>` contains white spaces, it must be quoted using the double quote. An example of this would be `$ns at $time "$ns trace-annotate "Event A happened""`. This annotation appears in the nam window and is used to control playing of nam by events.

`$ns set-animation-rate <timestep>` causes nam to set the animation playout rate to the given timestep value. Time is in seconds, with optional prefixes (for example, 1 is 1 second, or 2ms is 0.002 seconds).

**Part X**

**Other**

## Chapter 46

# Educational use of NS and NAM

This chapter is about using *ns* and *nam* for educational purpose. *ns* is a discrete event simulator and supports various flavors of TCP, many different models of unicast and multicast routing, along with different multicast protocols. It supports mobile networking including local and satellite networks. It also supports applications like web caching. And *ns* uses *nam*, an animation tool, developed in Tcl/Tk, to visualize the simulation packet traces which is created by running *ns* scripts. Thus *ns* and *nam* could be used together to easily demonstrate different networking issues in a classroom environment. In this chapter we'll talk mostly about an educational scripts' database that we have developed. We'll also talk about how to use *nam* to run *namtrace* files.

### 46.1 Using NS for educational purposes

We have developed a web-based interface specifically to cater to the above mentioned educational need of using *ns* in the classrooms. This web-interface is serviced by a database of *ns* scripts that could be used for classroom demonstrations and/or other educational purposes. It can be found at [http://www.isi.edu/nsnam/script\\_inv](http://www.isi.edu/nsnam/script_inv). This page also serves as an interface for uploading or submitting similar scripts to the inventory. So even though we currently have only a few scripts in the inventory to start with, we hope that the inventory will eventually grow in size with script contributions from all of you. In the following paragraphs we shall talk more about this educational scripts' index webpage.

#### 46.1.1 Installing/building/running *ns*

In order to run the educational scripts mentioned in the previous section, you would need to have a running version of *ns* in your machine. The homepage for *ns* is located at <http://www.isi.edu/nsnam/ns>. See *ns-build* page at <http://www.isi.edu/nsnam/ns/ns-build.html> for downloading and building *ns* in your machine. If you want to know about using *ns* to write/run scripts, visit *ns* tutorial for beginners at <http://www.isi.edu/nsnam/ns/tutorial/index.html>.

#### 46.1.2 The educational scripts' inventory page:

The educational script inventory page is located at [http://www.isi.edu/nsnam/script\\_inv](http://www.isi.edu/nsnam/script_inv). It may be used either to search, browse and download one or more simulation scripts (and/or related files like the *namtrace*, screenshot, webpage describing whatever is being demonstrated through the simulation) from the inventory or to submit simulation scripts to the inventory. We discuss both the options in the following paragraphs:

**SEARCH/VIEW/DOWNLOAD NS SCRIPTS:** You could search the database using keywords by going to the “Search database” page. You could also browse through the entire database by going to the “View database” page. The search function is very basic at the present and we hope to extend it as the database begins to grow in size. Each script entry in the database has the following information:

**Name of the script**

**Name of the author, author’s email-id and webpage(if provided)**

**Description of what the simulation does.**

**ns version required to run the script**

**Any other comments about script and**

**The category of the script** Currently we have categories of Application, Transport (TCP and others), Routing (unicast and multicast), Multicast protocols, Queue management, Wireless and Others (to include any other category).

**Other related files** At the right hand down corner of each entry there might be links to a NamTrace file, a screenshot and a webpage for the simulation script, if these files/information have been submitted by the author along with the script.

In order to download any script or any of the related files, simply left-click on the filename and a download dialog box will appear where you can provide the path to download the file to.

**SUBMIT NS SCRIPTS TO INVENTORY:** In case you have simulation scripts suitable for classroom demonstrations, you could submit them to the inventory. You have to **ATLEAST** submit the following in order to successfully upload your script:

**Name of author**

**Author’s E-mailid**

**Name of script (and path to the location of your script) to contribute**

**Brief description of script**

**Version of NS script uses**

**Category for the script**

You may **OPTIONALLY** provide the following along with the above required fields:

**Author’s WebPage**

**Namtrace file (namdump for your script simulation)**

**Screenshot (an image of your nam screen)**

**Webpage (pointer to webpage you may have for the script)**

**Other comments, if any**

*Important:* We suggest that you provide the namtrace file along with your script since many users might want to use the namtrace only, for visualization, and download script only when they want to make any changes to the simulation.

## 46.2 Using NAM for educational purposes

In order to visualize an ns simulation, you need to have the NAM tool installed. You could either simply download the nam binary for your platform or download the nam distribution and build in your machine. The link for getting nam binaries as well as nam source is <http://www.isi.edu/nsnam/nam> which also happens to be the nam homepage.

*Steps to use nam in powerpoint:* After opening powerpoint, go to “Slide Show” (visible on the top menu) and click on “action buttons”. Select the type of button you prefer. This would create a button for your slide. Clicking on your button will pop up an “Action Setting” window. Inside the window, there is a place called “Run Program” where you can define your nam program to run.

# Bibliography

- [1] C. Alaettinoglu, A.U. Shankar, K. Dussa-Zeiger, and I. Matta. Design and implementation of MaRS: A routing testbed. *Internetworking: Research and Experience*, 5:17–41, 1994.
- [2] Sandeep Bajaj, Lee Breslau, Deborah Estrin, Kevin Fall, Sally Floyd, Padma Haldar, Mark Handley, Ahmed Helmy, John Heidemann, Polly Huang, Satish Kumar, Steven McCanne, Reza Rejaie, Puneet Sharma, Kannan Varadhan, Ya Xu, Haobo Yu, and Daniel Zappala. Improving simulation for network research. Technical Report 99-702b, University of Southern California, March 1999. (revised September 1999).
- [3] Paul Barford and Mark Crovella. Generating representative web workloads for network and server performance evaluation. In *Proceedings of the ACM SIGMETRICS*, pages 151–160, June 1998.
- [4] L.S. Brakmo, S. O'Malley, and L.L. Peterson. TCP vegas: New techniques for congestion detection and avoidance. In *Proceedings of the ACM SIGCOMM*, pages 24–35, October 1994.
- [5] L.S. Brakmo, S. O'Malley, and L.L. Peterson. TCP vegas: New techniques for congestion detection and avoidance. Technical Report TR 94 04, Department of Computer Science, The University of Arizona, Tucson, AZ 85721, February 1994.
- [6] R. Brown. Calendar queues: A fast  $O(1)$  priority queue implementation for the simulation event set problem. *Communications of the ACM*, 31(10):1220–1227, October 1988.
- [7] Pei Cao and Chengjie Liu. Maintaining strong cache consistency in the World-Wide Web. In *Proceedings of the IEEE ICDCS*, pages 12–21, May 1997.
- [8] N. Christin, J. Liebeherr, and T. Abdelzaher. A quantitative assured forwarding service. In *Proceedings of IEEE INFOCOM 2002*, volume 2, pages 864–873, New York, NY, June 2002.
- [9] S. Deering, D. Estrin, D. Farinacci, V. Jacobson, Ching-Gung Liu, and L. Wei. An architecture for wide-area multicast routing. Technical Report USC-SC-94-565, Computer Science Department, University of Southern California, Los Angeles, CA 90089., 1994.
- [10] Anja Feldmann, Anna C. Gilbert, Polly Huang, and Walter Willinger. Dynamics of IP traffic: A study of the role of variability and the impact of control. pages 301–313, Cambridge, MA, USA, August 1999.
- [11] S. Floyd, V. Jacobson, C. Liu, S. McCanne, and L. Zhang. A reliable multicast framework for light-weight sessions and application level framing. In *Proceedings of the ACM SIGCOMM*, pages 342–356, August 1995.
- [12] H. T. Friis. A note on a simple transmission formula. *Proc. IRE*, 34, 1946.
- [13] H. W. Hethcote. The mathematics of infectious diseases. *SIAM Review*, 42(4):599–653, October 2000.
- [14] A. Heybey. *Netsim Manual*. MIT, 1989.
- [15] R. Jain. *The Art of Computer Systems Performance Analysis*. John Wiley and Sons, Inc., 1996.
- [16] Pierre L'Ecuyer. Good parameters and implementations for combined multiple recursive random number generators. *Operations Research*, 47(1):159–164, 1999.

- [17] Pierre L'Ecuyer. Software for uniform random number generation: Distinguishing the good and the bad. In *Proceedings of the 2001 Winter Simulation Conference*, pages 95–105, December 2001.
- [18] Pierre L'Ecuyer, Richard Simard, E. Jack Chen, and W. David Kelton. An object-oriented random number package with many long streams and substreams. *Operations Research*, 2001.
- [19] A. Legout and E.W. Biersack. PLM: Fast convergence for cumulative layered multicast transmission schemes. In *Proceedings of the ACM SIGMETRICS*, Santa Clara, CA, U.S.A., June 2000.
- [20] J. Liebeherr and N. Christin. JoBS: Joint buffer management and scheduling for differentiated services. In *Proceedings of IWQoS 2001*, pages 404–418, Karlsruhe, Germany, June 2001.
- [21] J. Liebeherr and N. Christin. Rate allocation and buffer management for differentiated services. *Computer Networks*, 40(1):89–110, September 2002.
- [22] M. Mathis and J. Mahdavi. Forward acknowledgement: Refining TCP congestion control. In *Proceedings of the ACM SIGCOMM*, August 1996.
- [23] M. Mathis, J. Mahdavi, S. Floyd, and A. Romanov. *TCP Selective Acknowledgement Options*, RFC 2018 edition, 1996.
- [24] S. McCanne and S. Floyd. ns—Network Simulator. <http://www-mash.cs.berkeley.edu/ns/>.
- [25] S. McCanne and V. Jacobson. The bsd packet filter: A new architecture for user-level packet capture. pages 259–269, January 1993.
- [26] John Ousterhout. Scripting: Higher-level programming for the 21st century. *IEEE Computer*, 31(3):23–30, March 1998.
- [27] S.K. Park and R.W. Miller. Random number generation: Good ones are hard to find. *Communications of the ACM*, 31(10):1192–1201, October 1988.
- [28] Peter Piatek, Jeremy Ethridge, Mandeep Baines, and Farhan Shallwani. *A Network Simulator, Differentiated Services Implementation*. Open IP, Nortel Networks, 2000.
- [29] T. S. Rappaport. *Wireless communications, principles and practice*. Prentice Hall, 1996.
- [30] D. Waitzman, C. Partridge, and S.E. Deering. *Distance Vector Multicast Routing Protocol*, RFC 1075 edition, 1988.
- [31] Dina Katabi, Mark Handley, and Charlie Rohrs. Congestion Control for High Bandwidth-Delay Product Networks ACM Sigcomm 2002