



There's lots of data. Gene micro-arrays, once done completely by hand are now churned out by armies of robots. There's more sequence data than ever, they even have a USB stick for it. Of course, there's the one everyone is familiar with, web search.

















To make the rest of the work more concrete, we'll describe briefly the data-flow / streaming framework RaftLib which is used for all of our experiments. We'll start by talking about this simple "sum" application which was first used as an example of a data flow application by Dennis (doi: 10.1007/3-540-06859-7_145).

```
Stream Processing
template< typename A,
          typename B,
         kypename C > class sum : public raft::kernel
public:
  sum() : raft::kernel()
   4
      input.addPort< A >( "input_a" );
      input.addPort< B >( "input_b" );
     output.addPort< C >( "sum" );
  }
  virtual raft::kstatus run()
      /** look at mem on head of queue for a & b, no copy **/
     auto a( input[ "input_a" ].pop_s< A >() );
      auto b( input[ "input_b" ].pop_s< B >() );
     /** allocate mem directly on queue **/
      auto c( output[ "sum" ].allocate_s< C >() );
      (*c) = (*a) + (*b);
      /** mem automatically freed upon scope exit **/
      return( raft::proceed );
  }
};
```

The constructor (there are more efficient ways to declare ports, these used for clarity) declares two input ports "input_a" and "input_b," and one output port "sum." The second function "run()" is the worker which is called by the scheduler. It takes data from two input ports when it is available and pops an item from each input port and writes the sum to the output port. The return value indicates that nothing has happened to warrant exiting the program, although the program will exit on its own with it is provable that there is no further input available.



The constructor (there are more efficient ways to declare ports, these used for clarity) declares two input ports "input_a" and "input_b," and one output port "sum." The second function "run()" is the worker which is called by the scheduler. It takes data from two input ports when it is available and pops an item from each input port and writes the sum to the output port. The return value indicates that nothing has happened to warrant exiting the program, although the program will exit on its own with it is provable that there is no further input available.



The constructor (there are more efficient ways to declare ports, these used for clarity) declares two input ports "input_a" and "input_b," and one output port "sum." The second function "run()" is the worker which is called by the scheduler. It takes data from two input ports when it is available and pops an item from each input port and writes the sum to the output port. The return value indicates that nothing has happened to warrant exiting the program, although the program will exit on its own with it is provable that there is no further input available.



Example of Boyer-Moore string search topology



implementation of the aho-corasick as a string searching library



implementation of the aho-corasick as a string searching library



implementation of the aho-corasick as a string searching library



We can make all kinds of pipeline/task parallel topologies without explicit split / join. This is the strength of stream processing in that we break the mold of the fork/join model. This lack of explicit synchronization gives stream processing a unique ability to exploit extreme levels of parallelism.



At the top there's an example of a simple streaming application. Each stream can be modeled as a queue. At bottom is an example of the queue activity of our streaming application example. The x-axis is the queue position and the y-axis represents the # of cycles occupied within each time frame. The front that is stable is what we're interested in, can we find this quickly? One way to do that s to use queueing models, but they require some idea of service rate.



At the top there's an example of a simple streaming application. Each stream can be modeled as a queue. At bottom is an example of the queue activity of our streaming application example. The x-axis is the queue position and the y-axis represents the # of cycles occupied within each time frame. The front that is stable is what we're interested in, can we find this quickly? One way to do that s to use queueing models, but they require some idea of service rate.



Does buffer sizing have an impact on overall throughput? YES! The front from the previous slide corresponds to approximately 80 kB on this slide. Any smaller and we stifle performance. Too much larger and we start loosing performance, but why?



50 MB at the far right, green dots are page faults, blue and orange dots are L1/L2 misses respectively. The basic premise is that you end up with quite a bit of locality with small buffers that are fairly cacheable, but going above a certain size eliminates the possibility that most of the buffer can end up in the cache. Big is often good for performance but too big is bad. On shared systems with lots of executing threads this can be very bad.



Non-blocking service rate can also be useful for partitioning an application between compute resources. Offline heuristics work pretty well for providing a starting partition, they don't work well online. Most are too slow, in general partitioning is NP-Hard....There are plenty of decent heuristics, we're going to focus elsewhere so just keep in mind that this is another potential usage of service rate.



In prior work I introduced using gain/loss flow models for calculating the throughput through a queueing network. The one thing that we couldn't get at the time was the *mu* on the slide (orange), our method of online service rate determination enables the use of this method during execution for things like thread migration decisions.



Here's the old way of figuring out how fast a compute kernel could execute outside of its network. Each kernel is characterized on its intended compute platform with its intended environment. It takes time.



Every time we change the assignment of a kernel to a compute resource, we have to re-characterize it. This takes a lot of time.



for huge compute graphs individual characterization (necessary for accurate modeling) is really not feasible.



The monitor thread takes samples of non-blocking reads and writes from the queue it is observing. We process these as a small window, saving only very small bits of summary data which is used to estimate the service rate



We want to find the segment given by A, in the instant that the middle worker has an opening to add stars, then we can figure out how fast he can execute un-encumbered by the last worker.



In a high utilization real M/M/1 system, this is what it looks like. For the most part, the queue is highly occupied. Occasionally though we can find segments (colored in red) which are amenable for determining the service rate.



So how probable are these segments (red from previous slide) that are needed to determine the service rate online? Well, not to likely. As the service rate increases, the probability decreases. So we need small sampling frames to increase the likelihood that we'll see the segments in red.



So we need accurate timing. How accurate, well, we want as accurate as possible. One measure that we're interested in is back to back execution, over millions of executions (averaged) the rdtsc instruction has far less latency than the standard clock_gettime which is no huge surprise.



With the timer thread executing on a single core, the latency to access the updated timer values on other cores differ (especially when going to another socket). Every update on the local core invalidates the value on the remote core forcing a QPI access of the newly updated value. Prefetching seems the likely solution, however it doesn't quite fix things so we allocate memory on the other core's NUMA node and prefetch it so that there is no aliasing and the most up to date values are more likely to be in cache speeding access. This gets us closer to the 10nanosecond access time that we see on the local core.



The @ symbol in the bottom left hand corner is the minimal resolution of back to back timer calls (in this case RDTSC) averaged across cores. We need a stable time frame so that means that we need to go to the right (on the x-axis) towards larger multiples of the system timer. So, now we have an issue. We have to find the smallest time frame possible, but also the most stable time frame possible. This is done the first time RaftLib starts up, and a profile is written. At startup this profile is quickly verified (may change depending on the dynamic environment), and then the profile is used if it is acceptable (otherwise we search for the time again).



These are the raw values of reads that our instrumentation thread views. At first glance it appears that there is a nice front right where our expected nonblocking service rate is (red dashed line). The key is understanding that there are still hundreds of values above the red line. A quantile based approach is the obvious one, and it is the one we ultimately took. The issue with quantiles is that we can't take quantiles of an arbitrary distribution without saving lots of state. We only want to save a few values, so we need another approach. Observing that each one of these points is in fact a sum of non-blocked reads (although realizing that due to the fact that we use no locking or atomic accesses during the gathering of this data that there are many potential outliers within the data representing something less or greater than the non-blocking rates that we're searching for). Sums of observations of a random variable tend towards a Gaussian, and we can use this.



Using a closed form solution for the continuous Gaussian gives a solution to our data saving problem (we don't want to transmit continuously and we want something that will fit in a typical L1 or L2 cache line so that our instrumentation can be quick). But our data is still noisy.



Using a closed form solution for the continuous Gaussian gives a solution to our data saving problem (we don't want to transmit continuously and we want something that will fit in a typical L1 or L2 cache line so that our instrumentation can be quick). But our data is still noisy.



Applying a filter over only the previous 16 values gives a less noisy view (shown on the qq plot before and after)



Histogram view of the quantile we want, the rest below that is assumed to be other stuff (side effect of the atomic-less data collection)



Ok, this gets us almost to where we want. But it's not so stable



Using a streaming mean gets us to a more stable result



We use a Laplacian Gaussian Filter to filter the standard deviation of the quantile estimate, commonly used in edge detection to tell when a "stable" service rate has been found. The above is a plot of the Laplacian filtered standard deviation over time, the dotted line shows the cut-off point. When we cut-off, we can re-start the instrumentation and find another service rate (which could have changed during execution).



Instrumentation provided service rates from micro-benchmark that shifts the service rate of B halfway through its execution (in elements, not time).



Going the other way, same benchmark...different shift.

RaftLib

C++ Streaming Template Library

Auto-parallelizes code

Manages resources, buffers, TCP links

GOAL: Automatically Optimized Online

software download: http://raftlib.io

Washington

University in St. Loui



SBS

33

Online Instrumentation

- Queue Occupancy (mean, histogram, etc.)
- Service Rate (non-blocking and actual throughput)
- Process Distribution**
- Less than 1% impact on processor load on average with our implementation
- Execution times not affected by instrumentation by any statistically significant measure
- Can be turned on and off dynamically through queue reallocation process
 34

**Currently only supported in experimental branch and with method of moments, eventually will migrate once I've explored using kernel methods vs. moments since the moments is a bit expensive to compute still.



