

BookSim 2.0 User's Guide

Nan Jiang, George Michelogiannakis, Daniel Becker, Brian Towles and William J. Dally

May 7, 2013

Contents

1	Introduction	1
2	Getting started	2
2.1	Downloading and building the simulator	2
2.2	Running a simulation	2
2.3	Simulation output	2
3	Example	3
4	Configuration parameters	3
4.1	Topologies	6
4.2	Physical sub-networks	7
4.3	Routing algorithms	7
4.4	Flow control	7
4.5	Router organizations	7
4.5.1	The input-queued router	7
4.5.2	The event-driven router	8
4.6	Allocators	8
4.7	Traffic	9
4.7.1	Injection mode	9
4.7.2	Request-reply traffic	9
4.7.3	Traffic patterns	9
4.8	Simulation parameters	10
A	Random number generation	10

1 Introduction

This document describes the use of the BookSim interconnection network simulator. The simulator is designed as a companion to the textbook “Principles and Practices of Interconnection Networks” (PPIN) published by Morgan Kaufmann (ISBN: 0122007514) and it is assumed that its reader is familiar with the material covered in that text.

This user guide is fairly brief as, with most simulators, the best way to learn and *understand* the simulator is to study the code. Most of the simulator’s components are designed to be modular so tasks such as adding a new routing algorithm, topology, or router microarchitecture should not require a complete redesign of the code. Once you have downloaded the code, compiled it, and run

a simple example (Section 2), the more detailed examples of Section 3 give a good overview of the capabilities of the simulator. A list of configuration options is provided in Section 4 for reference.

2 Getting started

2.1 Downloading and building the simulator

The latest official release of the simulator can be checked out from our subversion (SVN) repository. Make sure a subversion client is installed on your machine; under UNIX, you can use the following command to check out a working copy:

```
svn co https://nocs.stanford.edu/cgi-bin/svn.cgi/booksim2.0
```

The simulator files should now be in the `booksim2.0/trunk/src/` directory. The simulator itself is written in C++ and has been specifically tested with the GNU G++ compiler (version ≥ 3). The front end of the simulator uses LEX and YACC generated parser to process the simulator configuration file; however, unless you plan on making changes to the front end parser, LEX and YACC are not needed. The `Makefile` should be edited so that the first few lines reflect the correct paths to the tools for your particular system. The default `Makefile` should work on the Stanford Leland machines. Type `make` to build the simulator.

A note for Windows users: The above instructions have been tested to work with Cygwin 1.7.18.

2.2 Running a simulation

The simulator is invoked using the following command line:

```
./booksim [configfile]
```

The parameter `configfile` is a file that contains configuration information for the simulator. So, for example, to simulate the performance of a simple 8×8 torus (8-ary 2-cube) network on uniform traffic, a configuration such as the one shown in Figure 1 could be used. This particular example configuration file can be found in the `examples/torus88` directory.

In addition to specifying the topology, the configuration file also contains basic information about the routing algorithm, flow control, and traffic. This simple example uses dimension-order routing and four virtual channels. The `injection_rate` parameter is added to tell the simulator to inject (on average) 0.15 packets per simulation cycle per node. Packet size defaults to a single flit. Also, any line of the configuration file that begins with `//` is treated as a comment and ignored by the simulator. A detailed list of configuration parameters is given in Section 4. Any parameters not specified by the user will take on default values. The default value for every parameter in the simulator is specified in the file `booksim_config.cpp`.

2.3 Simulation output

Continuing our example, running the torus simulation produces the output shown in Figure 2. Each simulation has three basic phases: warm up, measurement, and drain. The length of the warm up and measurement phases is a multiple of a basic sample period (defined by `sample_period` in the configuration). As shown in the figure, the current latency and throughput (rate of accepted packets) for the simulation is printed after each sample period. The overall throughput is determined by the lowest throughput of all the destination in the network, but the average throughput is also displayed.

```
// Topology
topology = torus;
k        = 8;
n        = 2;

// Routing
routing_function = dim_order;

// Flow control
num_vcs = 4;

// Traffic
traffic      = uniform;
injection_rate = 0.15;
```

Figure 1: Example configuration file for simulating a 8-ary 2-cube network.

After the warm up periods have passed (default to $3 \times \text{sample_period}$), the simulator prints the “Warmed up” message and resets all the simulation statistics. Then, the measurement phase begins and statistics continue to be reported after each sample period. The measurement phase typically last for last least 3 sample periods. Once the measurement periods have passed, all the measurement packets are drained from the network before final latency and throughput numbers are reported. Details of the configuration parameters used to control the length of the simulation phases are covered in Section 4.8.

3 Example

One of the most basic performance measures of any interconnection network is its latency versus offered load. Figure 3 shows a simple configuration file for making this measurement in a 8-ary 2-mesh network under the transpose traffic pattern. This configuration was used to generate a single data point in Figure 25.2 in PPIN. The particular configuration accounts for some small delays and pipelining of the input-queued router and also introduces a small input speedup to account for any inefficiencies in allocation. By running simulations for many increments of `injection_rate`, the average latency curve can be found. Then, to compare the performance of dimension-order routing against several other routing algorithms, for example, the `routing_function` option can be changed.

4 Configuration parameters

All information used to configure a simulation is passed through a configuration file as illustrated by the example in Section 2.2. This section lists the major configuration parameters for the simulator. Additional description can be found in the configuration parameter class file `booksim_config.cpp`. A user can incorporate additional options by changing the this file.

```
BEGIN Configuration File: examples/torus88

....

// Topology
topology = torus;
k = 8;
n = 2;
// Routing
routing_function = dim_order;
// Flow control
num_vcs = 2;
// Traffic
traffic = uniform;
injection_rate = 0.15;

END Configuration File: examples/torus88
Class 0:
Packet latency average = 33.2807
minimum = 7
maximum = 73

....

Warmed up ...Time used is 2000 cycles

....

Draining all recorded packets ...
Draining remaining packets ...
Time taken is 5104 cycles
===== Overall Traffic Statistics =====
===== Traffic class 0 =====
Packet latency average = 33.6964 (1 samples)
minimum = 7 (1 samples)
maximum = 84 (1 samples)

....

Total run time 1.01837
```

Figure 2: Simulator output from running the `examples/torus88` configuration file.

```
// Topology

topology = mesh;
k = 8;
n = 2;

// Routing
routing_function = dor;

// Flow control
num_vcs      = 8;
vc_buf_size  = 8;
wait_for_tail_credit = 1;

// Router architecture
vc_allocator = islip;
sw_allocator = islip;
alloc_iters  = 1;

credit_delay  = 2;
routing_delay = 0;
vc_alloc_delay = 1;
sw_alloc_delay = 1;

input_speedup    = 2;
output_speedup   = 1;
internal_speedup = 1.0;

// Traffic
traffic = transpose;
packet_size = 20;

// Simulation
sim_type = latency;

injection_rate = 0.005;
```

Figure 3: A typical configuration file (`examples/mesh88_lat`) for creating a latency versus offered load curve for a 8-ary 2-mesh network.

4.1 Topologies

The `topology` parameter determines the underlying topology of the network. There is also a set of numerical parameters that describes the size of the networks.

<code>k</code>	Network radix, the number of routers per dimension
<code>n</code>	Network dimension
<code>c</code>	Network concentration, the number of nodes sharing a single router. Typically set to 1, ≥ 1 only in networks that has concentration (i.e. <code>cmesh</code>).
<code>x</code>	(NoC simulations only) The number of routers in the X dimension. Used to calculate channel latency between routers.
<code>y</code>	(NoC simulations only) The number of routers in the y dimension. Used to calculate channel latency between routers.
<code>xr</code>	(NoC simulations only) For networks that have $c \geq 1$, the number of nodes in the x direction per router. Used to calculate channel latency between routers.
<code>yr</code>	(NoC simulations only) For networks that have $c \geq 1$, the number of nodes in the y direction per router. Used to calculate channel latency between routers.

The channel latency of the network must be configured within the source code of the topology files. All topologies by default have channel latency of 1 cycles. Topologies available in BookSim are:

<code>fly</code>	A k -ary n -fly (butterfly) topology. The <code>k</code> parameter determines the network's radix and the <code>n</code> parameter determines the network's dimension. Note: a k -ary 1-fly is essentially a single radix- k router, useful for testing.
<code>mesh</code>	A k -ary n -mesh (mesh) topology. The <code>k</code> parameter determines the network's radix and the <code>n</code> parameter determines the network's dimension.
<code>torus</code>	A k -ary n -cube (torus) topology. The <code>k</code> parameter determines the network's radix and the <code>n</code> parameter determines the network's dimension.
<code>cmesh</code>	Concentrated mesh topology is a k -ary n -mesh topology with multiple nodes sharing a single router. The <code>c</code> determines the concentration. The <code>cmesh</code> topology has the option that turns on "express channels" as described in by default these channels are turned off.
<code>fat tree</code>	Fat Tree topology with 3 levels. Nodes are routers are arranged in a tree structure but the number of links between levels stays constant. At the bottom <code>k</code> nodes shares a level 0 router.
<code>flattened butterfly</code>	A topology based on the paper "Flattened butterfly: a cost-efficient topology for high-radix networks" ISCA 2007
<code>dragonfly</code>	A topology based on the paper "Technology-driven, highly-scalable dragonfly topology." ISCA 2008
<code>quad tree</code>	A quad tree topology.
<code>tree 4</code>	
<code>anynet</code>	A topology based on an user input file specifying connectivity of nodes and routers.

4.2 Physical sub-networks

The `physical_subnetworks` parameter defines the number of physical sub-networks present in the network (defaults to one). All sub-networks receive the same configuration parameters and thus are identical. Traffic sources maintain an injection queue for each sub-network. The packet generation process is unaffected. It enqueues generated packets into the proper sub-network queue according to a division function in the traffic manager. At every cycle, flits at the head of each queue attempt to be injected. Traffic destinations can eject one flit from each sub-network each cycle.

4.3 Routing algorithms

The `routing_function` parameter selects a routing algorithm for the topology. Many routing algorithms need multiple virtual channels for deadlock freedom. In addition to `routefunc.cpp`, some topologies source files include additional routing functions. Also, the simulator code is structured so that additional routing algorithms can be added with minimal changes to the overall simulator (see the `routefunc.cpp` file in the simulator's source code).

4.4 Flow control

The simulator supports basic virtual-channel flow control with credit-based backpressure.

`num_vcs` The number of virtual channels per physical channel.

`vc_buf_size` The depth of each virtual channel in flits.

`wait_for_tail_credit` If non-zero, do not reallocate a virtual channel until the tail flit has left that virtual channel. This conservative approach prevents a dependency from being formed between two packets sharing the same virtual channel in succession.

4.5 Router organizations

The simulator also supports two different router microarchitectures. The input-queued router follows the general organization described in PPIN while the event-driven router is modeled after the router used in the Avici TSR and described in U.S. Patent 6,370,145. The microarchitecture is selected using the `router` option. Also, both routers share a small set of options.

`credit_delay` The processing delay (in cycles) for a credit. Does not include the wire delay for transmitting the credit.

`internal_speedup` An arbitrary speedup of the internals of the routers over the channel transmission rate. For example, a speedup 1.5 means that, on average, 1.5 flits can be forwarded by the router in the time required for a single flit to be transmitted across a channel. Also, the configuration parser expects a floating point number for this field, so integer speedups should also include a decimal point (e.g. "2.0").

4.5.1 The input-queued router

The input-queued router (`router = iq`) follows the pipeline described in PPIN of route computation, virtual-channel allocation, switch allocation, and switch traversal. There are several options specific to the input-queued router.

<code>input_speedup</code>	An integer speedup of the input ports in space. A speedup of 2, for example, gives each input two input ports into the crossbar. Access to these ports is statically allocated based on the virtual channel number: virtual channel v at input i is connected to port $i \cdot s + (v \bmod s)$ for an input speedup of s .
<code>output_speedup</code>	An integer speedup of the output ports in space. Similar to <code>input_speedup</code>
<code>routing_delay</code>	The delay (in cycles) of route computation.
<code>hold_switch_for_packet</code>	
<code>speculative</code>	Enable speculative switch allocation (i.e., allow switch allocation to occur in parallel with VC allocation for header flits).
<code>alloc_iters</code>	For the <code>islip</code> , <code>pim</code> and <code>select</code> allocators, allocation can be improved by performing multiple iterations of the algorithm; this parameter controls the number of iterations to be performed for both switch and VC allocation.
<code>arb_type</code>	If the VC or switch allocator is a separable input- or output-first allocator, this parameter selects the type of arbiter to use.
<code>sw_allocator</code>	The type of allocator used for switch allocation. See Section 4.6 for a list of the possible allocators.
<code>sw_alloc_delay</code>	The delay (in cycles) of switch allocation.
<code>vc_allocator</code>	The type of allocator used for virtual-channel allocation. See Section 4.6 for a list of the possible allocators.
<code>vc_alloc_delay</code>	The delay (in cycles) of virtual-channel allocation.

4.5.2 The event-driven router

The event-driven router (`router = event`) is a microarchitecture designed specifically to support a large number of virtual channels (VCs) efficiently. Instead of continuously polling the state of the virtual channels, as in the input-queued router, only changes in VC state are tracked. The efficiency then comes from the fact that the number of state changes per cycle is constant and independent of the number of VCs.

4.6 Allocators

Many of the allocators used in the simulator are configurable (see the input-queued router in Section 4.5.1) and several allocation algorithms are available.

<code>max_size</code>	Maximum-size matching.
<code>islip</code>	iSLIP separable allocator.
<code>pim</code>	Parallel iterative matching separable allocator.
<code>loa</code>	Lonely output allocator.
<code>wavefront</code>	Wavefront allocator.

`separable_input_first` Separable input-first allocator.

`separable_output_first` Separable output-first allocator.

`select` Priority-based allocator. Allocation is performed as in iSLIP, but with preference towards higher priority packets.

4.7 Traffic

4.7.1 Injection mode

The rate at which packets are injected into the simulator is set using the `injection_rate` option. The simulator's cycle time is a flit cycle, the time it takes a single flit to be injected at a source, and the injection rate is specified in packets per flit cycle. For example, setting `injection_rate = 0.25` means that each source injects a new packet in one out of every four simulator cycles. The unit of `injection_rate` can optionally be changed to flits per cycle by setting `injection_rate_uses_flits` to 1.

The injection process can further be specified as either Bernoulli process (`injection_process = bernoulli`) or an on-off process (`injection_process = on_off`). The burstiness of the latter is controlled via the `burst_alpha` and `burst_beta` parameters. See PPIN Section 24.2.2 for a description of the on-off process and its parameters.

4.7.2 Request-reply traffic

By default, all packets that are injected into the network have the same, fixed length. The number of flits per packet is set using the `const_flits_per_packet` option.

Alternatively, the simulator can be configured to generate request-reply traffic. In this mode, the traffic manager injects read and write requests into the network; when a destination node receives such a request packet, it returns a reply packet of the same type. Injection of reply packets has priority over injection of new requests packets. Consequently, the effective overall network load is generated by both the injected requests and the automatically generated replies. Request-reply traffic ignores the `const_flits_per_packet` option; instead, packet sizes are determined by the `{read|write}_{request|reply}_size` options. Furthermore, the mapping of packet types to VCs can be customized using the `{read|write}_{request|reply}_{begin|end}_vc` options.

4.7.3 Traffic patterns

The simulator also supports several different traffic patterns that are specified using the `traffic` option. To describe these patterns, we use the same notation of PPIN Section 3.2: s_i (d_i) denotes the i^{th} bit of the source (destination) address whereas s_x (d_x) denotes the x^{th} radix- k digit of the source (destination) address. The bit length of an address is $b = \log_2 N$, where N is the number of nodes in the network.

<code>uniform</code>	Each source sends an equal amount of traffic to each destination (<code>traffic = uniform</code>).
<code>bitcomp</code>	Bit complement. $d_i = \neg s_i$.
<code>bitrev</code>	Bit reverse. $d_i = s_{b-i-1}$.
<code>shuffle</code>	$d_i = s_{i-1 \bmod b}$.

<code>transpose</code>	$d_i = s_{i+b/2 \bmod b}$.
<code>tornado</code>	$d_x = s_x + \lceil k/2 \rceil - 1 \bmod k$.
<code>neighbor</code>	$d_x = s_x + 1 \bmod k$.
<code>randperm</code>	Random permutation. A fixed permutation traffic pattern is chosen uniformly at random from the set of all permutations. The seed used to generate this permutation is set by the <code>perm_seed</code> option. So, randomly selecting values for <code>perm_seed</code> gives a random sampling of permutations while a fixed value of <code>perm_seed</code> allows the same permutation to be used for several experiments.

4.8 Simulation parameters

The duration and other aspects of a simulation are controlled using the set of simulation parameters.

<code>sim_type</code>	A simulation can either focus on <code>throughput</code> or <code>latency</code> . The key difference between these two types is that a <code>latency</code> simulation will wait for all measurement packets to drain before ending the simulation to ensure an accurate latency measurement. In <code>throughput</code> simulations, this final drain step is eliminated to allow simulation of networks operating beyond their saturation point.
<code>sample_period</code>	The sample period is expressed in simulator cycles and is used as a multiplier when specifying the warm-up length of a simulation and the maximum number of samples. Also, intermediate statistics are displayed once every <code>sample_period</code> cycles. This is only applicable in injection mode.
<code>warmup_periods</code>	The length of the simulator warm up expressed as a multiple of the <code>sample_period</code> . After warming up, all statistics counters are reset. This is only applicable in injection mode.
<code>max_samples</code>	The total length of simulation expressed as a multiple of the <code>sample_period</code> . This is only applicable in injection mode.
<code>latency_thres</code>	If the sampled latency of the current simulation exceeds <code>latency_thres</code> , the simulation is immediately ended.
<code>sim_count</code>	The number of back-to-back simulations to run for the given configuration. Useful for creating ensemble averages of particular statistics.
<code>seed</code>	A random seed for the simulation.
<code>print_activity</code>	At the end of a simulation using <code>iq_router</code> , print out the activity for buffer, switch, and channel of the network.
<code>watch_file</code>	Specific flits can have their "watch" status turn on. Require input a file which has flit id listed. 1 id per line.

A Random number generation

The simulator uses Knuth's integer and floating point pseudorandom number generators. These algorithms and their explanations appear in "The Art of Computer Programming: Seminumerical Algorithms".