

RAFTLIB

Presented by: Jonathan Beard
To: C++Now 2016

RAFTLIB

Alternate Titles:

This thing I started on a plane

RAFTLIB

Alternate Titles:

This thing I started on a plane
What's this RaftLib Thingy?

RAFTLIB

Alternate Titles:

This thing I started on a plane

What's this RaftLib Thingy?

OMG, Another Threading Library...

RAFTLIB

Alternate Titles:

This thing I started on a plane

What's this RaftLib Thingy?

OMG, Another Threading Library...

Why I hate parallel programming

RAFTLIB

Alternate Titles:

This thing I started on a plane

What's this RaftLib Thingy?

OMG, Another Threading Library...

Why I hate parallel programming

A self help guide for pthread anxiety



All thoughts, opinions are my own.
RaftLib is not a product of ARM Inc.
Please don't ask about ARM products or
strategy. I will scowl and not answer.

Thank you 😎

ABOUT ME

my website

<http://www.jonathanbeard.io>



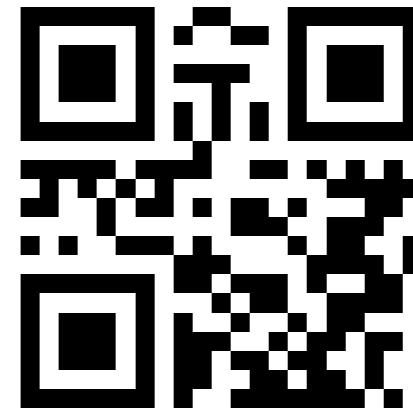
slides at

<http://goo.gl/cwT5UB>

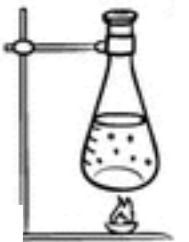
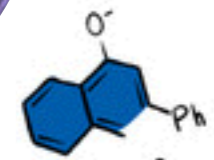
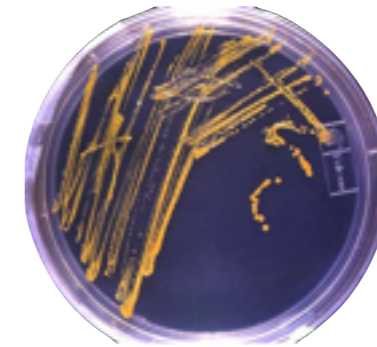
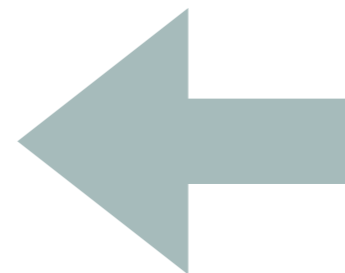
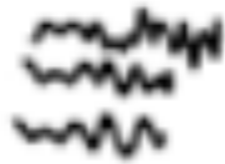
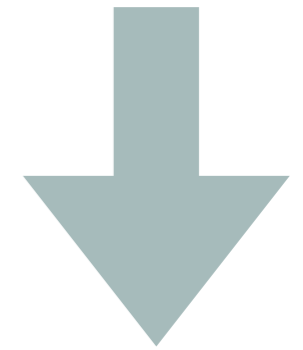
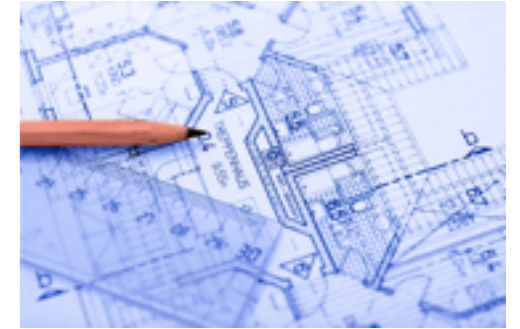
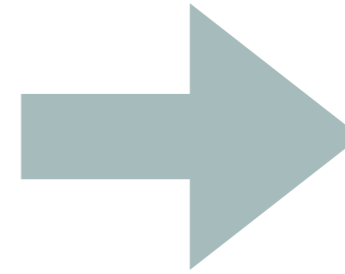
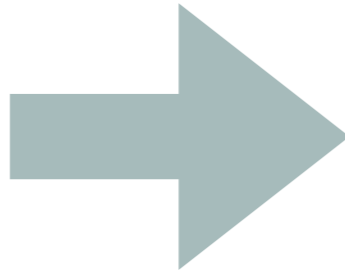


project page

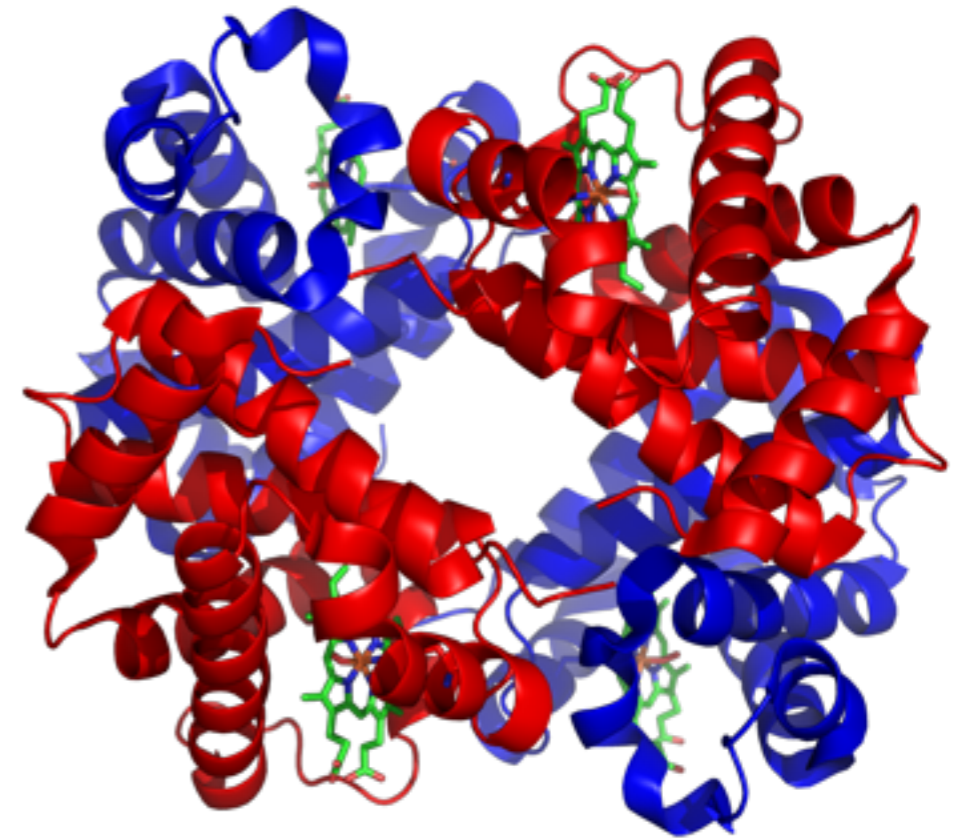
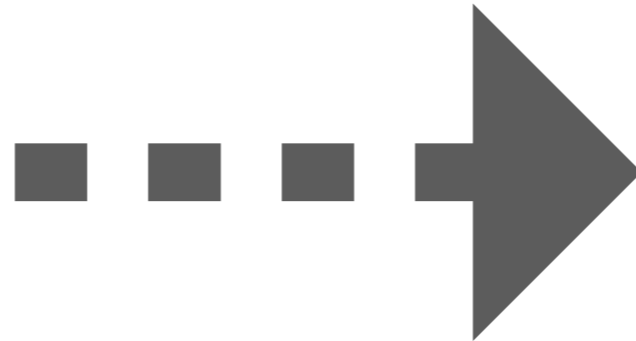
raftlib.io



WHERE I STARTED



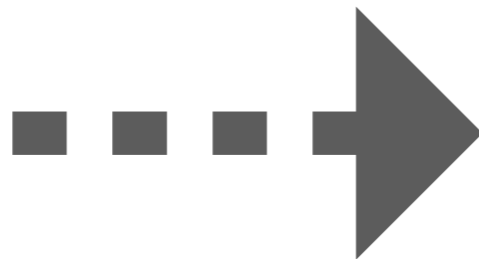
WHERE I STARTED



WHERE I STARTED



WHERE I STARTED



```
const uint8_t fsm_table[STATES][CHARS]=  
{  
    /* 0 - @      */ { [9]=0x21 },  
    /* 1 - NAME   */ { [0 ... 8]=0x11 , [11]=0x32 },  
    /* 2 - SEQ    */ { [1 ... 2]=0x42 , [5]=0x42 , [8]=0x42 , [11]=0x13 },  
    /* 3 - \n     */ { [1 ... 2]=0x42 , [5]=0x42 , [8]=0x42 , [10]=0x14 },  
    /* 4 - NAME2  */ { [0 ... 8]=0x54 , [11]=0x15 },  
    /* 5 - SCORE  */ { [0 ... 10]=0x65 , [11]=0x16 },  
    /* 6 - SEND/SO */ { [9] = 0x71 , [11] = 0x16 }  
};
```

NECESSITY DRIVES IDEAS

Java
Storm Obj-C R
Swift OpenMP Go
Scala
X10 Fortran Perl
C++ Javascript Ruby
Python Chapel Rust
C Pascal Spark MPI
Ada



NECESSITY DRIVES IDEAS

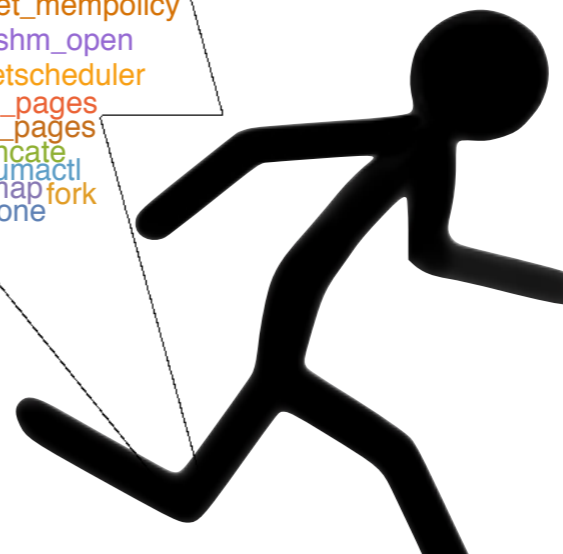
Java
Storm Obj-C R
Swift OpenMP Go
Scala
X10 Fortran Perl
C++ Javascript Ruby
Python Chapel Rust
C Pascal Spark MPI
Ada

sched_setscheduler
pthread_setaffinity
pthread_getaffinity
posix_memadvise

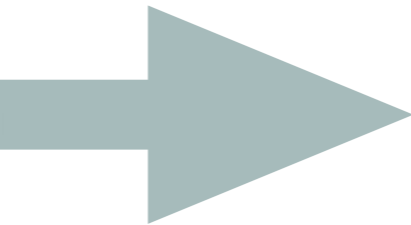
sched_setaffinity
sched_getaffinity

unlink
pthread_create
set_mempolicy
shm_open

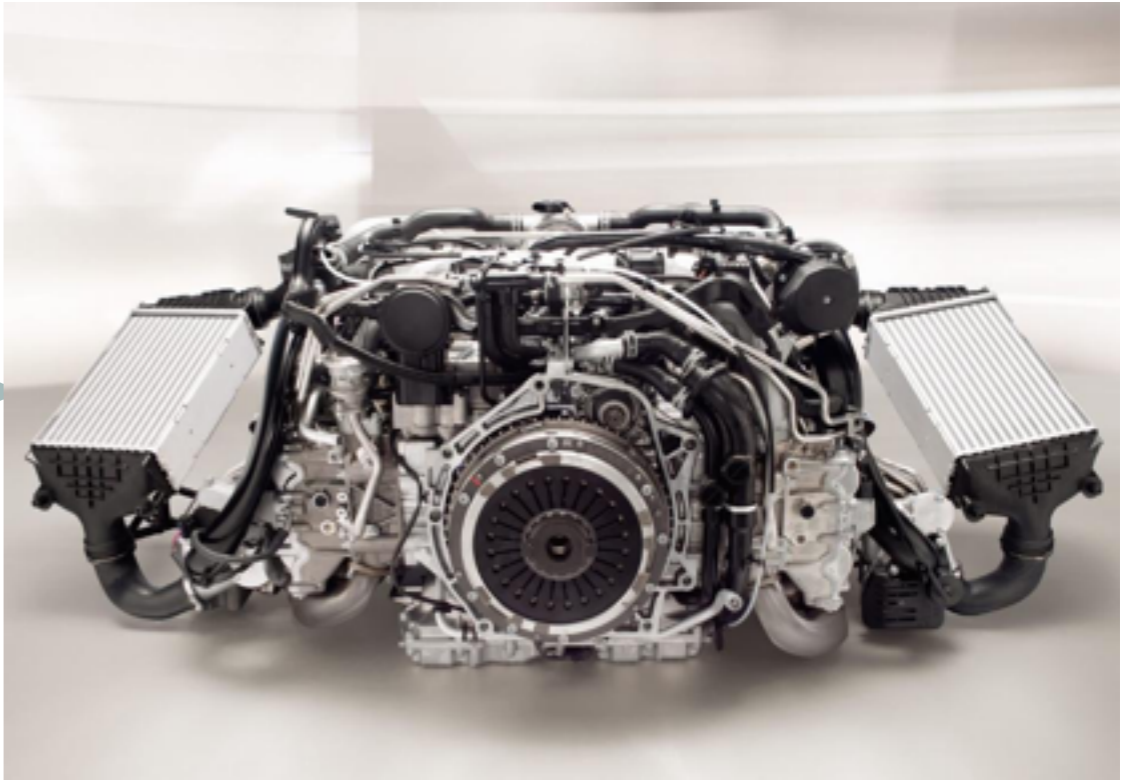
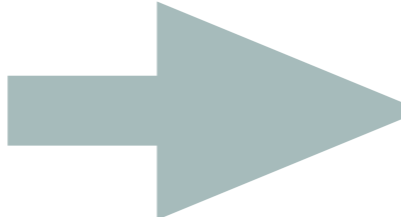
sched_getscheduler
migrate_pages
move_pages
truncate
numactl
mmap fork
clone



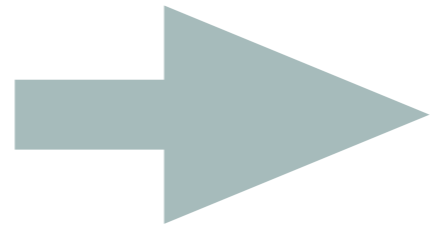
AN (PERHAPS BAD) ANALOGY



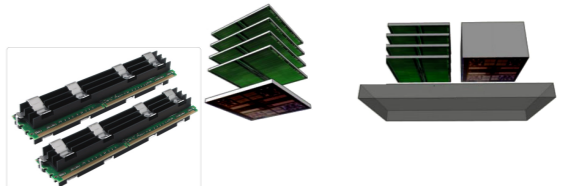
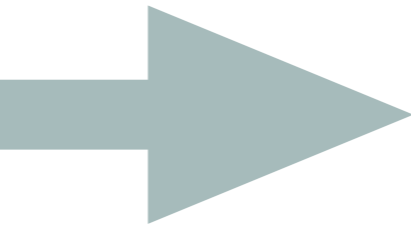
AN (PERHAPS BAD) ANALOGY



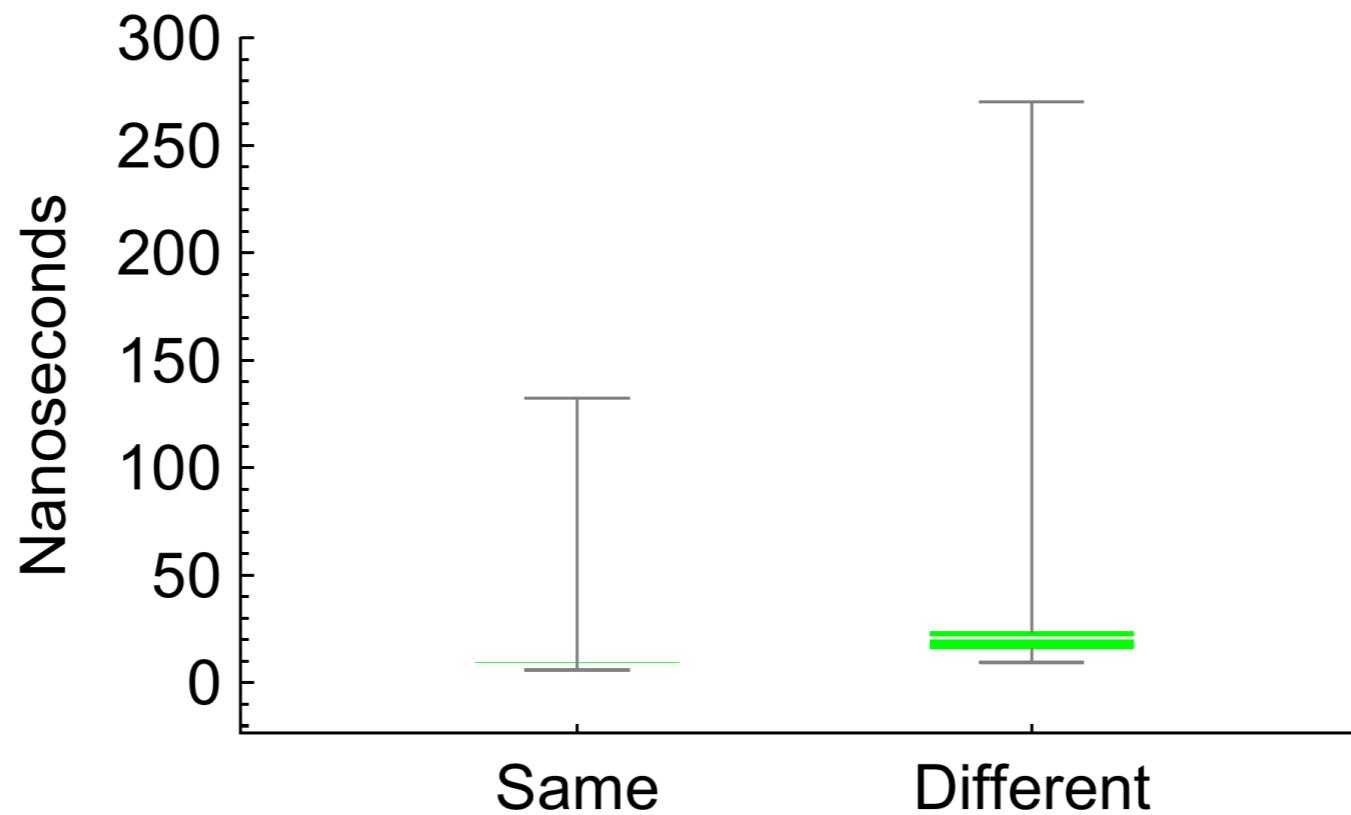
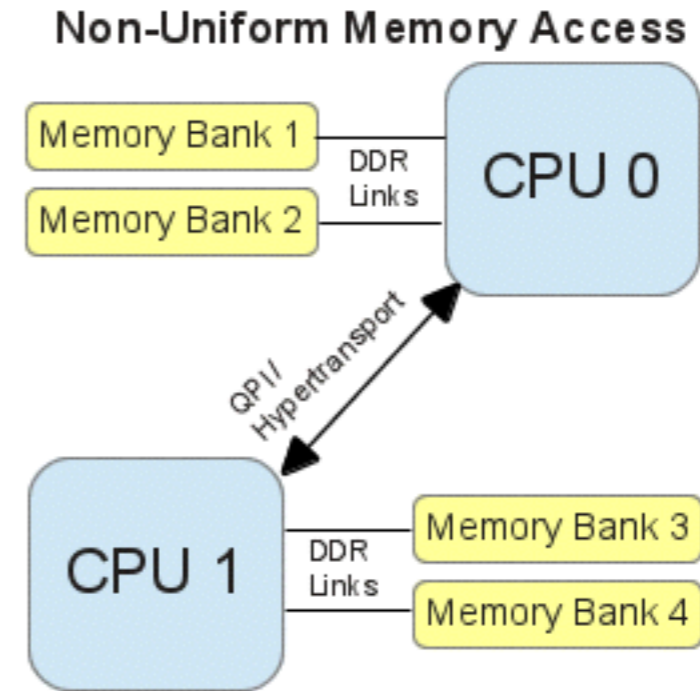
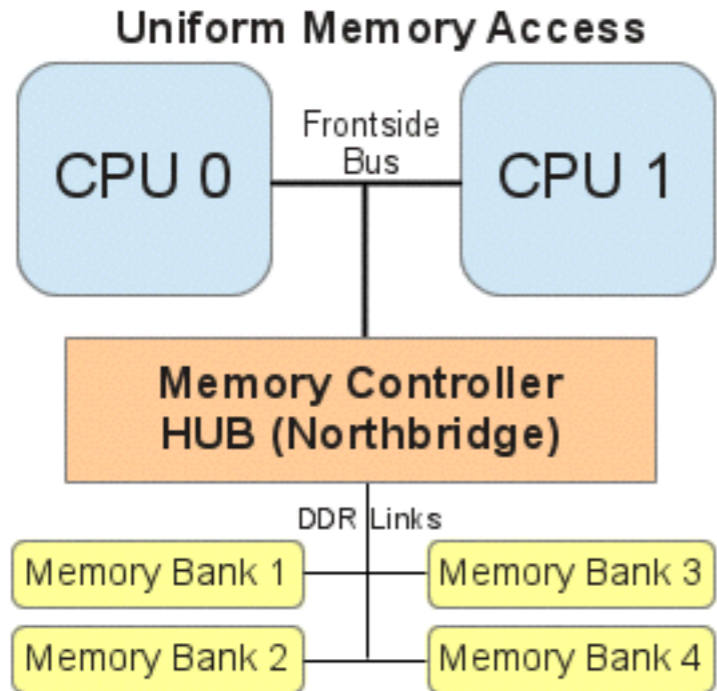
AN (PERHAPS BAD) ANALOGY



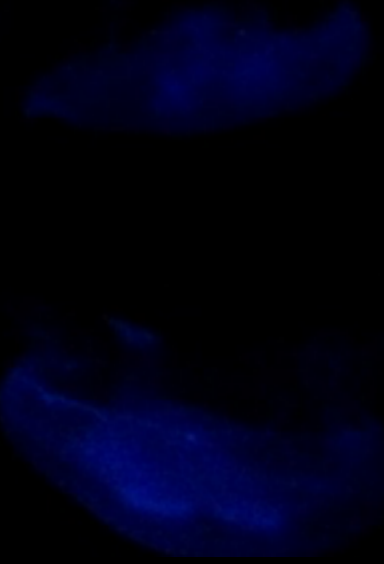
ANALOGY PART TWO



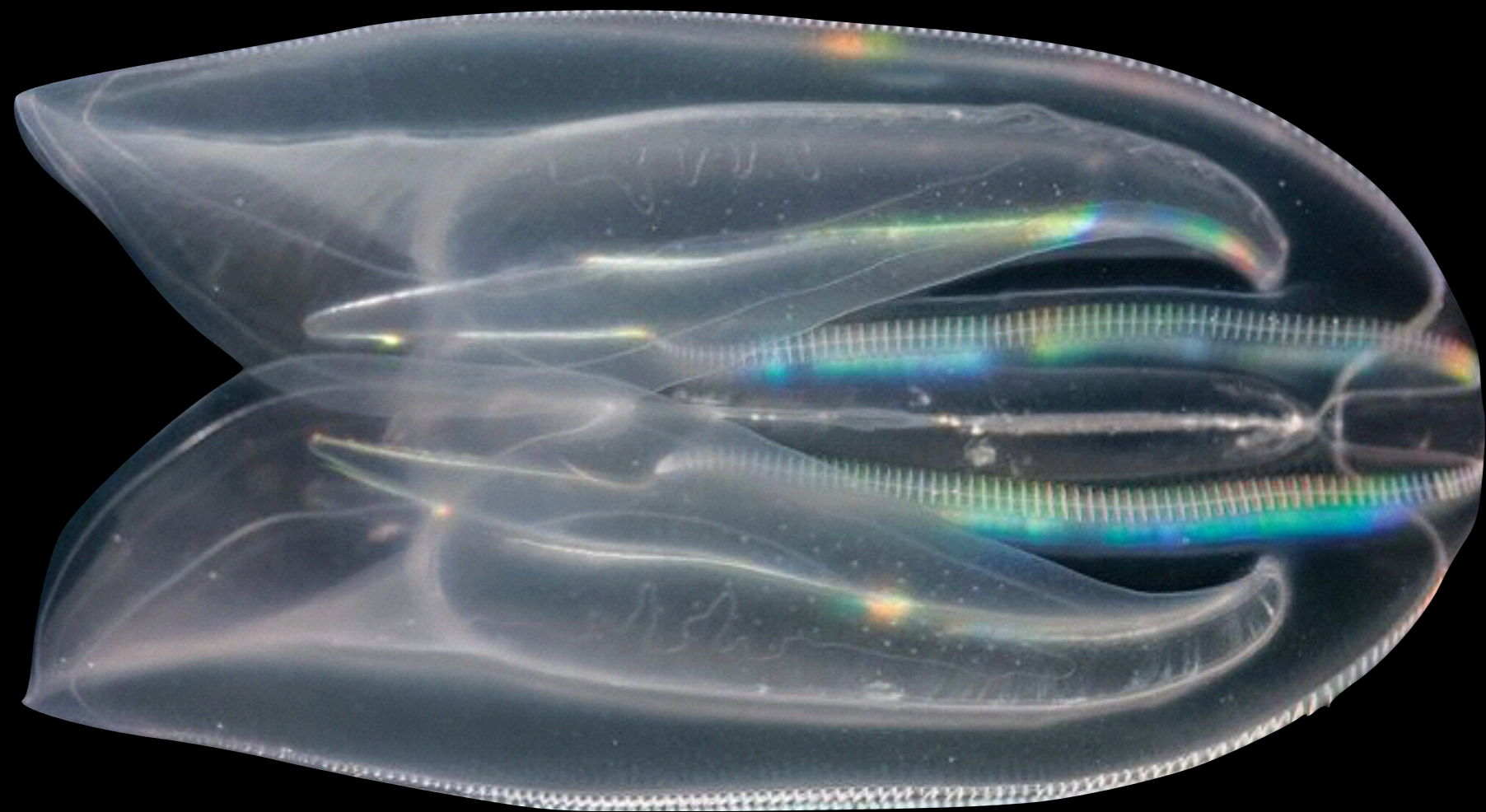
TOPOLOGY



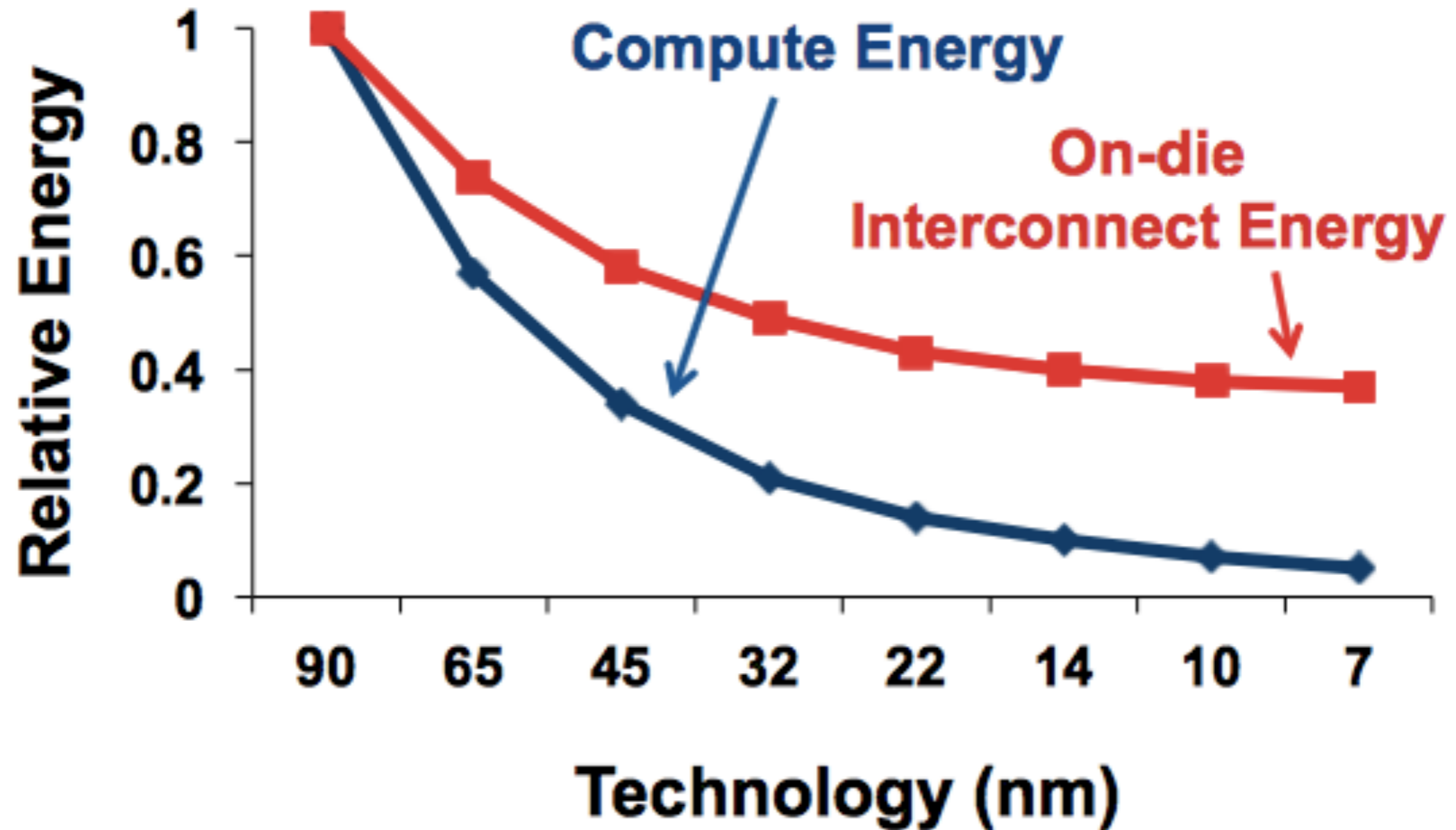
THE JELLYFISH



THE FIRST ORGANISM TO OVERLAP ACCESS AND EXECUTION



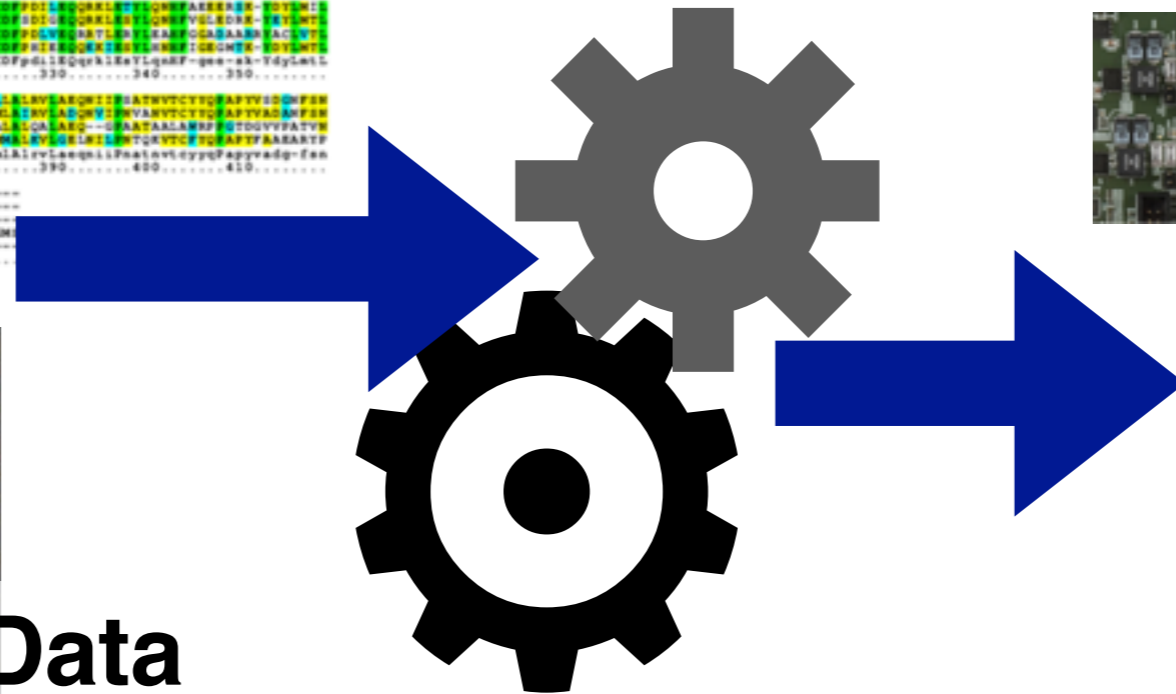
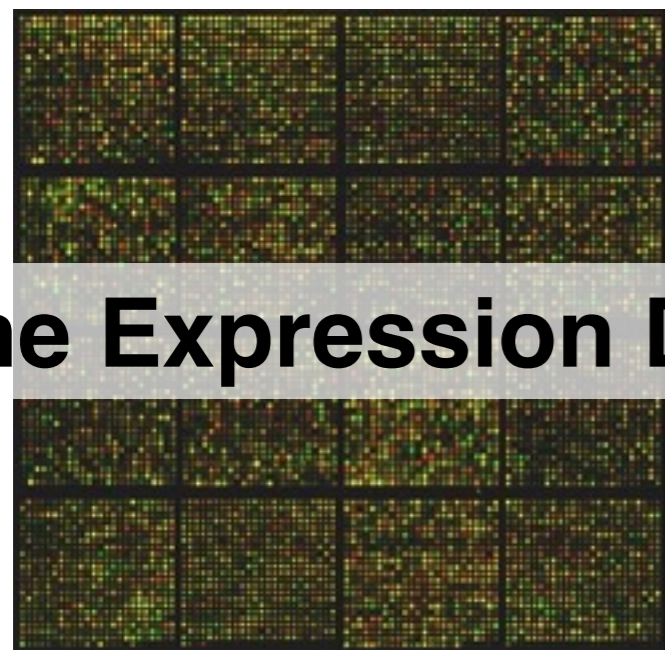
DATA MOVEMENT DOMINATES



Source: Shekhar Borkar, Journal of Lightwave Technology, 2013

I SHOULDN'T HAVE TO CARE

Multiple Sequence Alignment



Gene Expression Data

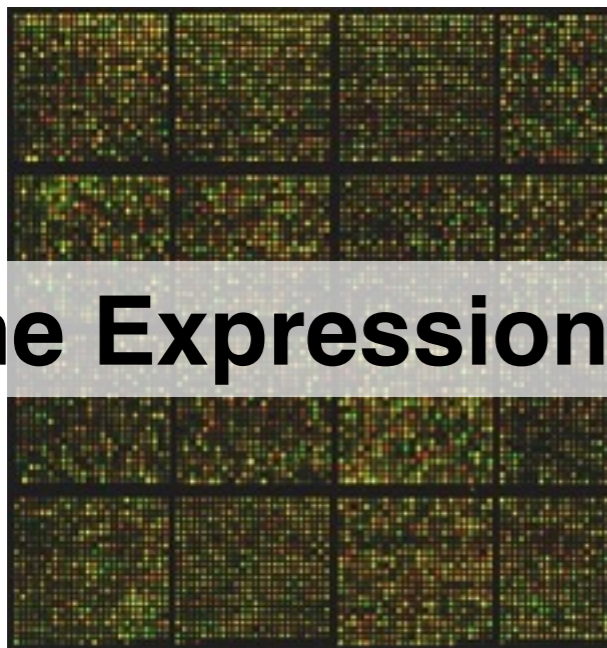
I SHOULDN'T HAV



Multiple Align



Gene Expression

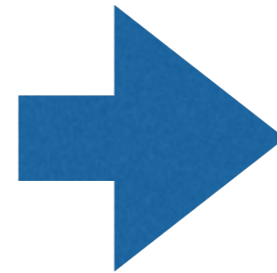


HARDWARE

1984

Cray X-MP/48

\$19 million / GFLOP

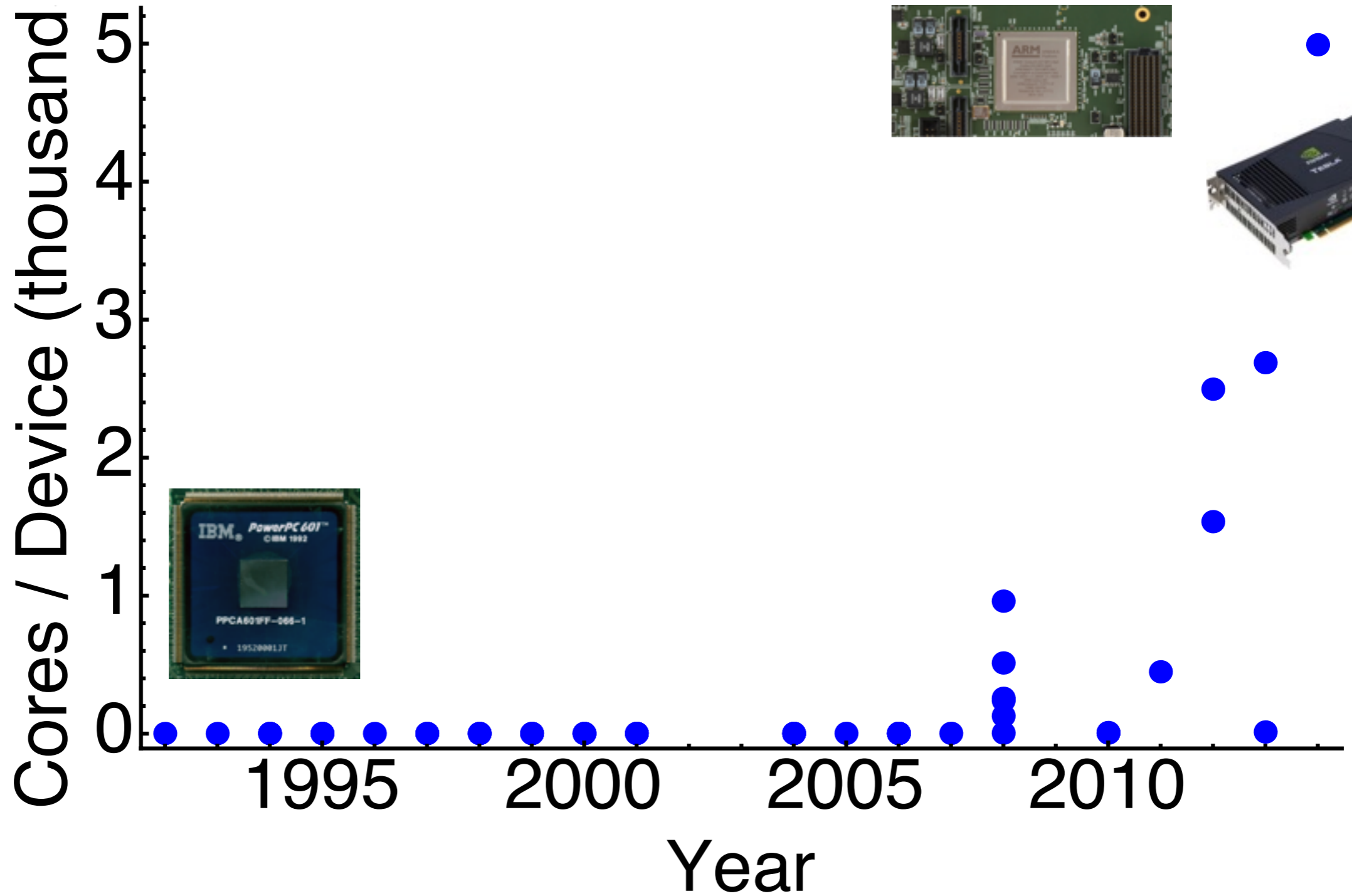


2015

\$.08 / GFLOP



CORES PER DEVICE



FINANCIAL INCENTIVE



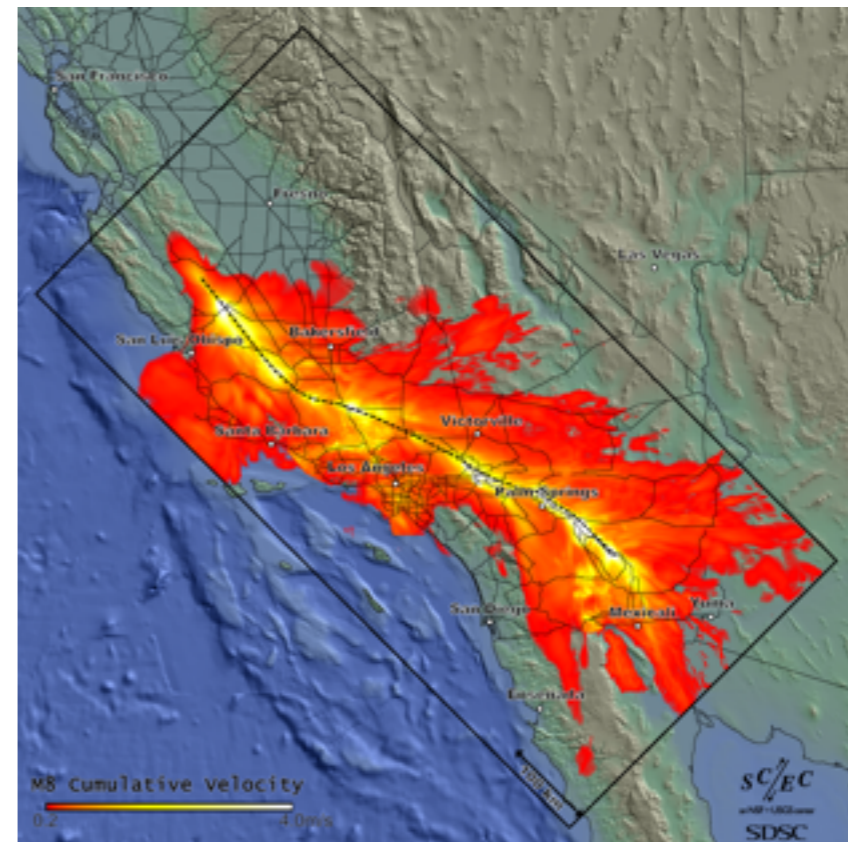
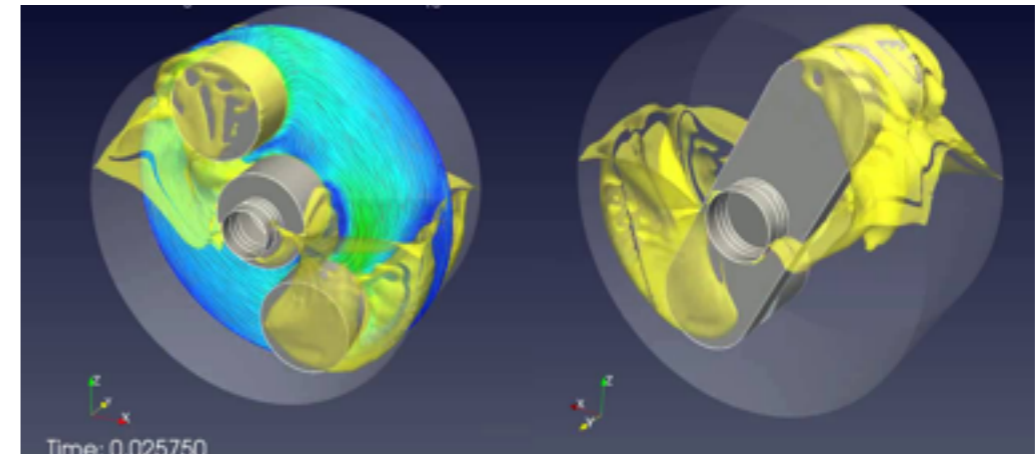
Sequential JS: \$4-7/line

Sequential Java: \$5-10/line

Embedded Code: \$30-50/line

HPC Code: \$100/line

WHY YOU SHOULD CARE



PRODUCTIVITY / EFFICIENCY AN EQUALIZER

- Titan SC estimates for porting range from 5, to > 20 million USD
- Most code never really optimized for machine topology, wasting \$\$ (time/product) and energy
- **Getting the most out of what you have is an equalizer**



Google Cloud Platform



Java
Storm Obj-C R
Go
Swift OpenMP Scala
X10 Fortran Perl
C++ Javascript Ruby
Python Chapel Rust
C Pascal Spark MPI
Ada

HAVES AND HAVE NOTS

Gov't / Big Business

- *Lots of \$\$*
- *Can hire the best people*
- *Can acquire the rest*
- *Plenty of compute resources*

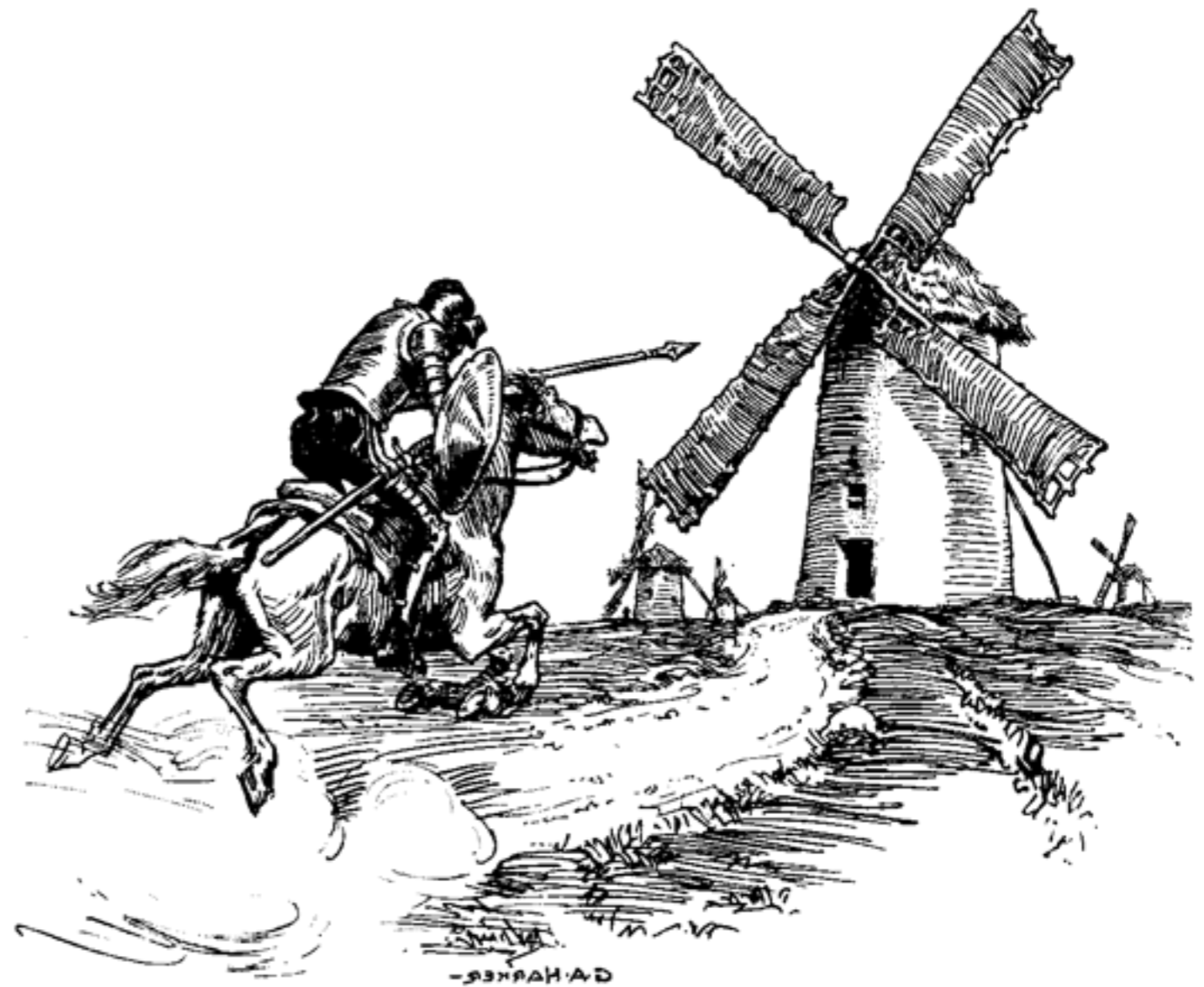
Start-up / Small Government

- *Not a lot of \$\$*
- *Often can't hire the best people*
- *Left to the mercy of cloud providers*



AN IDEA

Let's make
computers super
fast, and easy to
program



WHERE TO START

Brook for GPUs: Stream Computing on Graphics Hardware

Ian Buck Tim Foley Daniel Horn Jeremy Sugerman Kayvon Fatahalian Mike Houston Pat Hanrahan
Stanford University

Abstract

In this paper, we present Brook for GPUs, a system for general-purpose computation on programmable graphics hardware. Brook extends C to include simple data-parallel constructs, enabling the use of the GPU as a streaming coprocessor. We present a compiler and runtime system that abstracts and virtualizes many aspects of graphics hardware. In addition, we present an analysis of the effectiveness of the GPU as a compute engine compared to the CPU, to determine when the GPU can outperform the CPU for a particular algorithm. We evaluate our system with five applications, the SAXPY and SGEMV BLAS operators, image segmentation, FFT, and ray tracing. For these applications, we demonstrate that our Brook implementations perform comparably to hand-written GPU code and up to seven times faster than their CPU counterparts.

CR Categories: I.3.1 [Computer Graphics]: Hardware Architecture—Graphics processors D.3.2 [Programming Languages]: Language Classifications—Parallel Languages

Keywords: Programmable Graphics Hardware, Data Parallel Computing, Stream Computing, GPU Computing, Brook

1 Introduction

In recent years, commodity graphics hardware has rapidly evolved from being a fixed-function pipeline into having programmable vertex and fragment processors. While this new programmability was introduced for real-time shading, it has been observed that these processors feature instruction sets general enough to perform computation beyond the domain of rendering. Applications such as linear algebra [Krieger and Westermann 2003], physical simulation, [Harris et al. 2003], and a complete ray tracer [Purcell et al. 2002; Carr et al. 2002] have been demonstrated to run on GPUs.

Originally, GPUs could only be programmed using assembly languages. Microsoft's HLSL, NVIDIA's Cg, and OpenGL's GLSL allow shaders to be written in a higher level, C-like programming language [Microsoft 2003; IBM Research Center, Hawthorne, Research Center, Hawthorne, NY, 10532, USA; ggedik@us.ibm.com; bucu@us.ibm.com; thomas.j.watson@ibm.com; thomas.j.watson@ibm.com; thomas.j.watson@ibm.com; thomas.j.watson@ibm.com]. However, these languages do not assist the programmer in controlling other aspects of the graphics pipeline, such as allocating texture memory, loading shader programs, or constructing graphics primitives. As a result, the implementation of applications requires extensive knowledge of the latest graphics APIs as well as an understanding of the features and limitations of

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permission may be requested from Publications Dept., ACM, Inc., 1515 Broadway, New York, NY 10036 USA, fax +1 (212) 860-0481, or permission@acm.org.
© 2004 ACM 0730-0301/04/0800-1123 \$5.00

ABSTRACT

In this paper, we present SPADE – the System S declarative stream processing engine. System S is a large-scale, distributed data stream processing middleware under development at IBM T. J. Watson Research Center. As a front-end for rapid application development for System S, SPADE provides (1) an intermediate language for flexible composition of parallel and distributed data-flow graphs, (2) a toolkit of type-generic, built-in stream processing operators, that support scalar as well as vectorized processing and can seamlessly inter-operate with user-defined operators, and (3) a rich set of stream adapters to ingest/publish data from/to outside sources. More importantly, SPADE automatically brings performance optimization and scalability to System S applications. To that end, SPADE employs a code generation framework to create highly-optimized applications that run natively on the Stream Processing Core (SPC), the execution and communication substrate of System S, and take full advantage of other System S services. SPADE allows developers to construct their applications with fine granular stream operators without worrying about the performance implications that might exist, even in a distributed system. SPADE's optimizing compiler automatically maps applications into appropriately sized execution units in order to minimize communication overhead, while at the same time exploiting available parallelism. By virtue of the scalability of the System S runtime and SPADE's effective code generation and optimization, we can scale applications to a large number of nodes. Currently, we can run SPADE jobs on ≈ 500 processors within more than 100 physical nodes in a tightly connected cluster environment. SPADE has been in use at IBM Research to create real-world streaming appli-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
SIGMOD'08, June 9–12, 2008, Vancouver, BC, Canada.
Copyright 2008 ACM 978-1-60558-102-6/08/06 ...\$5.00.

modern hardware. In addition, the user is forced to express their algorithm in terms of graphics primitives, such as textures and triangles. As a result, general-purpose GPU computing is limited to only the most advanced graphics developers.

This paper presents Brook, a programming environment that provides developers with a view of the GPU as a streaming coprocessor. The main contributions of this paper are:

- The presentation of the Brook stream programming model for general-purpose GPU computing. Through the use of streams, kernels and reduction operators, Brook abstracts the GPU as a streaming processor.
- The demonstration of how various GPU hardware limitations can be virtualized or extended using our compiler and runtime system; specifically, the GPU memory system, the number of supported textures, and support for user-defined data structures.
- The presentation of a cost model for comparing GPU vs. CPU performance tradeoffs to better understand under what circumstances the GPU outperforms the CPU.

2 Background

2.1 Evolution of Streaming Hardware

Programmable graphics hardware dates back to the original programmable frame-buffer architectures [Foley et al. 1980]. One of the first examples of programmable graphics hardware was the Parallax mode (Molnar et al. 1989). These systems embedded graphics processing into the main processor, on the same bus as the main processor. Each rendering pass implements a SIMD instruction that performs a basic arithmetic operation and updates the frame buffer. The use of SIMD instructions in graphics hardware was first demonstrated in the OpenCL architecture [Woo et al. 1999] can be abstracted as a SIMD processor. Each rendering pass implements a SIMD instruction that performs a basic arithmetic operation and updates the frame buffer.

Using this abstract model, we present SPADE, a declarative stream processing engine. SPADE is a declarative stream processing engine that allows developers to express their algorithms in terms of streams and kernels. The stream processing engine uses operators to execute a single instruction, and a write to off-chip memory [Russell 1978; Krieger and Westermann 2003]. This is a major departure from traditional graphics hardware which executes small, fixed-size instructions and stores data to a memory. SPADE's stream processing engine is designed to be used in a distributed environment. The stream processing engine uses operators to execute a single instruction, and a write to off-chip memory [Russell 1978; Krieger and Westermann 2003]. This is a major departure from traditional graphics hardware which executes small, fixed-size instructions and stores data to a memory. SPADE's stream processing engine is designed to be used in a distributed environment.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permission may be requested from Publications Dept., ACM, Inc., 1515 Broadway, New York, NY 10036 USA, fax +1 (212) 860-0481, or permission@acm.org.
© 2004 ACM 0730-0301/04/0800-1123 \$5.00

WASHINGTON UNIVERSITY

SCHOOL OF ENGINEERING AND APPLIED SCIENCE

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

X-SIM AND X-EVAL:

TOOLS FOR SIMULATION AND ANALYSIS OF HETEROGENEOUS PIPELINED

ARCHITECTURES

by

Saurabh Gayen

Prepared under the direction of Professor Mark A. Franklin

SPADE: The System S Declarative Stream Processing Engine

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

Submitted to the School of Engineering and Applied Science

Washington University in partial fulfillment of the

requirements for the degree of

Master of Science in Computer Science

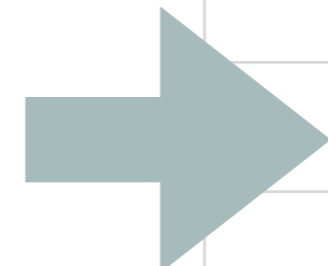
WHERE TO START

~~First Language~~

- Object oriented
- Multiple inheritance
- Simple types
- lots more



Library



Language
Java
C
C++
C#
Python

RAFTLIB

C++ Streaming Template Library

Simplifies parallelization of code

Abstracts the details

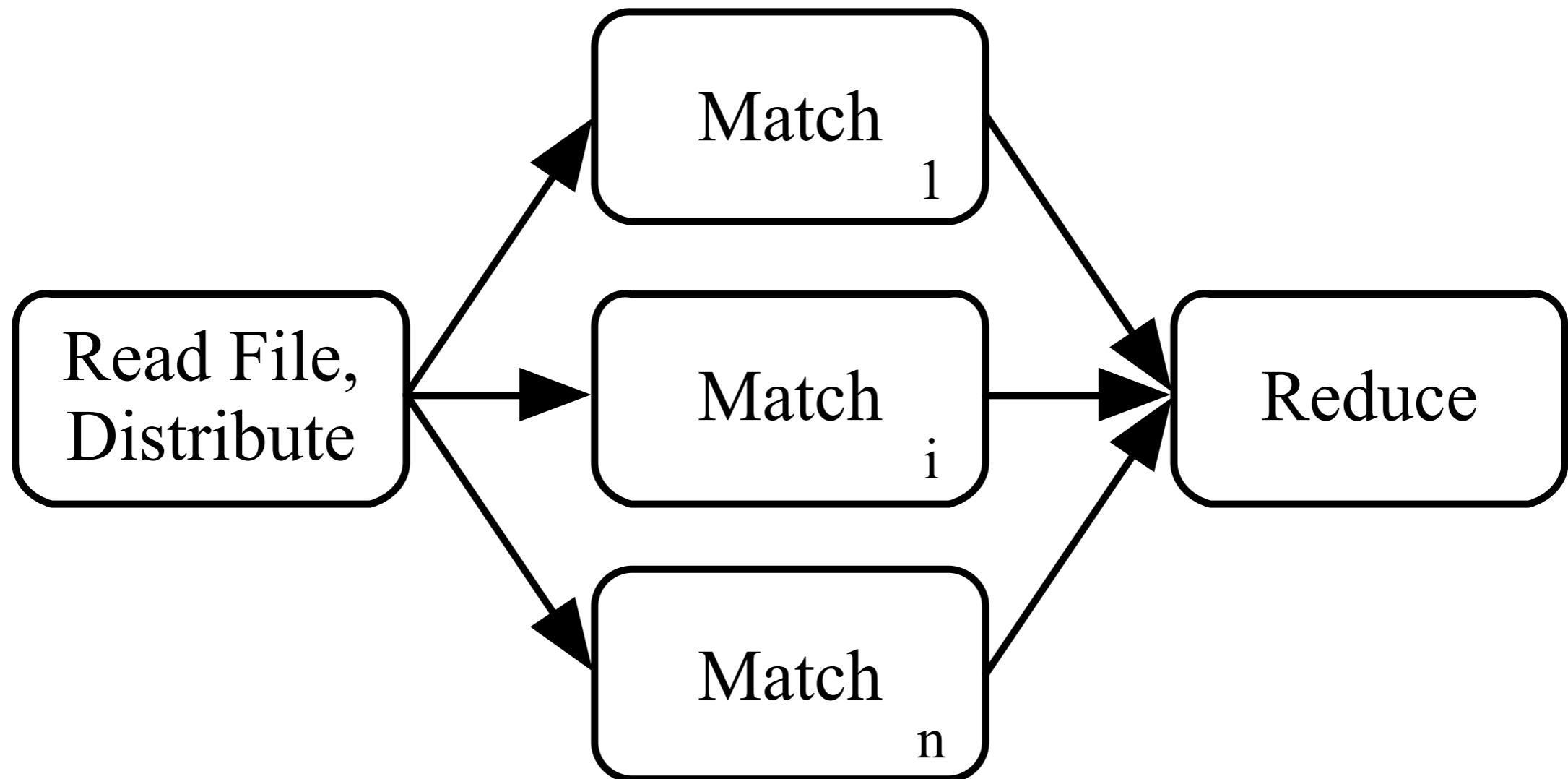
Auto-manages allocation types, data movement.

software download: <http://raftlib.io>

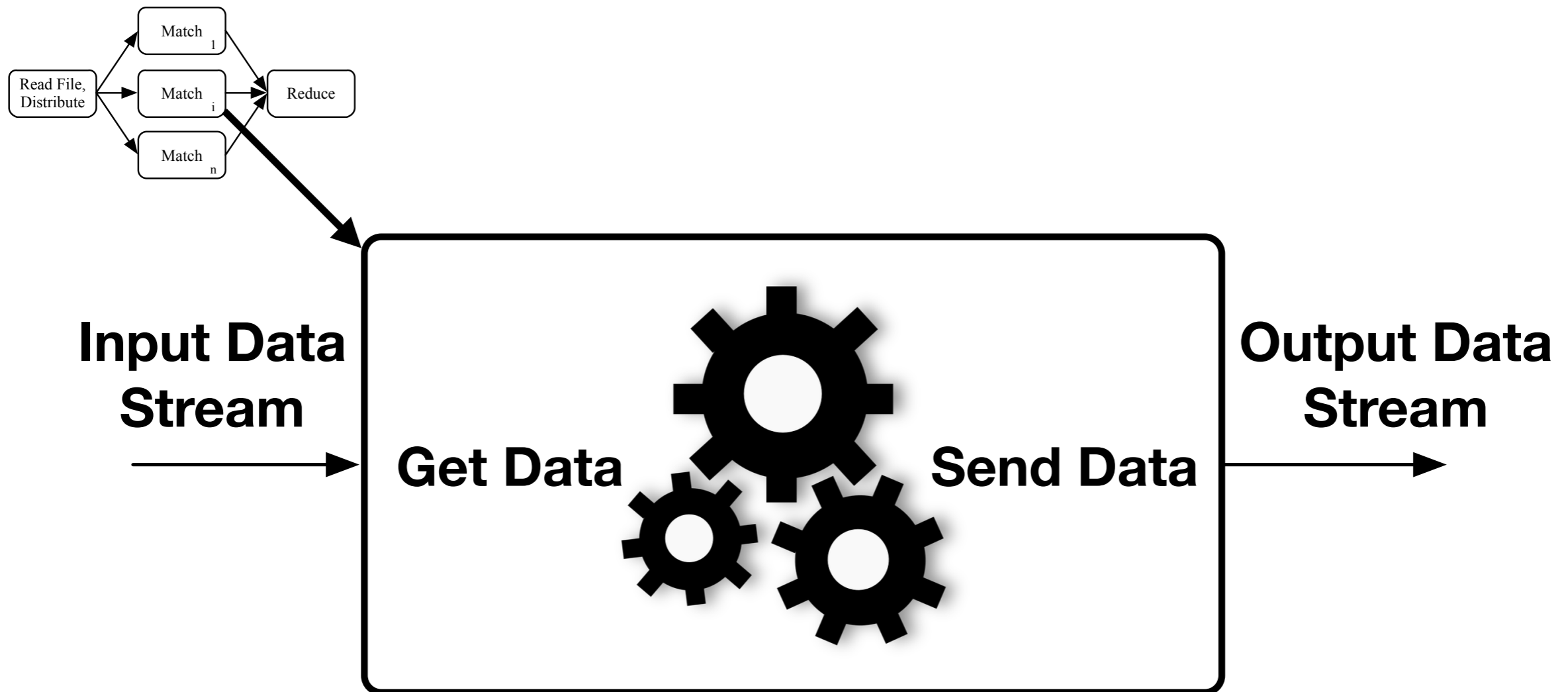
ISSUES THAT NEED SOLUTIONS

- Usable Generic API
- Where to run the code
- Where to allocate memory
- How big of allocations to start off with
- How many user vs. kernel threads
- ...Gilligan had a better chance of having a 3 hour tour
- ...than we have of fitting all the issues on a single slide

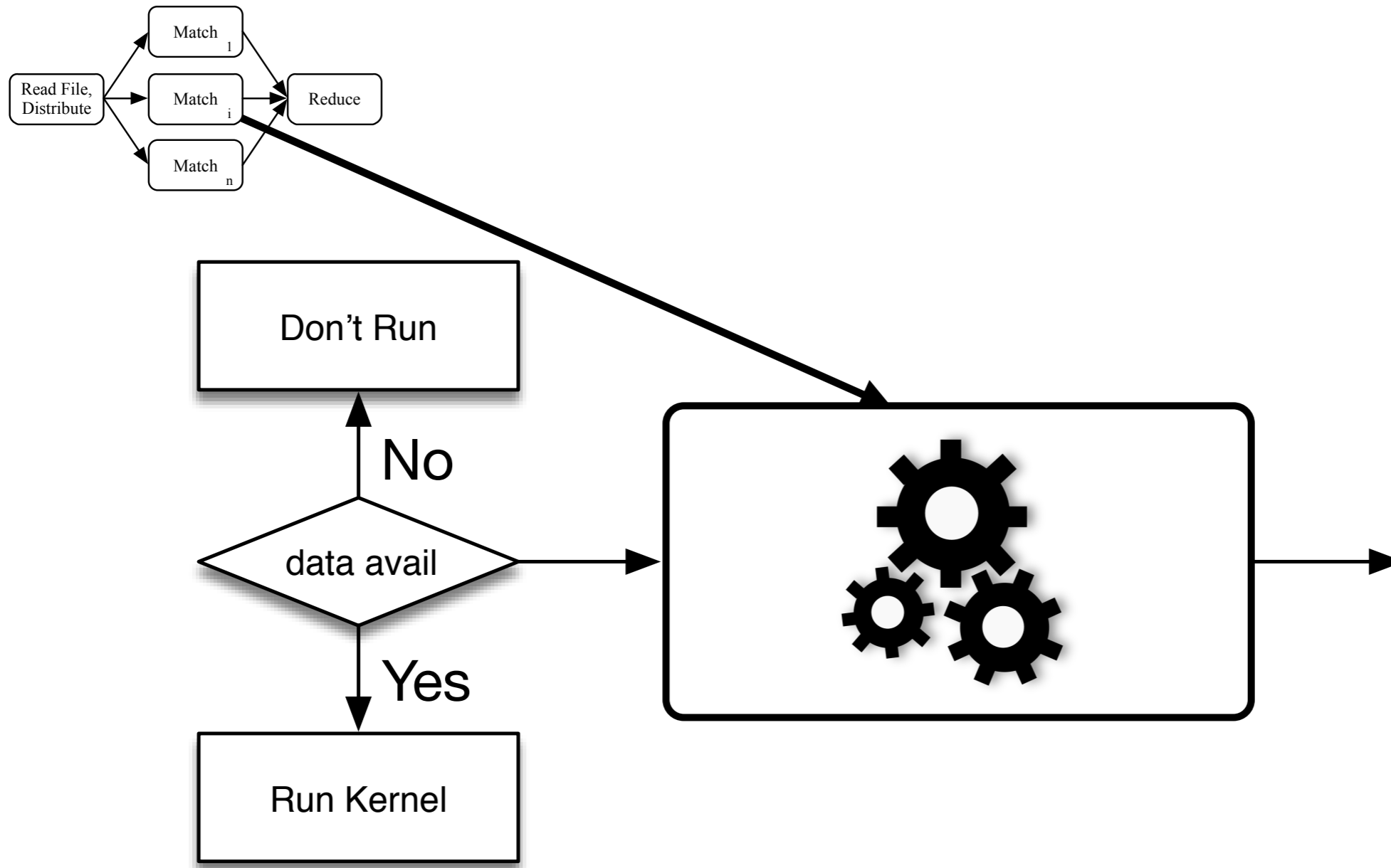
STREAM PROCESSING (STRING SEARCH EXAMPLE)



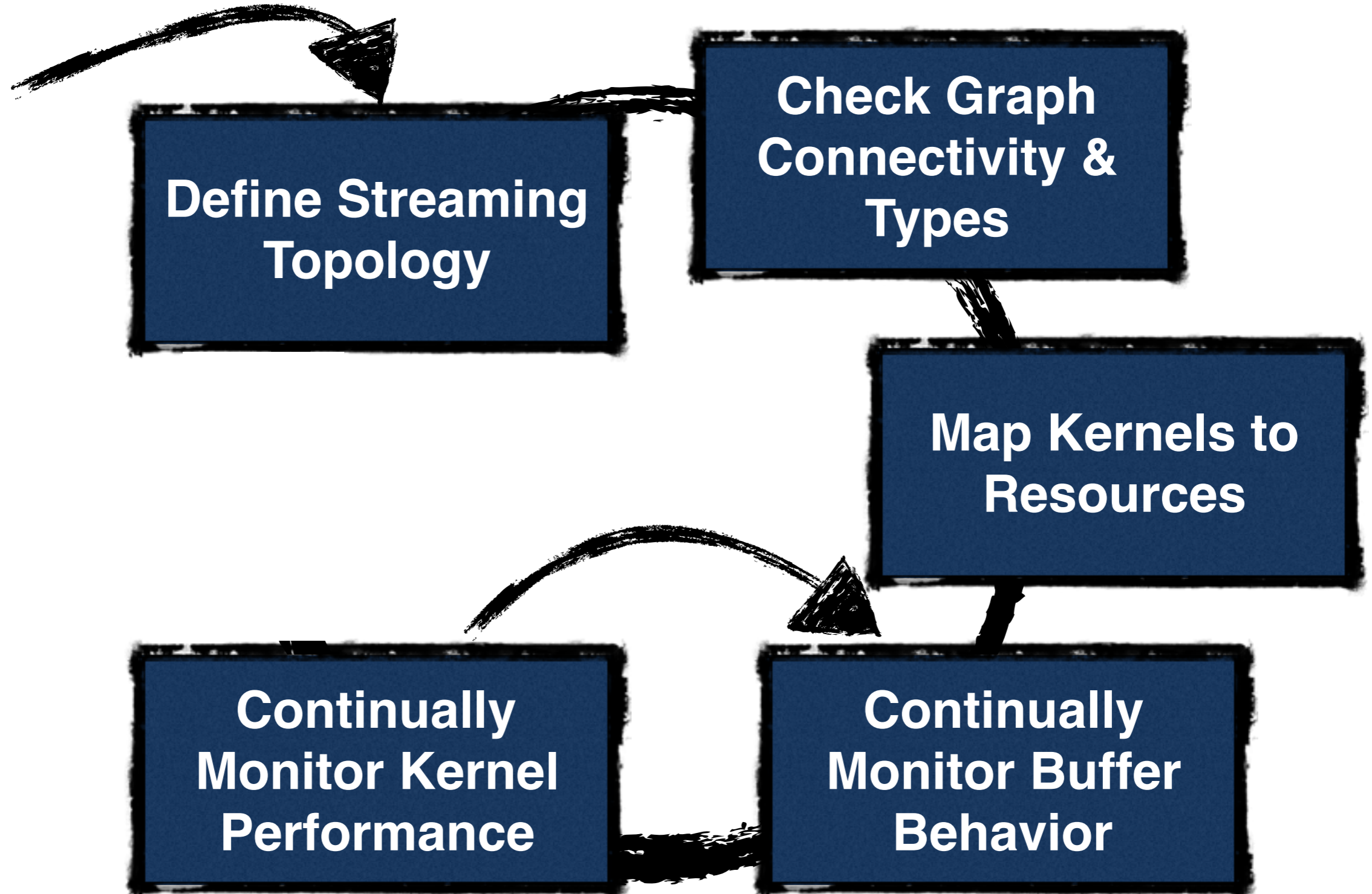
STREAM PROCESSING (STRING SEARCH EXAMPLE)



STREAM PROCESSING (STRING SEARCH EXAMPLE)



HOW IT WORKS (HIGH LEVEL)



Stream Processing (Search Ex)

```
template < class T > class search : public raft::kernel
{
public:
    search( const std::string && term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }

    search( const std::string &term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }

    virtual ~search() = default;

    virtual raft::kstatus run()
    {
        auto &chunk( input[ "0" ].template peek< T >() );
        auto it( chunk.begin() );
        do
        {
            it = std::search( it, chunk.end(),
                             term.begin(), term.end() );
            if( it != chunk.end() )
            {
                output[ "0" ].push( it.location() );
                it += 1;
            }
            else
            {
                break;
            }
        }
        while( true );
        input[ "0" ].unpeek();
        input[ "0" ].recycle( );
        return( raft::proceed );
    }
private:
    const std::size_t term_length;
    const std::string term;
};
```


Stream Processing (Search Ex)

```
template < class T > class search : public raft::kernel
{
public:
    search( const std::string && term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }

    search( const std::string &term ) : raft::kernel(),
                                       term_length( term.length() ),
                                       term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }

    virtual ~search() = default;

    virtual raft::kstatus run()
    {
        auto &chunk( input[ "0" ].template peek< T >() );
        auto it( chunk.begin() );
        do
        {
            it = std::search( it, chunk.end(),
                             term.begin(), term.end() );
            if( it != chunk.end() )
            {
                output[ "0" ].push( it.location() );
                it += 1;
            }
            else
            {
                break;
            }
        }
        while( true );
        input[ "0" ].unpeek();
        input[ "0" ].recycle( );
        return( raft::proceed );
    }
private:
    const std::size_t term_length;
    const std::string term;
};
```

Stream Processing (Search Ex)

```
template < class T > class search : public raft::kernel
{
public:
    search( const std::string && term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< T >( "0" );
    }
    search( const std::string && term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }
    virtual
    {
        search( const std::string &term ) : raft::kernel(),
                                             term_length( term.length() ),
                                             term( term )
        {
            input.addPort< T >( "0" );
            output.addPort< std::size_t >( "0" );
        }
        {
            break;
        }
    }
    while( true );
    input[ "0" ].unpeek();
    input[ "0" ].recycle( );
    return( raft::proceed );
}
private:
    const std::size_t term_length;
    const std::string term;
};
```

Stream Processing (Search Ex)

```
template < class T > class search : public raft::kernel
{
public:
    search( const std::string && term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }

    search( const std::string &term ) : raft::kernel(),
                                         term_length( term.length() ),
                                         term( term )
    {
        input.addPort< T >( "0" );
        output.addPort< std::size_t >( "0" );
    }

    virtual ~search() = default;
};
```

```
virtual raft::kstatus run()
{
    auto &chunk( input[ "0" ].template peek< T >() );
    auto it( chunk.begin() );
    do
    {
        it = std::search( it, chunk.end(),
                        term.begin(), term.end() );
        if( it != chunk.end() )
        {
            output[ "0" ].push( it.location() );
            it += 1;
        }
        else
        {
            break;
        }
    }
    while( true );
    input[ "0" ].unpeek();
    input[ "0" ].recycle();
    return( raft::proceed );
}
```

```
private:
    const std::size_t term_length;
    const std::string term;
};
```

Stream Processing (Search Ex)

```
template < class T > class search : public raft::kernel
{
public:
    search( const s virtual raft::kstatus run()
    {
        {
            input.addPor output.addPo
        }
        search( const s
        {
            input.addPor output.addPo
        }
        virtual ~search
        virtual raft::k
        {
            auto &chunk(
            auto it( chu
            do
            {
                it = std:
                if( it !=
                {
                    output
                    it +=
                }
                else
                {
                    break;
                }
            }
            while( true
            input[ "0" ]
            input[ "0" ]
            return( raft
        }
private:
    const std::size
    const std::string term;
};
```

Stream Processing (Search Ex)

```
int
main( int argc, char **argv )
{
    using chunk = raft::filechunk< 256 >;
    using fr     = raft::filereader< chunk, false >;
    using search = search< chunk >;
    using print  = raft::print< std::size_t, '\n' >;

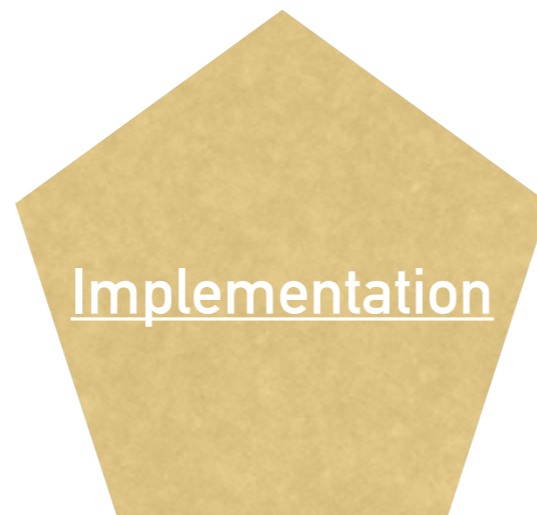
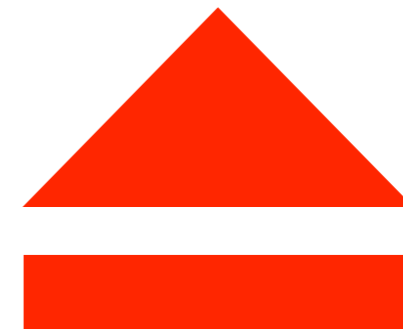
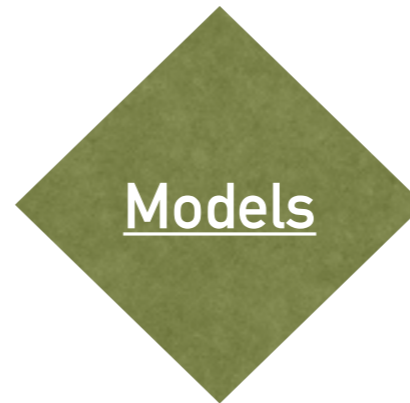
    const std::string term( "Alice" );
    raft::map m;

    fr read( argv[ 1 ], 1, term.length() );
    search find( term );
    print p;

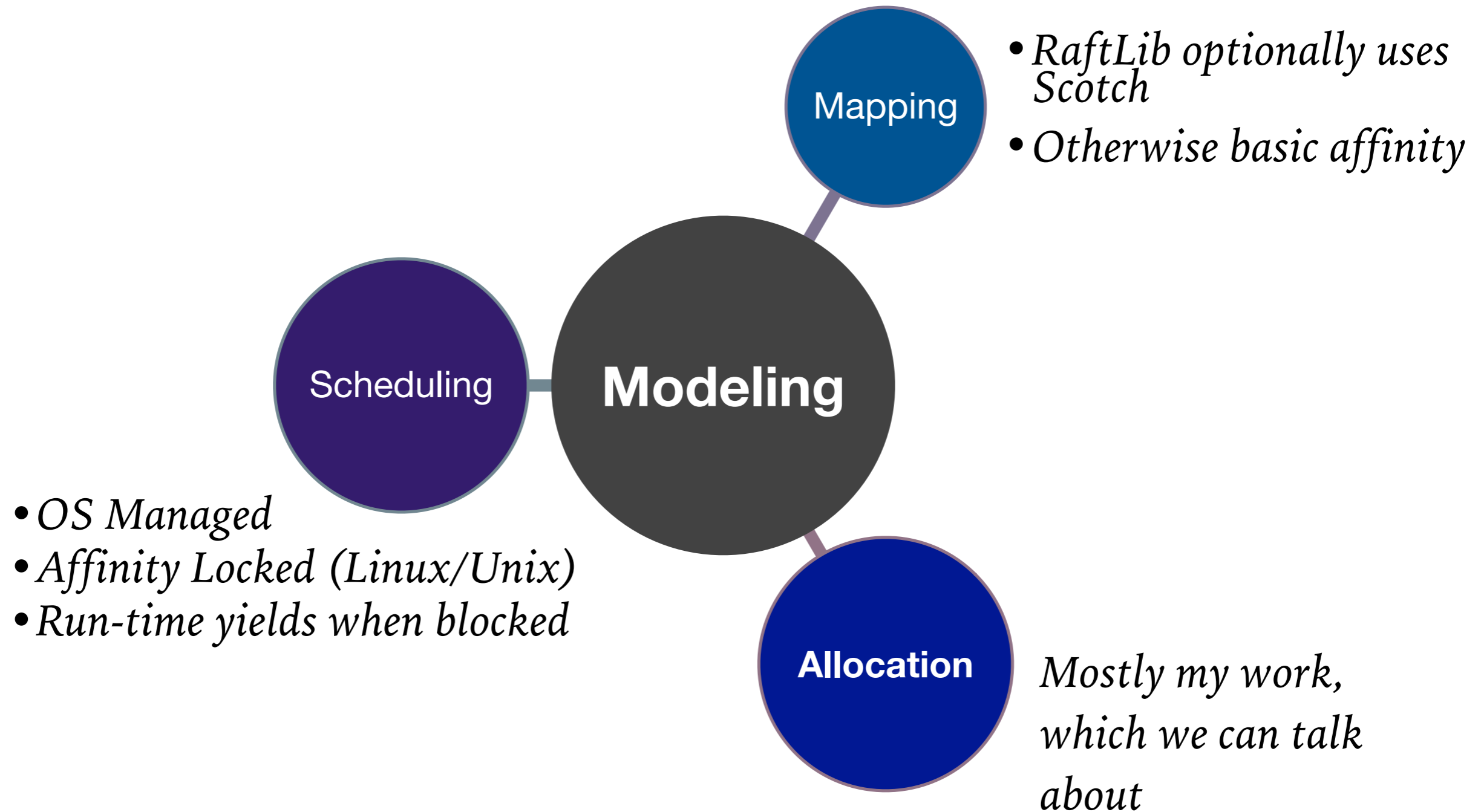
    m += read >> raft::order::out >> find >> raft::order::out >> p;
    m.exe();

    return( EXIT_SUCCESS );
}
```

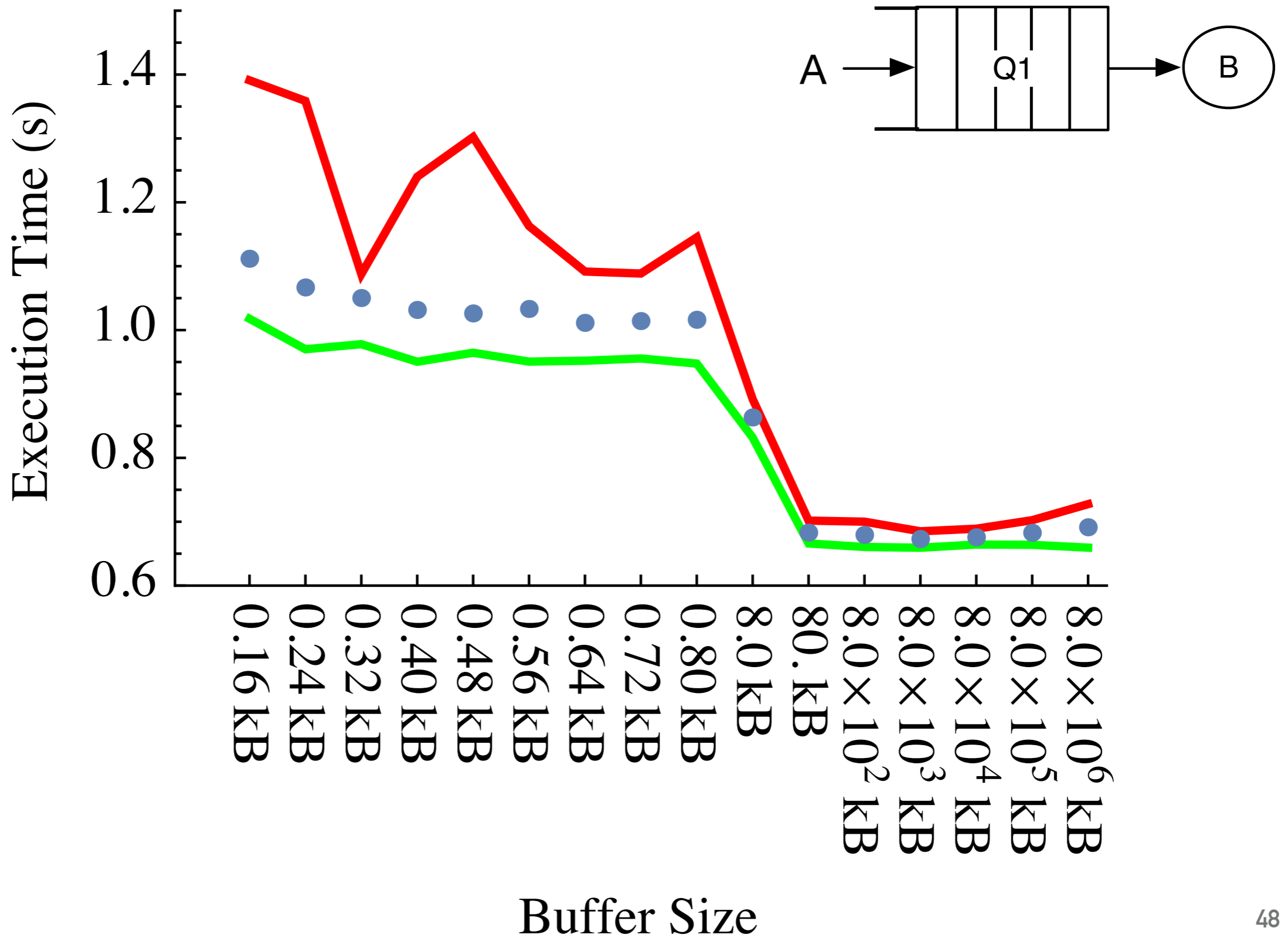
CHOOSE YOUR ADVENTURE



MODELING



MODELING ISSUES - SIZE OF Q1



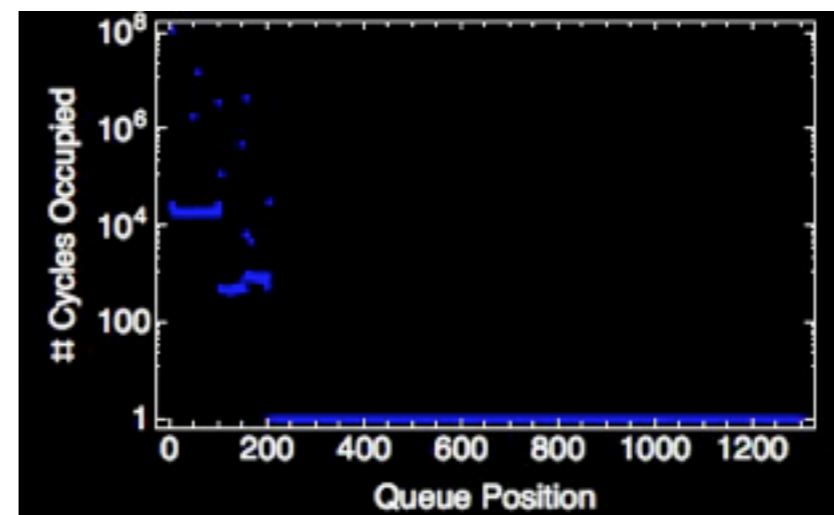
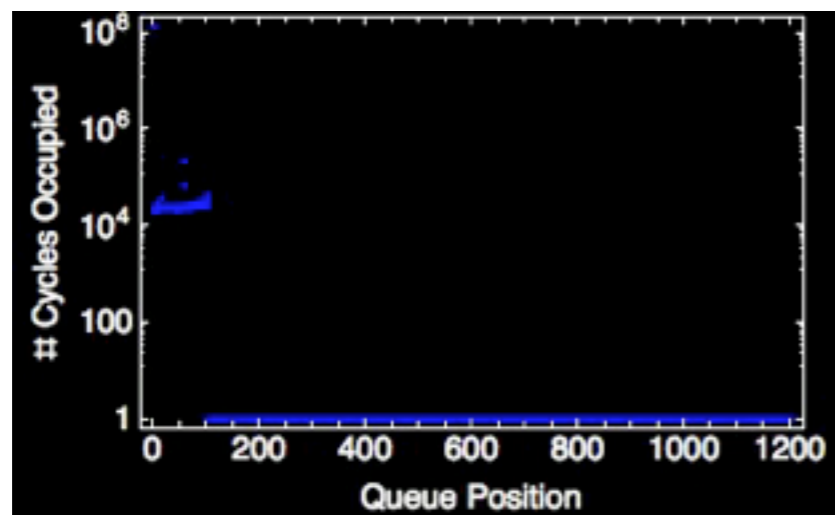
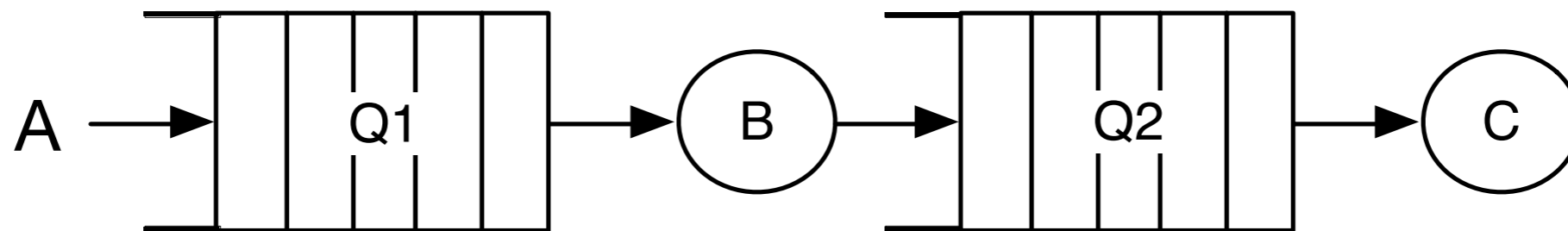
WHY DO WE NEED TO SOLVE

- Buffer allocations take time and energy
- Programmers are horrible at deciding (too many parameters)
- Hardware specific locations matter (NUMA)
- Re-allocating with an accelerator takes even more time (bus latency, hand-shakes, etc.)
- Must be solved in conjunction with partitioning/scheduling/
placement

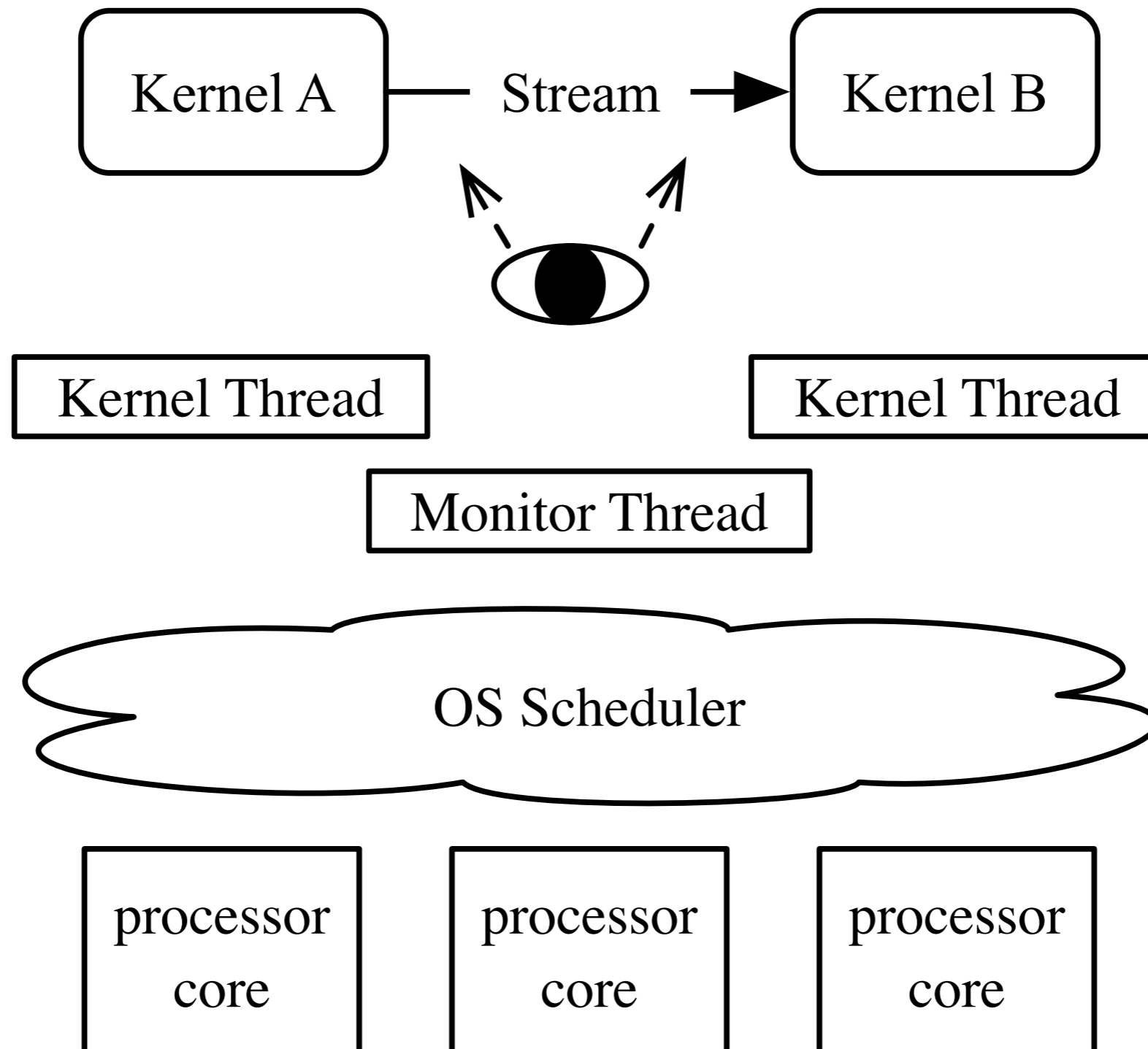
MODELING



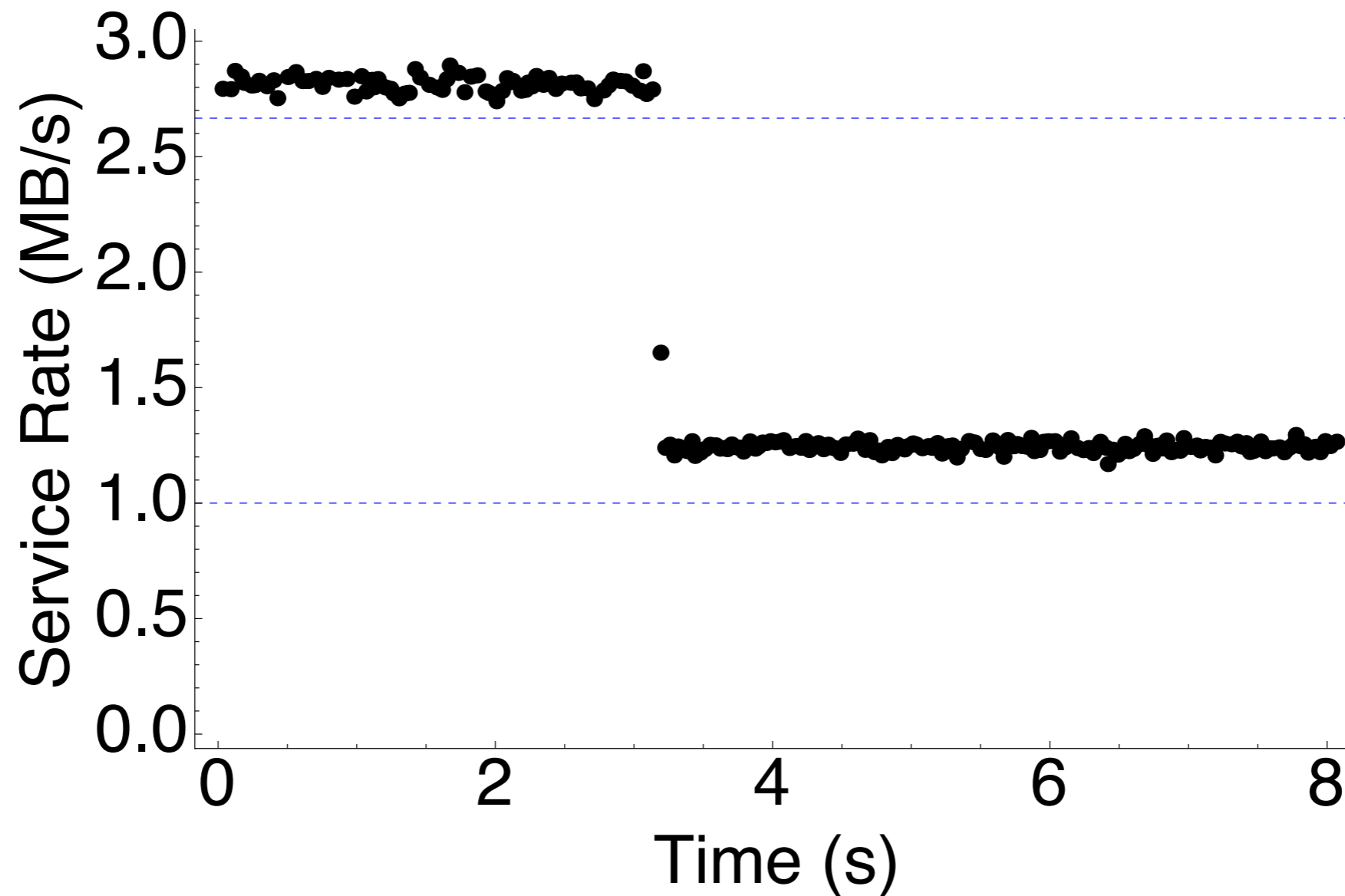
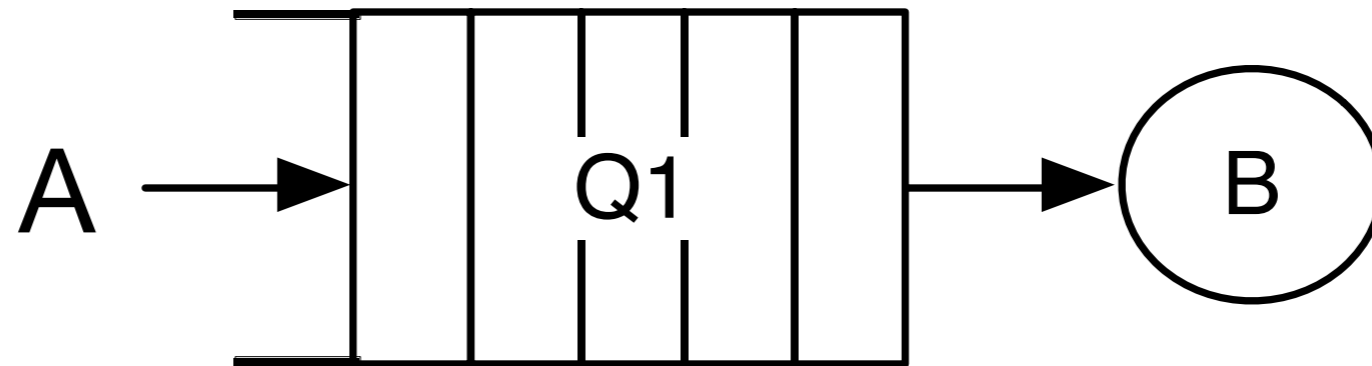
“Stream” is modeled as a Queue



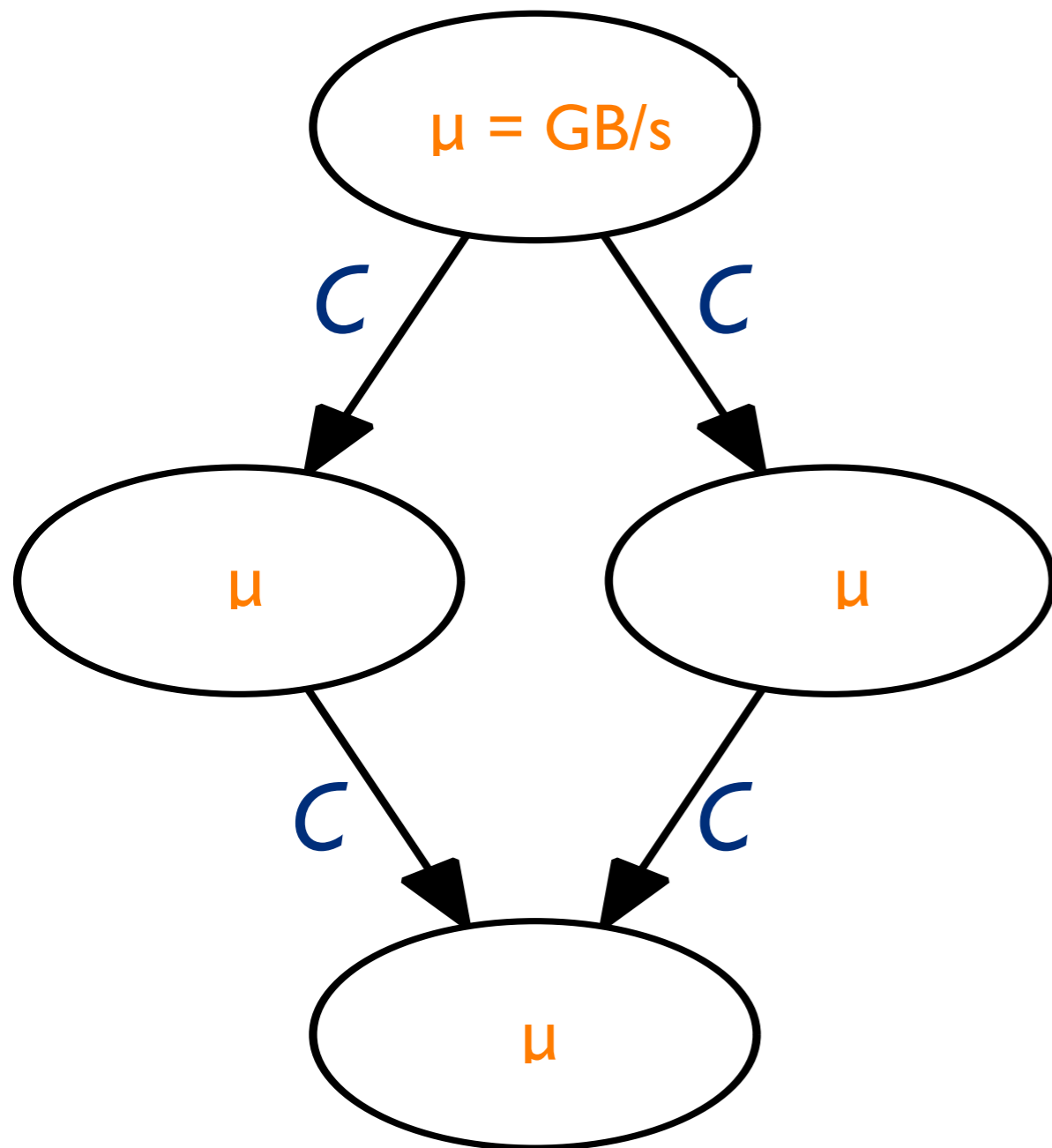
MONITOR



APPROXIMATE INSTRUMENTATION



SOLVE FOR THROUGHPUT QUICKLY



Use network flow model to quickly estimate flow within a streaming graph

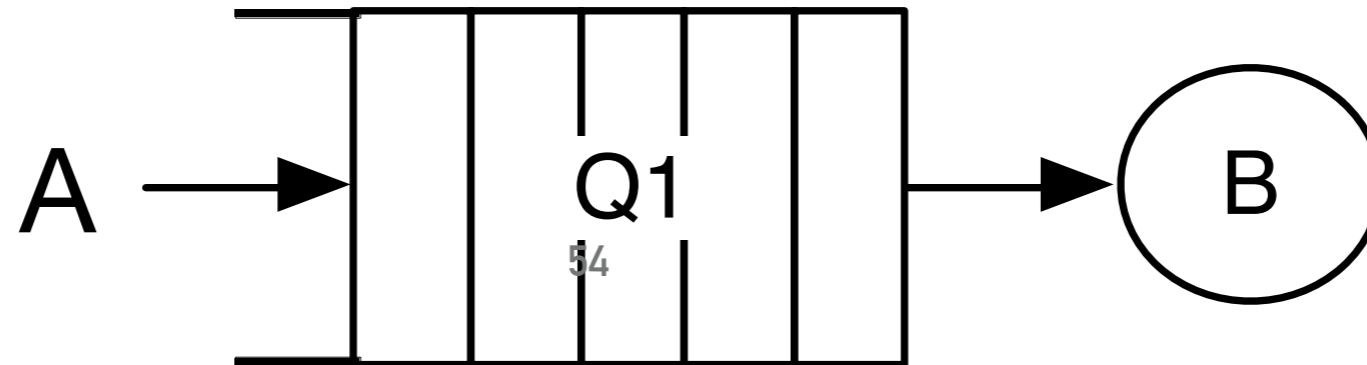
Decompose queueing network and solve each queueing station independently

HOW WOULD YOU GET THE RIGHT BUFFER SIZE?

Queueing Model to Solve for Buffer Size



?

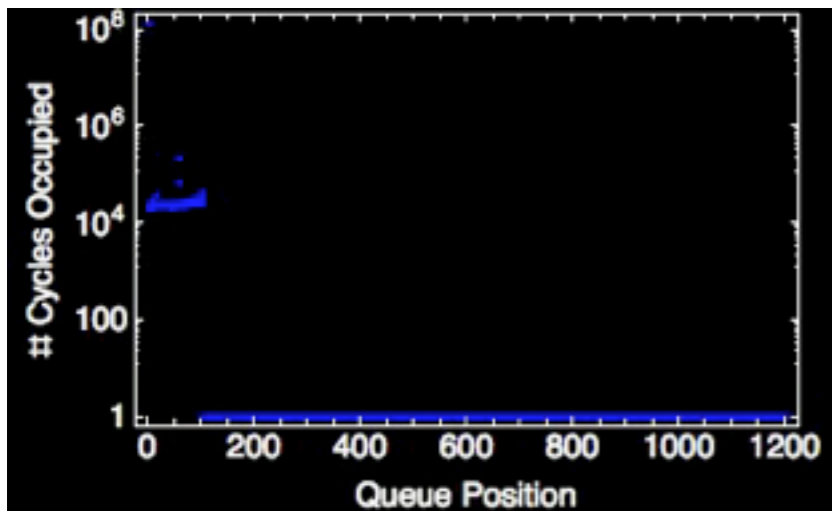
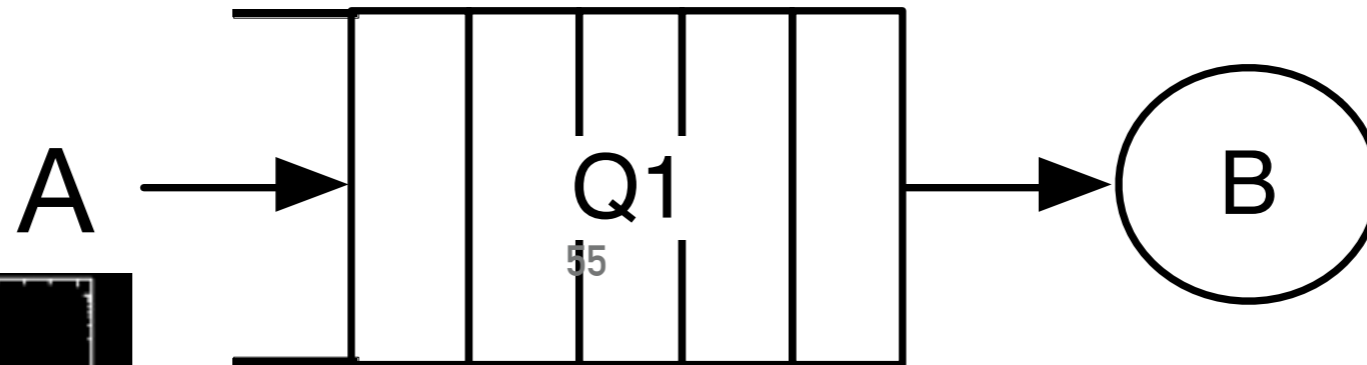


HOW WOULD YOU GET THE RIGHT BUFFER SIZE?

Queueing Model to Solve for Buffer Size



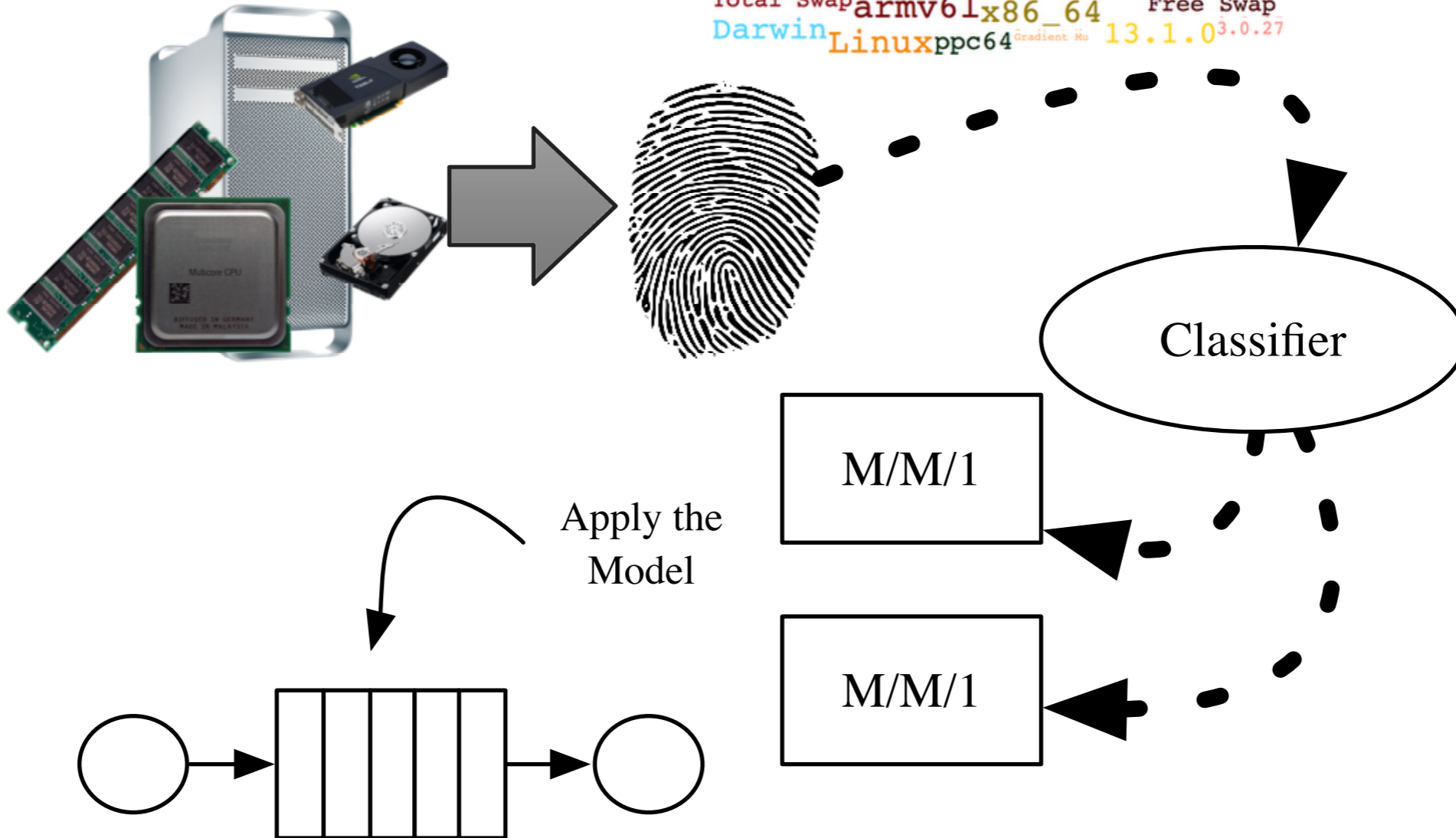
?



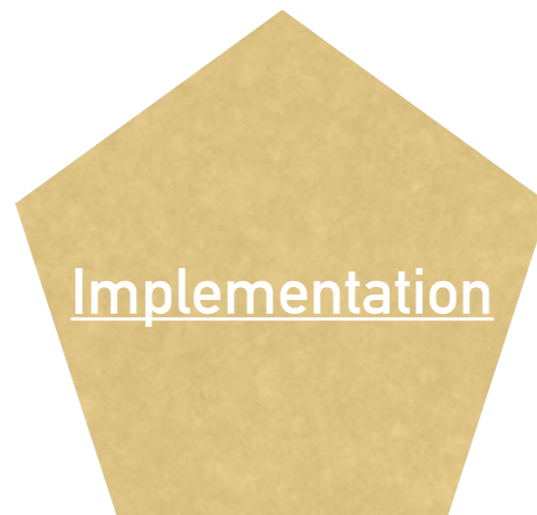
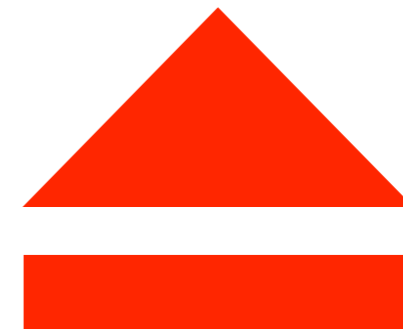
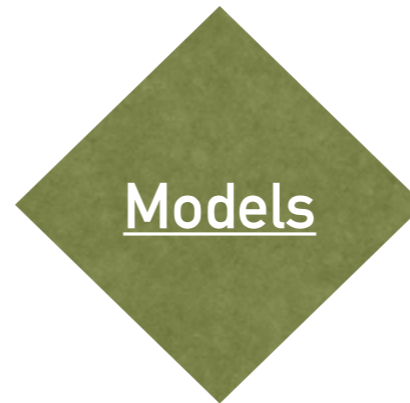
ML BASED MODEL SELECTION

```

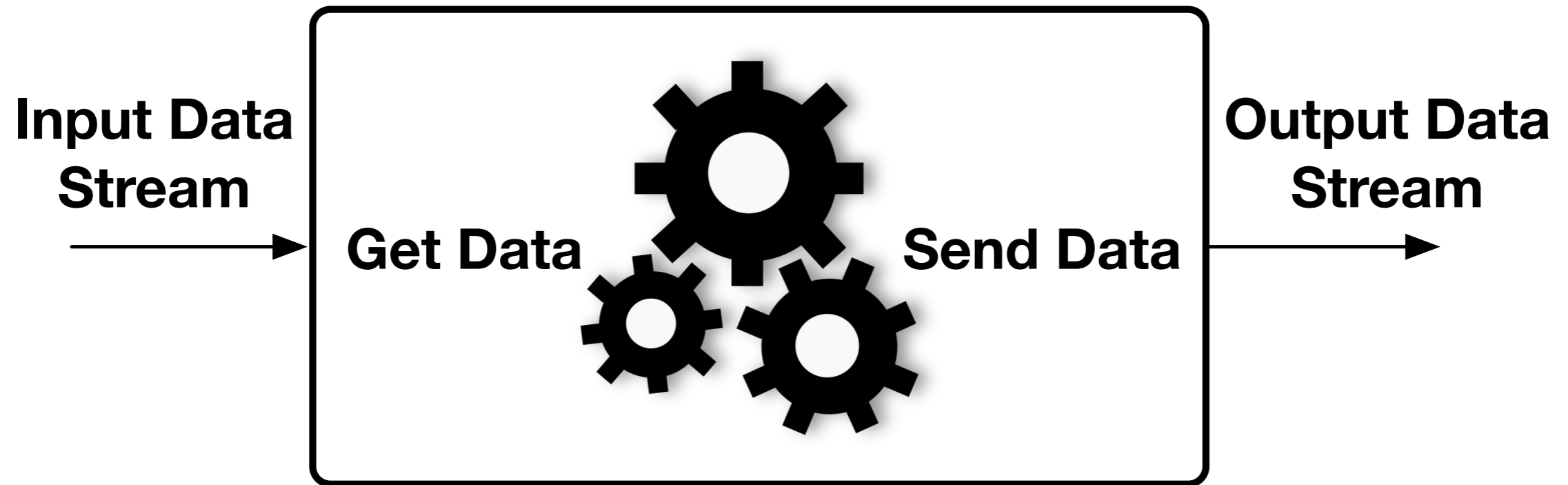
Dual Core AMD Opteron Processor 280 Level One D Cache Associativity
Dual Core AMD Opteron Processor 875 AMD Opteron Processor 6200
Level Three Cache Associativity Level Two Cache Associativity
Quad-Core AMD Opteron Processor 2387 Level Three Cache Line Size
2.6.32-358.23.2.el6.x86_64 Level Three Cache Size
Level One ICache Associativity x86_64 Number of Processors Free Ram
2.6.32-431.5.1.el6.x86_64
Quad-Core AMD Opteron Processor 2376 Level One DCache Size Gradient Lambda
traditional_with_pset_runqueue Level Two Cache Line Size
Intel Xeon CPU X5472 @ 3.00GHz Intel(R) Core i5 CPU M540 @ 2.53GHz
Dual-Core AMD Opteron Processor 2218 Level One ICache Line Size Arrival Process Service Time
Processor Frequency Level One ICache Line Size Level One DCache Line Size
Level One DCache Line Size Intel Xeon CPU E5345 @ 2.33GHz
Level Two Cache Size Dual-Core AMD Opteron Processor 2222 SE ARM1176JZF-S
Fifteen Min Load 3.13.0-24-PowerPc64-SMP Number of Processes Running
3.10.37+ Five Min Load SCHED_OTHER 3.2.14
One Min Load 2.6.32-358.11.1.el6.x86_64 Buffer Ram
Total Swap armv6l x86_64 Free Swap
Darwin Linux ppc64 Gradient Wu 13.1.0 3.0.27
    
```



CHOOSE YOUR ADVENTURE



COMPUTE KERNEL



COMPUTE KERNEL

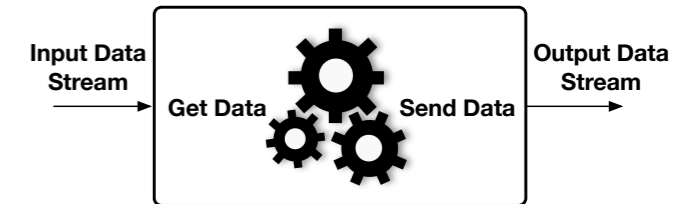
```
class akernel : public raft::kernel
{
public:
    akernel() : raft::kernel()
    {
        //add input ports
        input.addPort< /** type */ >( "x0", "x1", "x..." );
        //add output ports
        output.addPort< /** type */ >( "y0", "y1", "y..." );
    }

    virtual raft::kstatus run()
    {
        /** get data from input ports */
        auto &valFromX( input[ "x..." ].peek< /** type of "x..." */ >( ) );
        /** do something with data */

        const auto ret_val( do_something( valFromX ) );

        output[ "y..." ].push( ret_val );

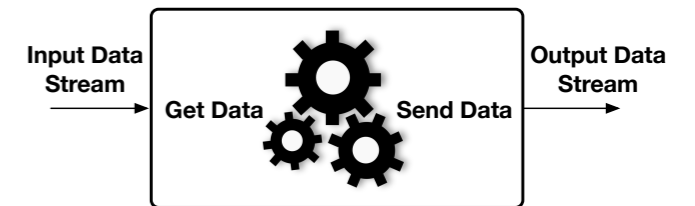
        input[ "x..." ].unpeek();
        input[ "x..." ].recycle();
        return( raft::proceed /** or stop */ );
    }
};
```



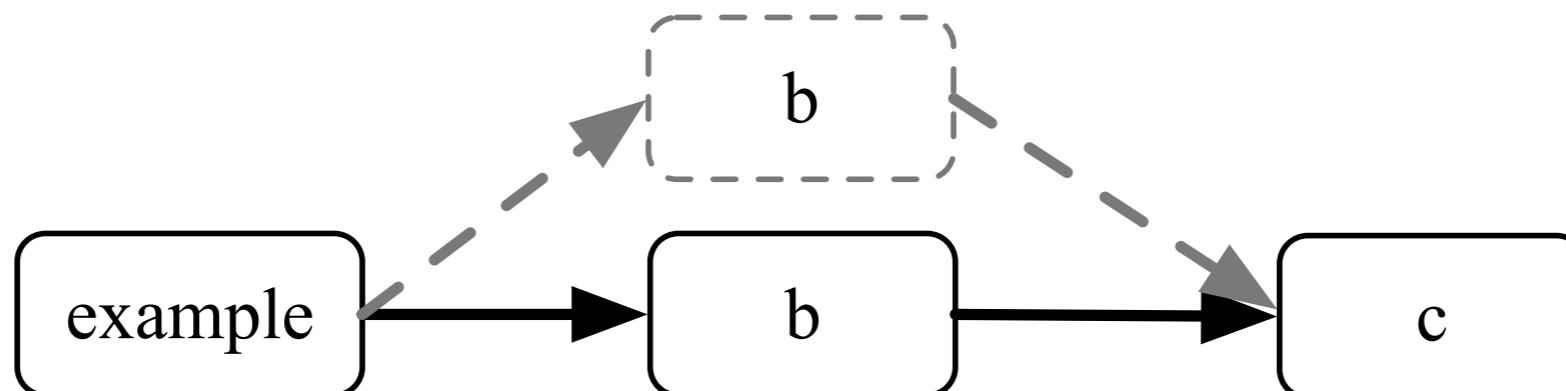
COMPUTE KERNEL

```
class example : public raft::parallel_k
{
public:
    example() : parallel_k()
    {
        input.addPort< /** some type */ >( "0" );
        /** add a starter output port */
        addPort();
    }

    /** implement virtual function */
    virtual std::size_t addPort()
    {
        return( (this)->addPortTo< /** type */ >( output /** direction */ ) );
    }
}
```



"..."

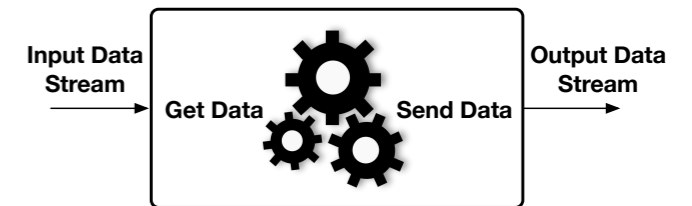


RECEIVING DATA

```
/**
 * return reference to memory on
 * in-bound stream
 */
template< class T >
T& peek( raft::signal *signal = nullptr )
```

```
template< class T >
autorelease< T, peekrange > peek_range( const std::size_t n )
```

- Returns object with “special access to stream”
- Operator `[]` overload returns auto pair
- Direct reference as in `peek()` for each element

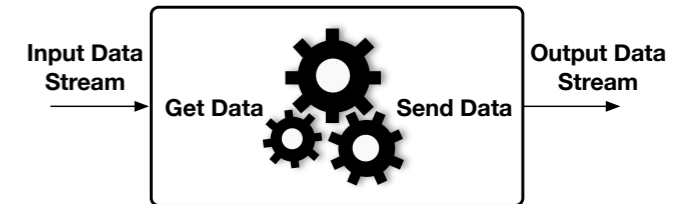
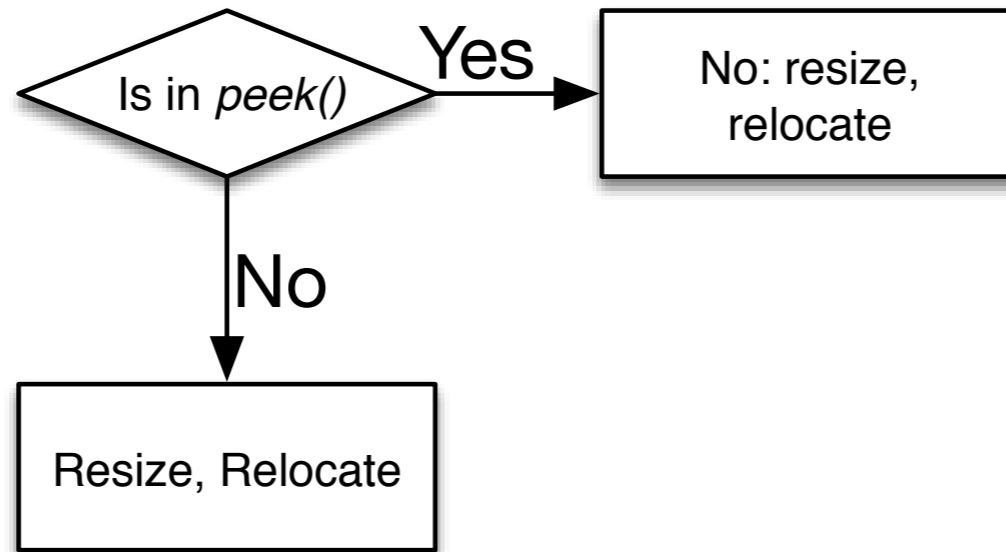


↓

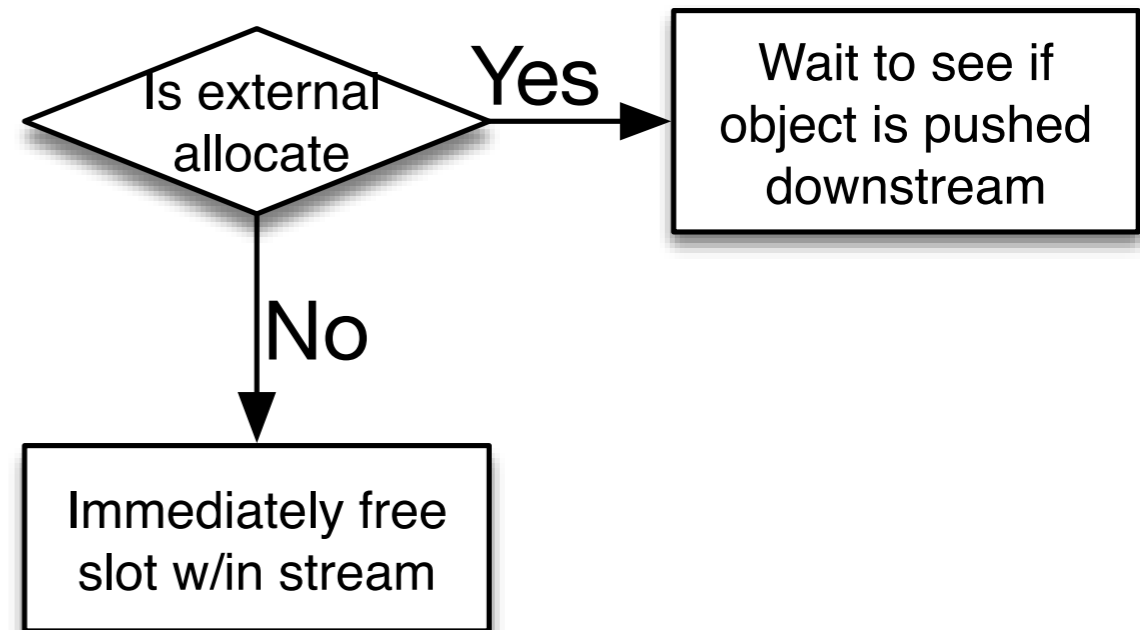
```
template < class T > struct autopair
{
    T &ele;
    Buffer::Signal &sig;
};
```

RECEIVING DATA

```
/**  
 * necessary to tell inbound stream we're done  
 * looking at it  
 */  
virtual  
void unpeek()
```



```
/**  
 * free up index in fifo, lazily deallocate large objects  
 */  
void  
recycle( const std::size_t range = 1 )
```



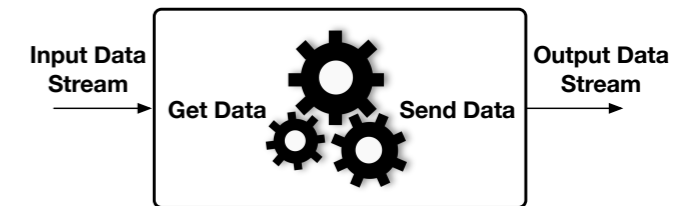
RECEIVING DATA

```
/**
 * these pops produce a copy
 */
template< class T >
void pop( T &item, raft::signal *signal = nullptr )

template< class T > using pop_range_t =
    std::vector< std::pair< T , raft::signal > >;

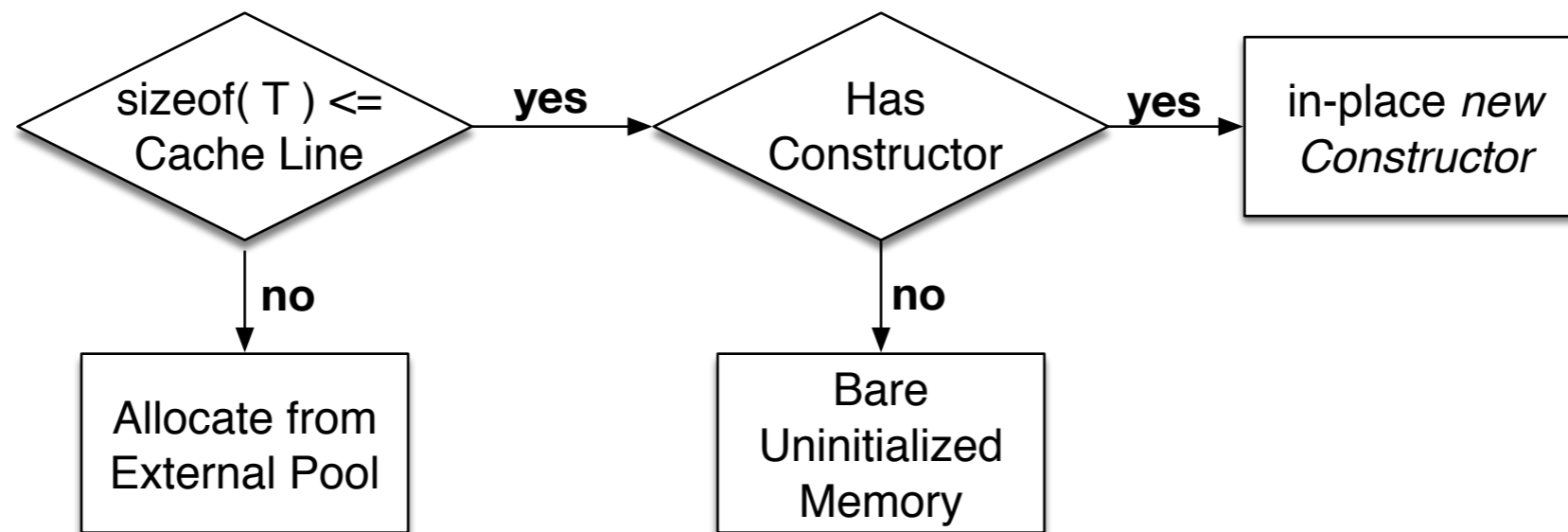
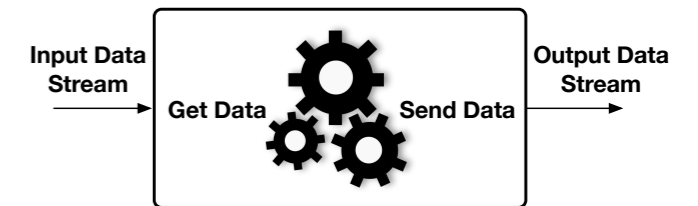
template< class T >
void pop_range( pop_range_t< T > &items,
               const std::size_t n_items )

/**
 * no copy, slightly higher overhead, "smart object"
 * implements peek, unpeek, recycle
 */
template< class T >
autorelease< T, poptype > pop_s()
```



SENDING DATA

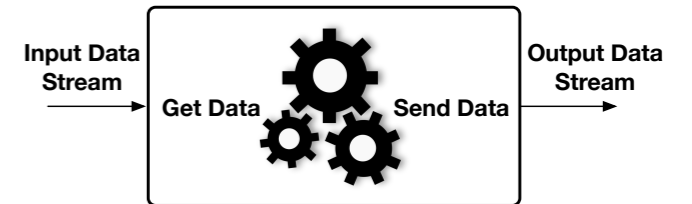
```
/** in-place allocation **/  
template < class T,  
          class ... Args >  
T& allocate( Args&&... params )
```



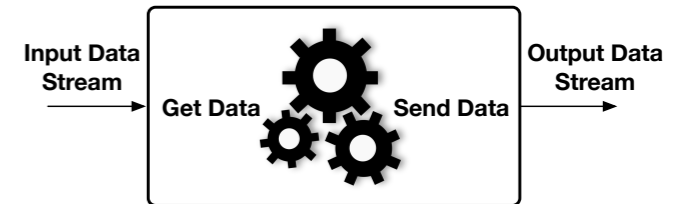
```
/** in-place alloc of range for fundamental types **/  
template < class T >  
auto allocate_range( const std::size_t n ) ->  
    std::vector< std::reference_wrapper< T > >
```


SENDING DATA

```
/** release data to stream **/  
virtual  
void send( const raft::signal = raft::none )  
  
/** release data to stream **/  
virtual  
void send_range( const raft::signal = raft::none )  
  
/** oops, don't need this memory **/  
virtual void deallocate()
```



SENDING DATA



```
/** multiple forms **/
template < class T >
void push( const T &item, const raft::signal signal = raft::none )
```

```
/** insert from container within run() function to stream **/
template< class iterator_type >
void insert(    iterator_type begin,
                iterator_type end,
                const raft::signal signal = raft::none )
```

INCLUDED KERNELS

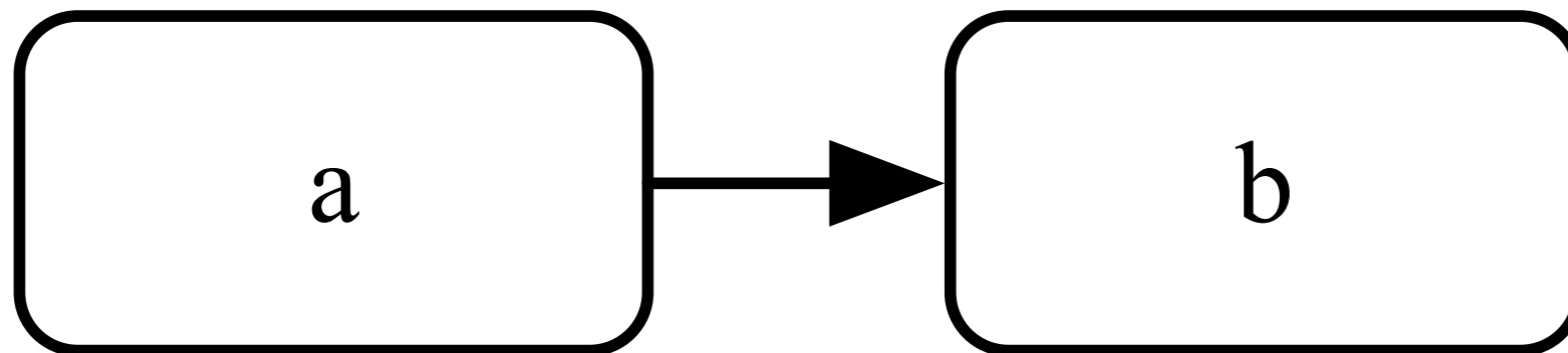
```
/**
 * thread safe print, specialization for '\n' vs. '\0'
 */
template< typename T, char delim = '\0' > class print

/** read from iterator to streams */
static
raft::readeach< T, Iterator >
read_each( Iterator &&begin,
           Iterator &&end )

/** write from iterator to streams */
template < class T, class BackInsert >
static
writeeach< T, BackInsert >
write_each( BackInsert &&bi )
{
    return( writeeach< T, BackInsert >( bi ) );
}
```

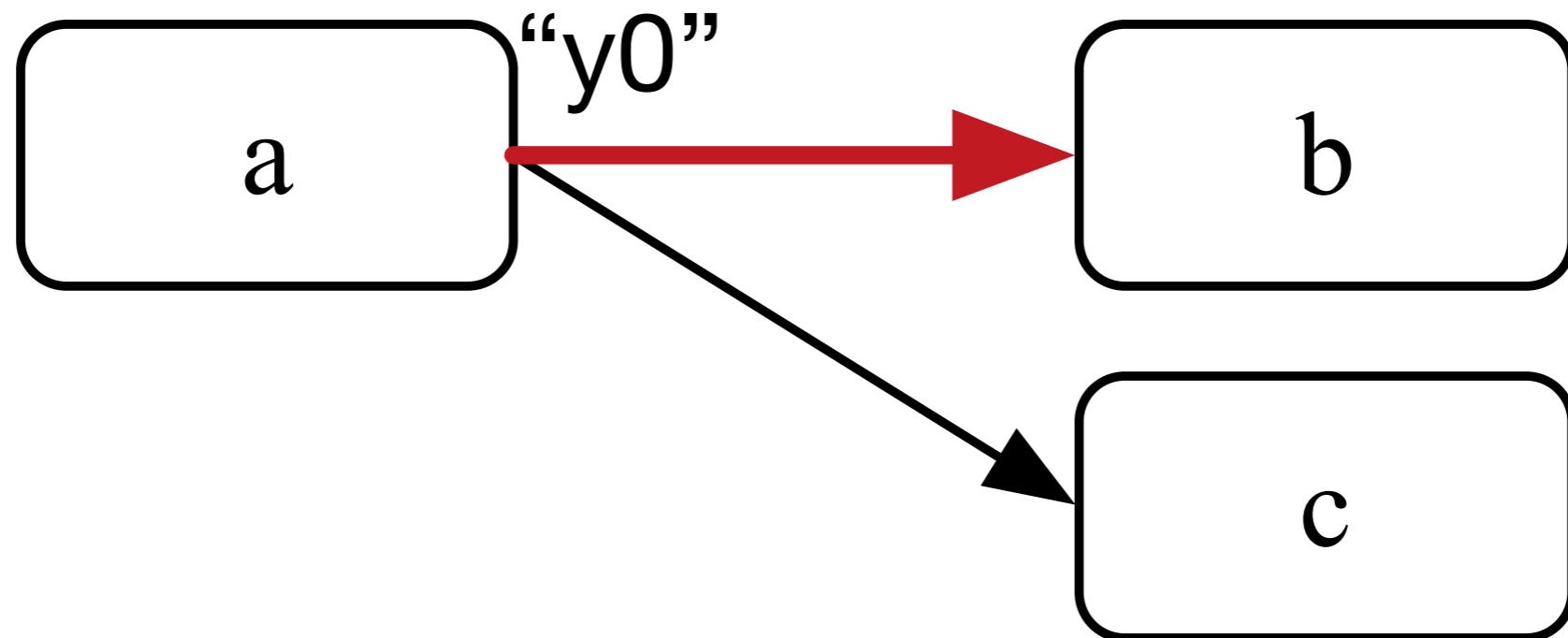
CONNECTING COMPUTE KERNELS

```
raft::map m;  
/** example only **/  
raft::kernel a, b;  
m += a >> b;
```



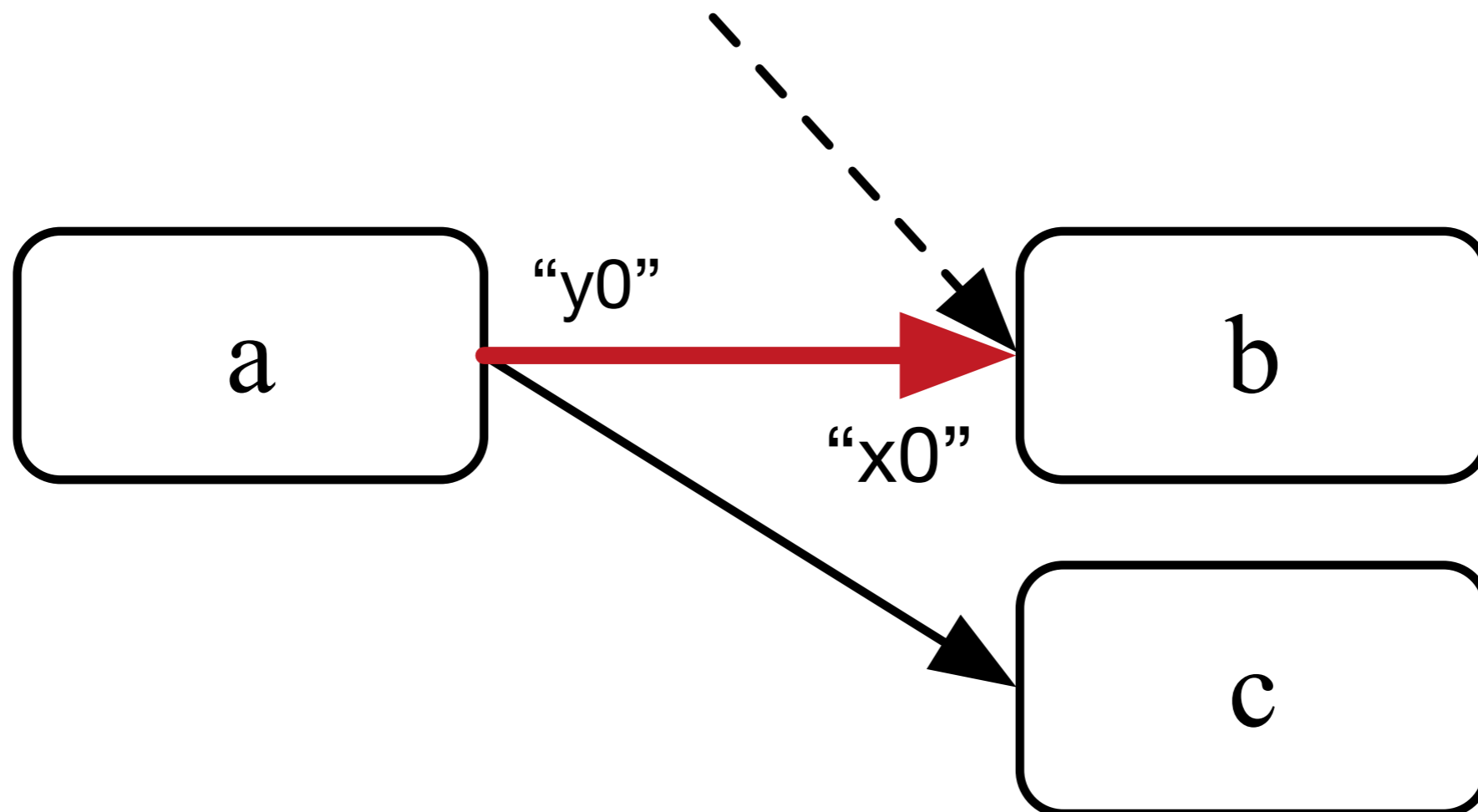
CONNECTING COMPUTE KERNELS

```
raft::map m;  
/** example only **/  
raft::kernel a, b;  
m += a[ "y0" ] >> b;
```



CONNECTING COMPUTE KERNELS

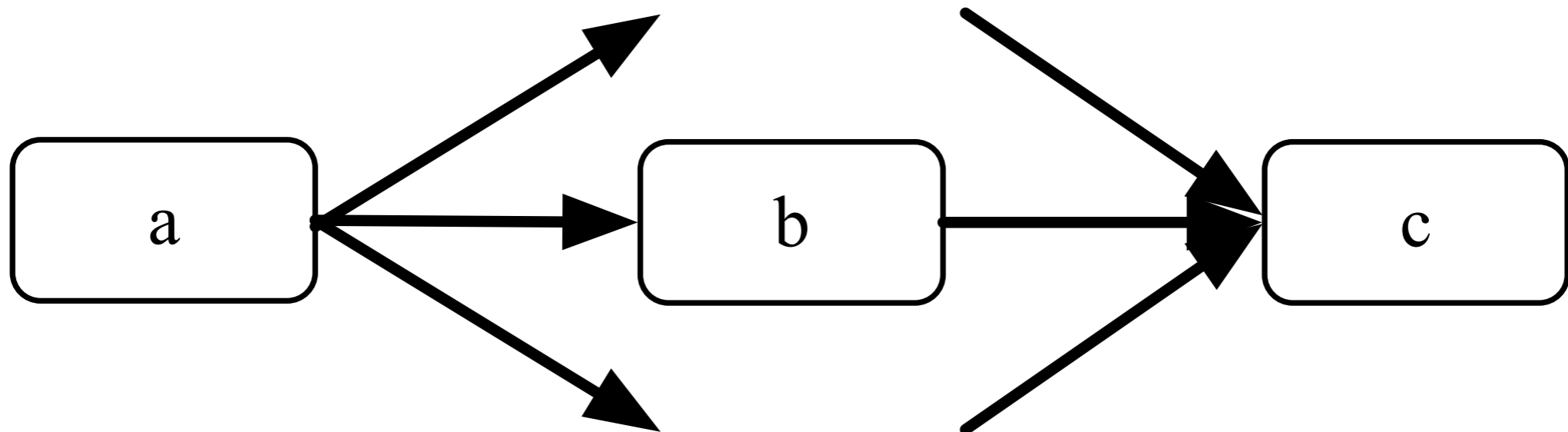
```
raft::map m;  
/** example only */  
raft::kernel a, b;  
m += a[ "y0" ] >> b[ "x0" ];
```



CONNECTING COMPUTE KERNELS

```
raft::map m;  
/** example only */  
raft::kernel a, b, c;  
m += a <= b >= c;
```

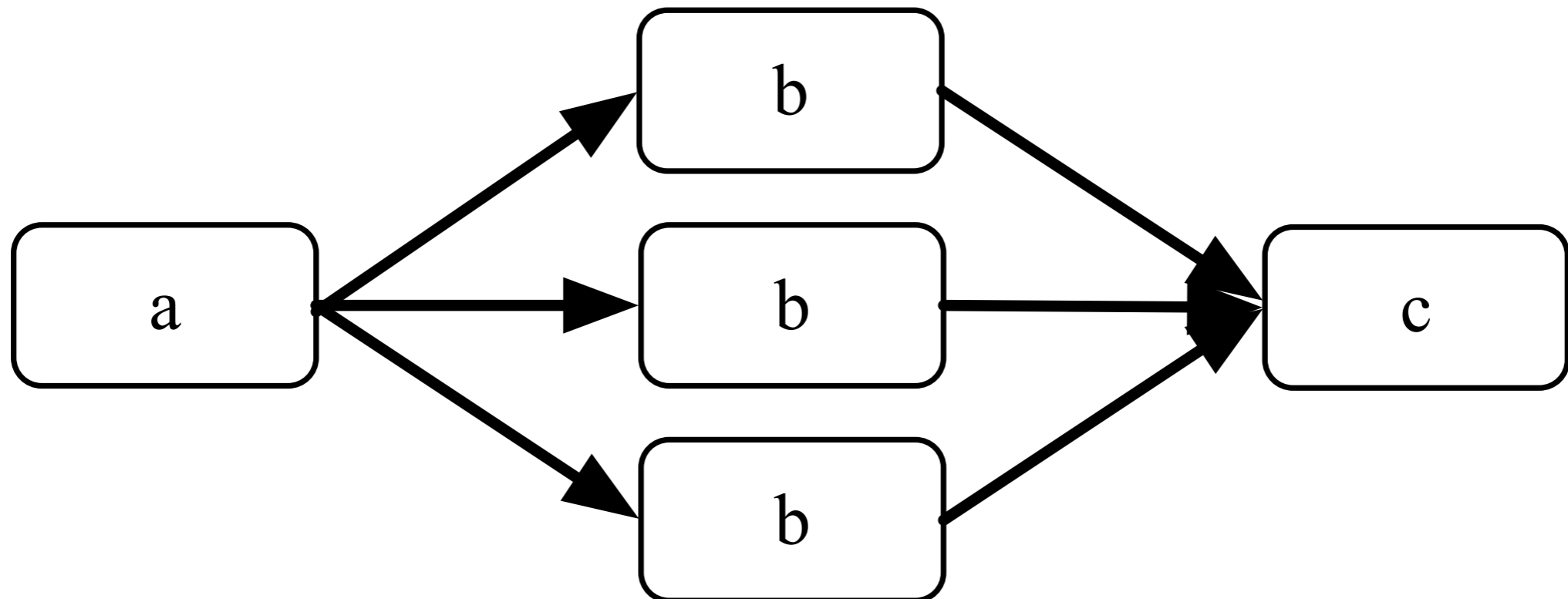
Topology user specifies



CONNECTING COMPUTE KERNELS

```
raft::map m;  
/** example only */  
raft::kernel a, b, c;  
m += a <= b >= c;
```

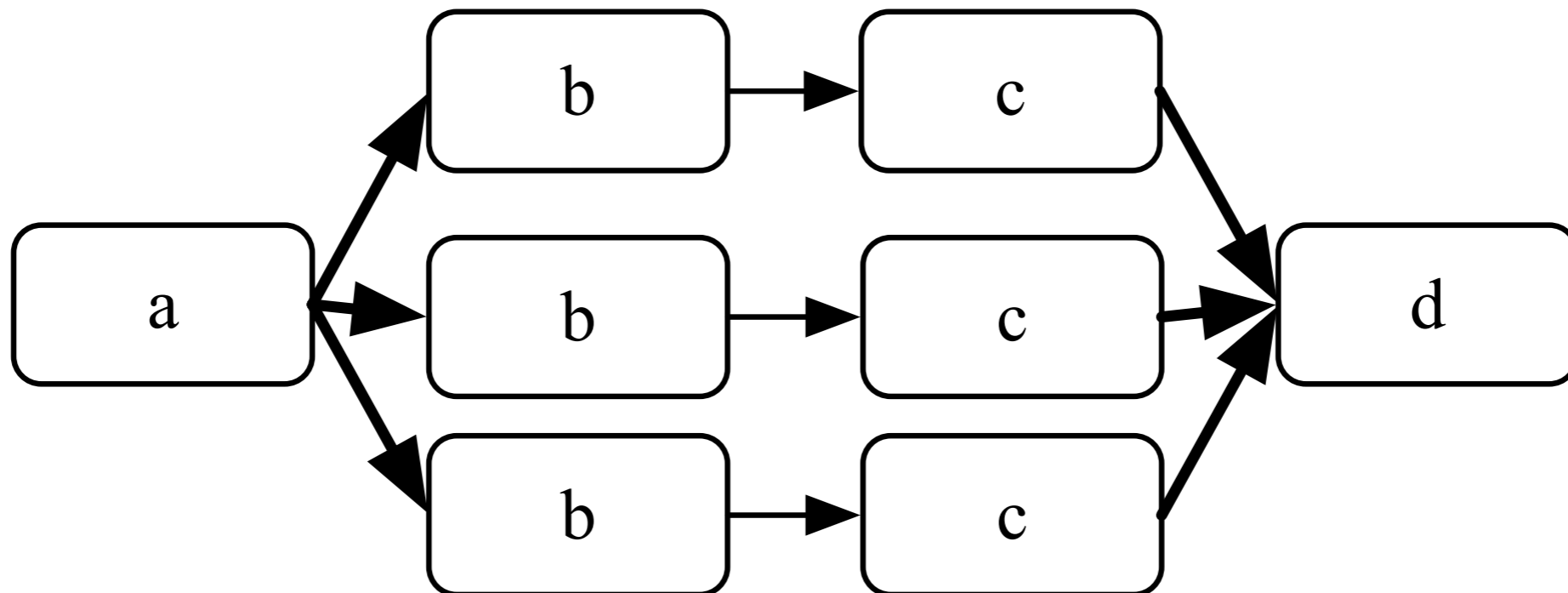
RaftLib Turns Into



CONNECTING COMPUTE KERNELS

```
raft::map m;  
/** example only */  
raft::kernel a, b, c, d;  
m += a <= b >> c >= d;
```

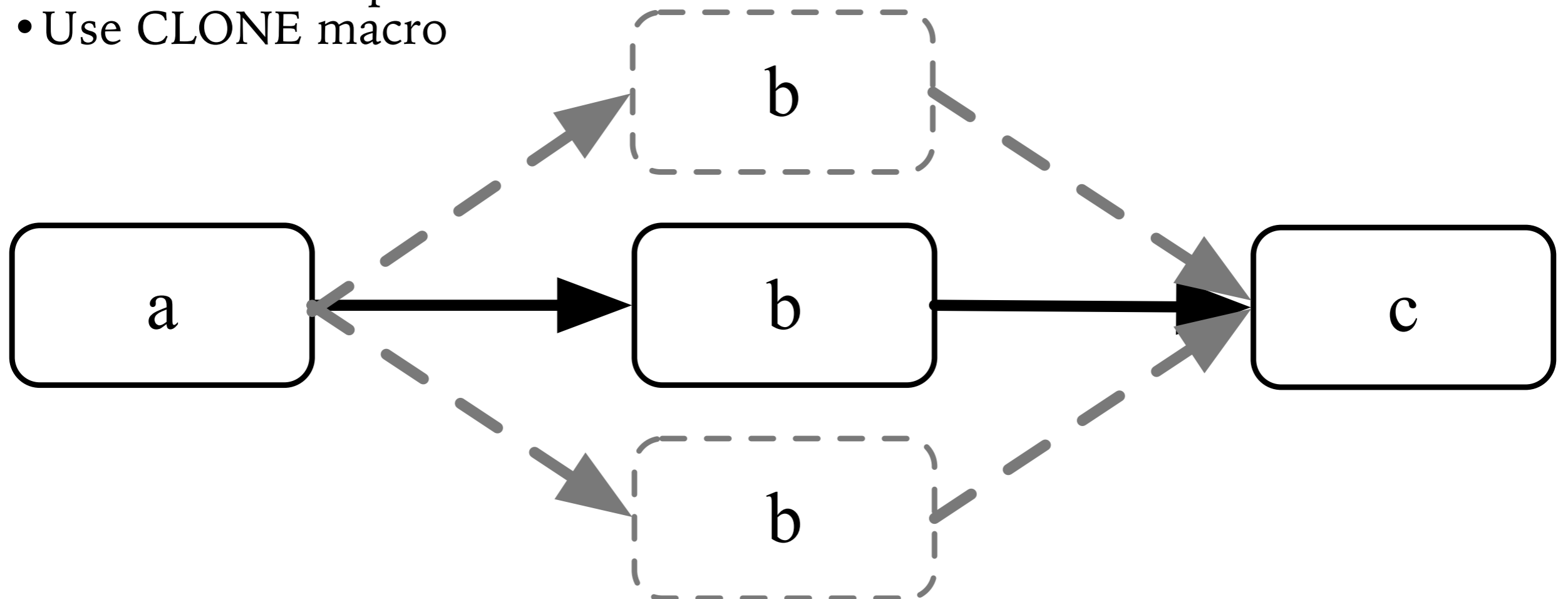
RaftLib Turns Into



CONNECTING COMPUTE KERNELS

```
raft::map m;  
/** example only **/  
raft::kernel a, b, c;  
m += a >> raft::order::out >> b >> raft::order::out >> c;
```

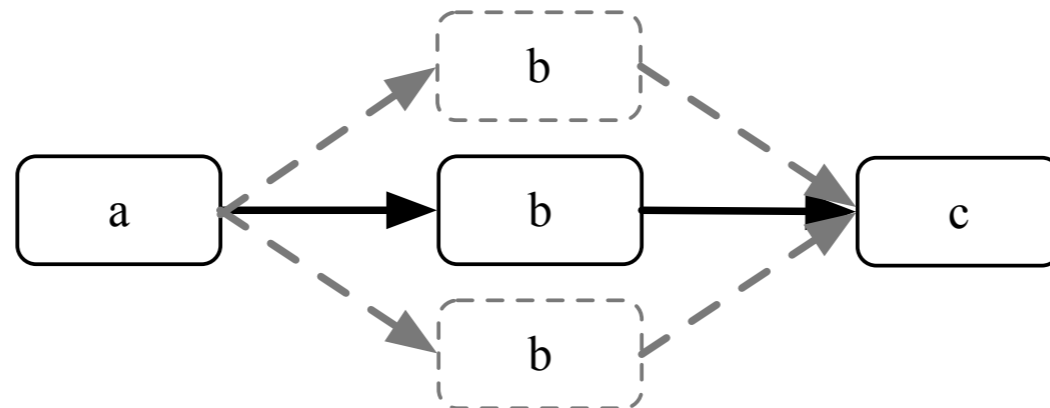
- Clone at **SESE** points
- Use CLONE macro



AUTO-PARALLELIZATION

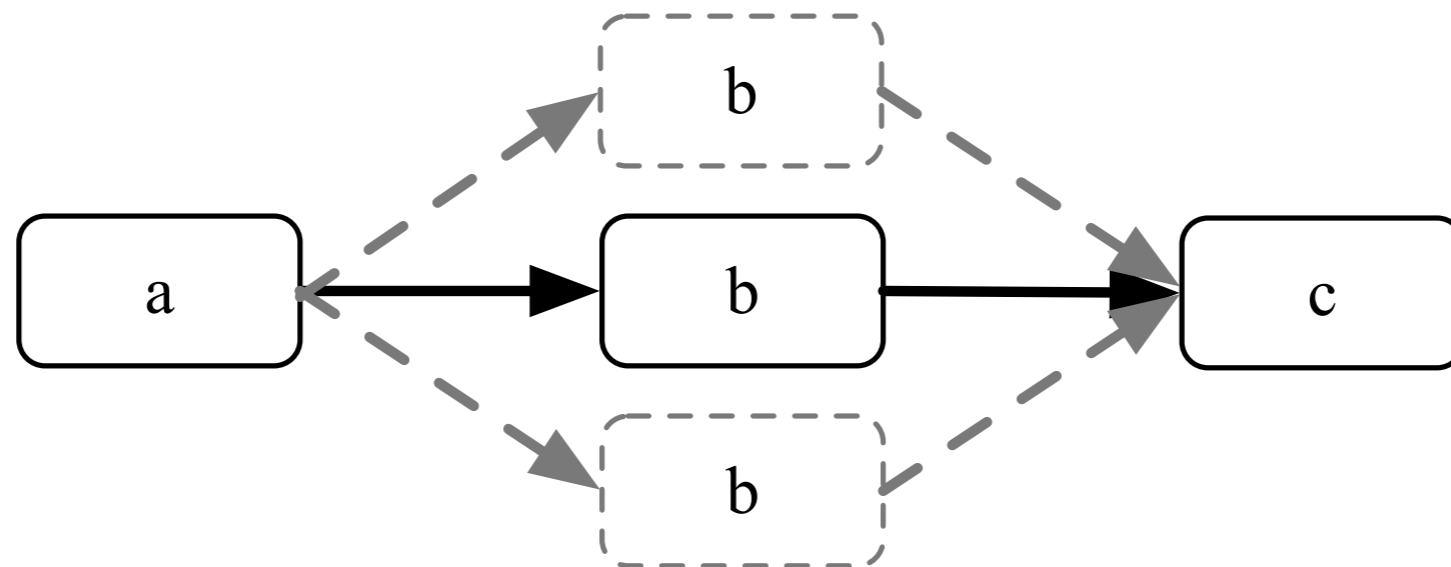
- 1) RaftLib figures out methodology
- 2) Runtime calls CLONE() macro of kernel “b”

```
#define CLONE()\nvirtual raft::kernel* clone()\n{\n    auto *ptr( \n        new typename std::remove_reference< decltype( *this ) >::type( ( *(\n            (typename std::decay< decltype( *this ) >::type * ) \n            this ) ) ) );\n    /** RL needs to dealloc this one **/\n    ptr->internal_alloc = true;\n    return( ptr );\n}
```



AUTO-PARALLELIZATION

- 3) Lock output port container of “a”
 - 4) Register new port
 - 5) Decide where to run it
 - 6) Allocate memory for stream
 - 7) Unlock output container of kernel “a”
- //do same on output side of “b”



CONNECTION TODO ITEMS

- Better SESE implementation
- Decide on syntax for set of kernels

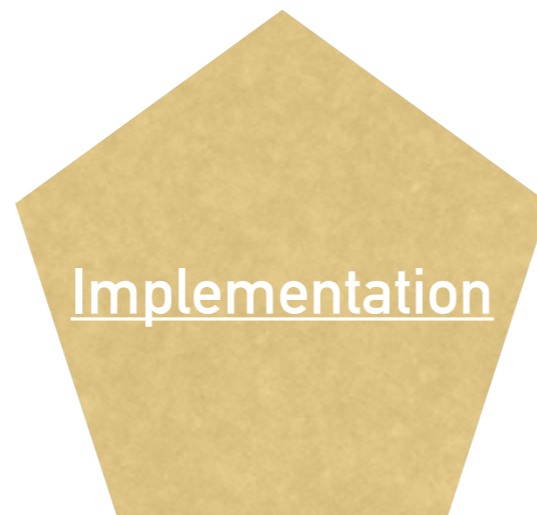
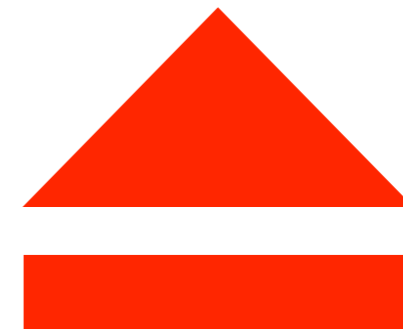
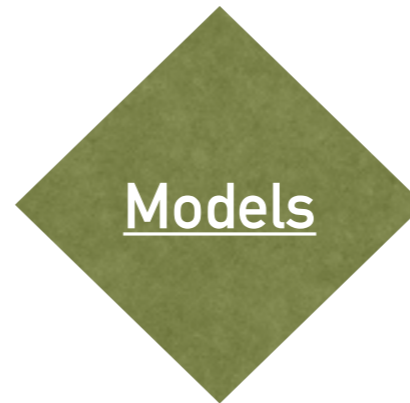
```
raft::map m;  
/** example only */  
raft::kernel a, b, c, d;  
m += src <= raft::kset( a, b, c, d ) >= dst;
```

- Address space stream modifier (new VM space)

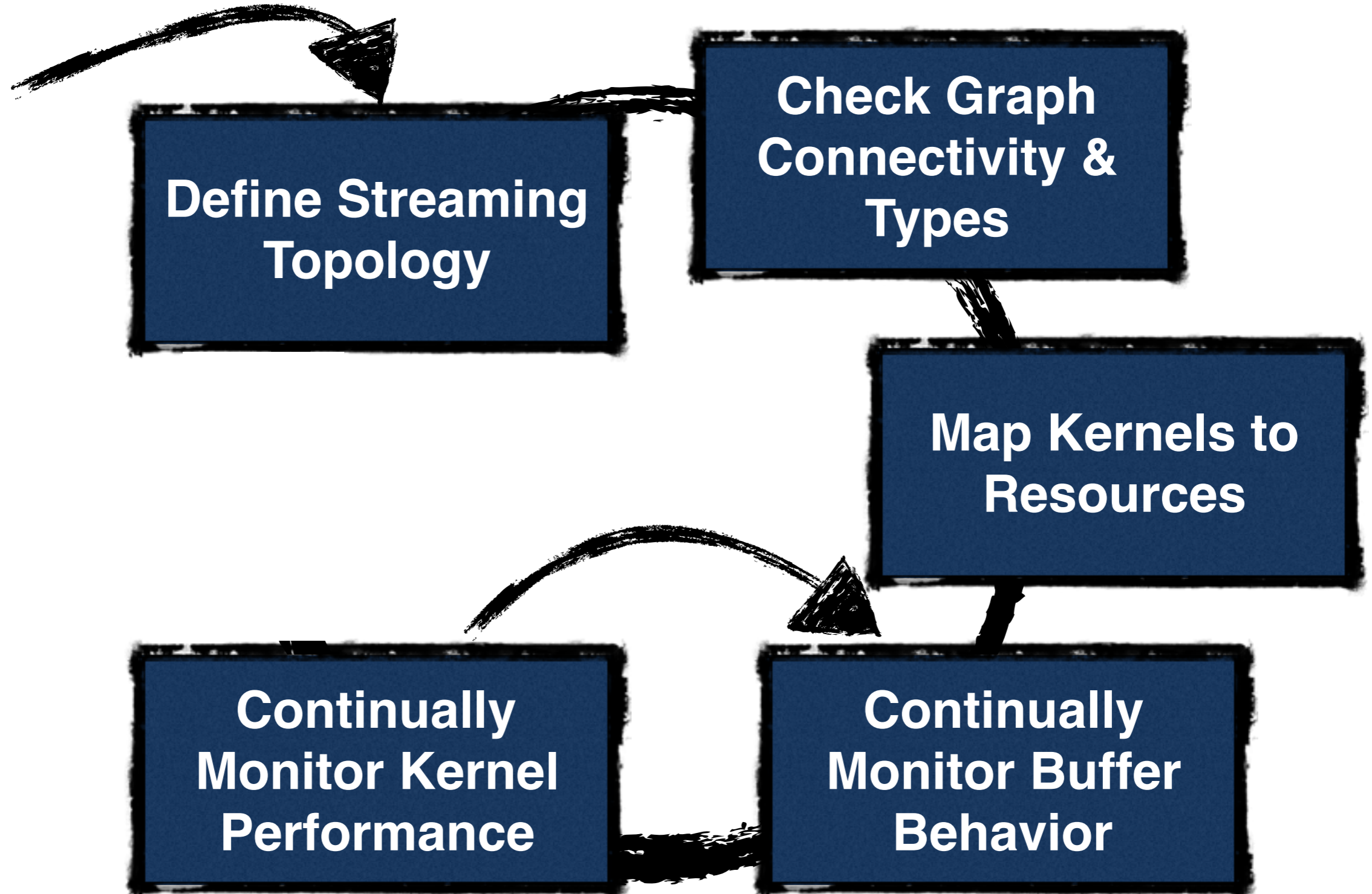
```
raft::map m;  
/** example only */  
raft::kernel a, b, c;  
m += a >> raft::vm::part >> b >> raft::vm::part >> c;
```

- Anything else?

CHOOSE YOUR ADVENTURE

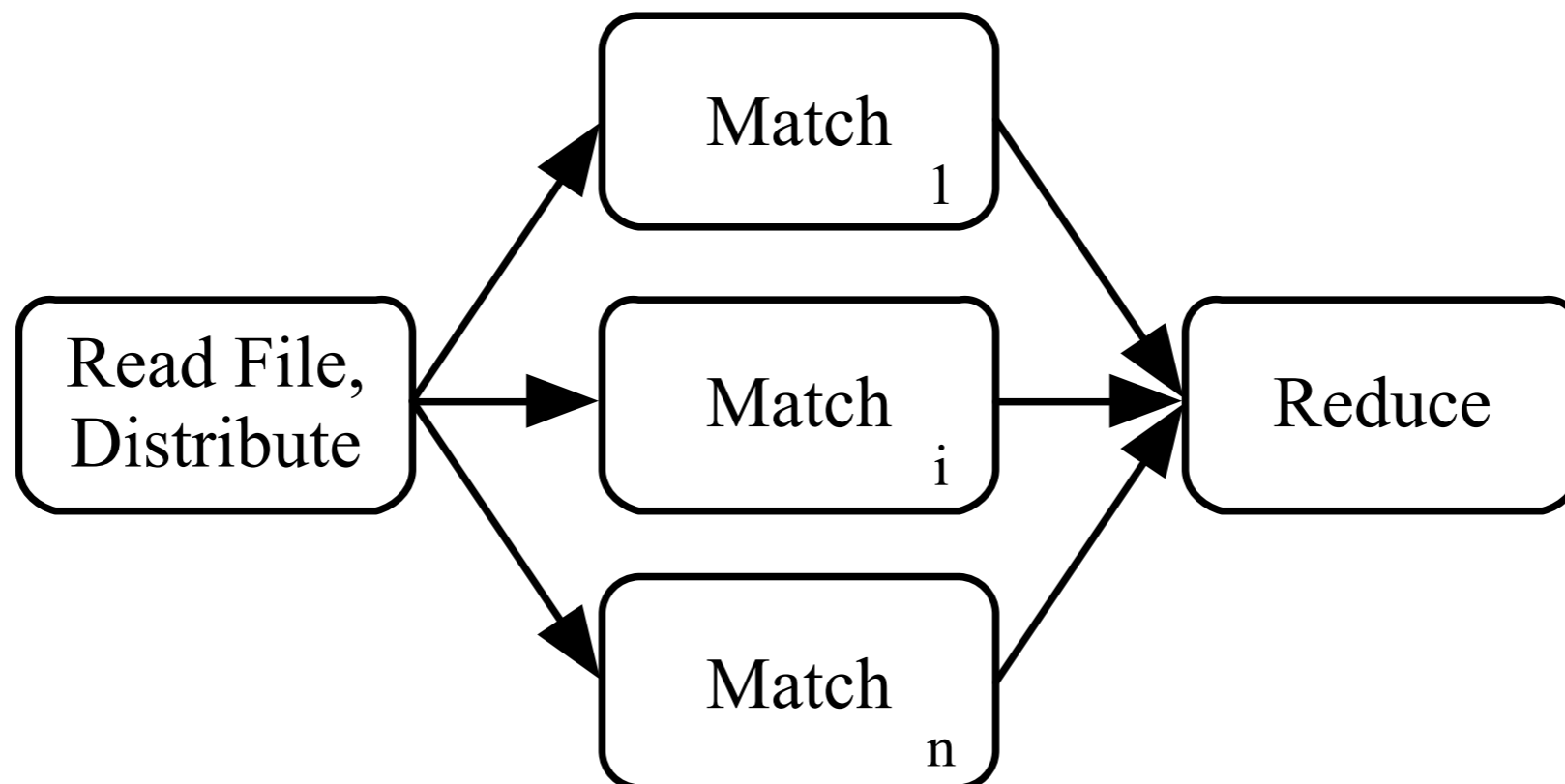


IMPLEMENTATION DETAILS



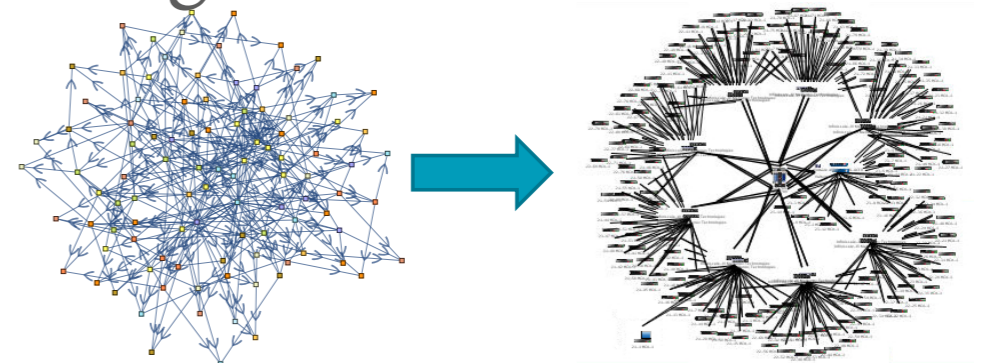
TOPOLOGY CHECK

- Add kernels to map
- Check type of each link (potential for type optimization)
- Handle static split/joins (produce any new kernels)
- DFS to ensure no unconnected edges



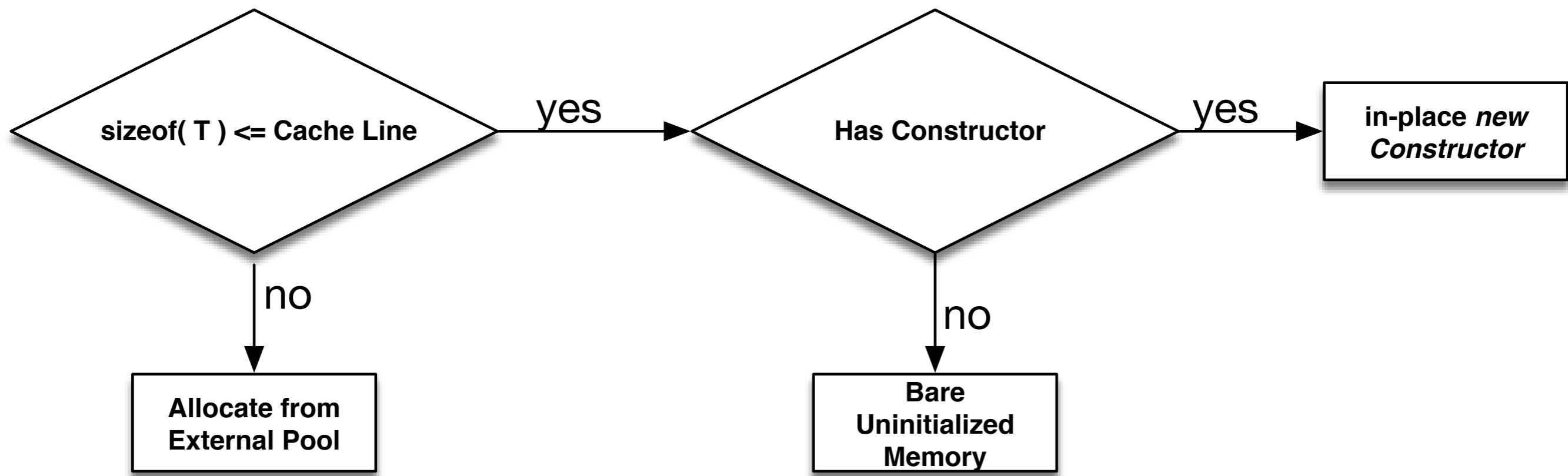
PARTITION (LINUX / UNIX)

- Take RaftLib representation, convert to Scotch format
- Use fixed communications cost at each edge for initial partition
- 2 main dimensions:
 - Flow between each edge in the application graph
 - Bandwidth available between compute resource
- Set affinity to partitioned compute cores
- Repartition as needed @ run-time
- TODO: incorporate OpenMPI utilities *hwloc* and *netloc* to get cross-platform hardware topology information



CHOOSE ALLOCATIONS

- Alignment
 - SIMD ops often require memory alignment, RaftLib takes an alignment by default approach for in-stream allocations
- In-stream vs. External Pool Allocate (template allocators)



MONITOR BEHAVIOR – ALLOCATORS

- Two options for figuring out optimal buffer size while running
 - model based (*discussed on modeling adventure path*)
 - branch & bound search
- separate thread, exits when app done
- pseudocode:

```
while( not done )
{
    if( queue_utilization > .5 )
    {
        queue->resize();
    }
    sleep( ALLOC_INTERVAL );
}
```

RUN-TIME LOCK FREE FIFO RESIZING

```
.....
enum access_key : key_t { allocate      = 0,
                           allocate_range = 1,
                           push         = 3,
                           recycle      = 4,
                           pop          = 5,
                           peek        = 6,
                           size        = 7,
                           N };
}
struct ThreadAccess
{
    union
    {
        std::uint64_t whole = 0; /** just in case, default zero **/
        dm::key_t      flag[ 8 ];
    };
    std::uint8_t padding[ L1D_CACHE_LINE_SIZE - 8 /** padding **/ ];
}
#if defined __APPLE__ || defined __linux
__attribute__((aligned( L1D_CACHE_LINE_SIZE )))
#endif
volatile thread_access[ 2 ];

/** std::memory_order_relaxed **/
std::atomic< std::uint64_t > checking_size = { 0 };

```

RUN-TIME LOCK FREE FIFO RESIZING

- Optimization...wait for the right conditions

```
if( rpt < wpt )
{
    //perfect to copy w/std::memcpy
}
```

- Factory allocators

```
/** allocator factory map */
std::map< Type::RingBufferType , instr_map_t* > const_map;

/** initialize some factories */
const_map.insert( std::make_pair( Type::Heap , new instr_map_t() ) );

const_map[ Type::Heap ]->insert(
    std::make_pair( false /** no instrumentation */,
                   RingBuffer< T, Type::Heap, false >::make_new_fifo ) );
const_map[ Type::Heap ]->insert(
    std::make_pair( true /** yes instrumentation */,
                   RingBuffer< T, Type::Heap, true >::make_new_fifo ) );
...many more
```

MONITOR BEHAVIOR – PARALLELIZATION

- Mechanics covered in interface, simple model here
- Run in separate thread, term on exit

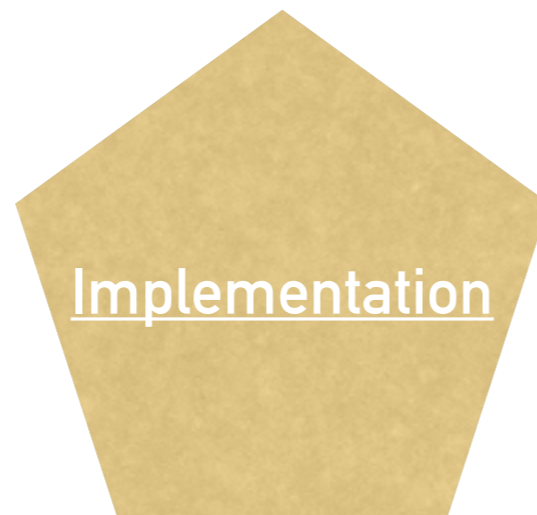
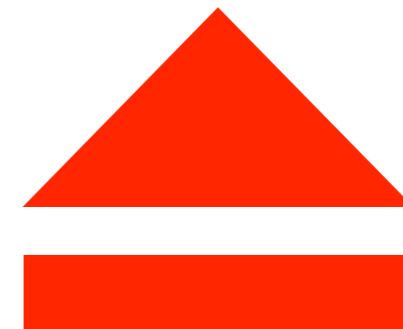
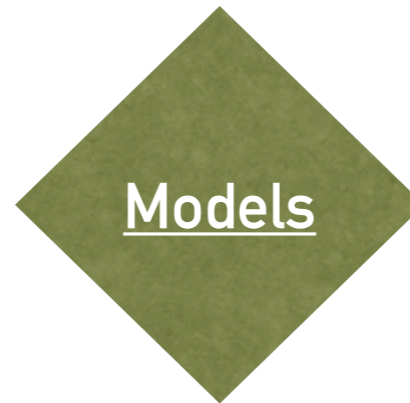
```
/** apply criteria **/
if( in_utilization > .5 && out_utilization < .5 )
{
    //tag kernel
    auto &tag( ag[ reinterpret_cast< std::uintptr_t >( kernel ) ] );
    tag += 1;
    if( tag == 3 )
    {
        dup_list.emplace_back( kernel );
    }
}

//after checking all kernels, handle duplication
```

IMPLEMENTATION TODO ITEMS

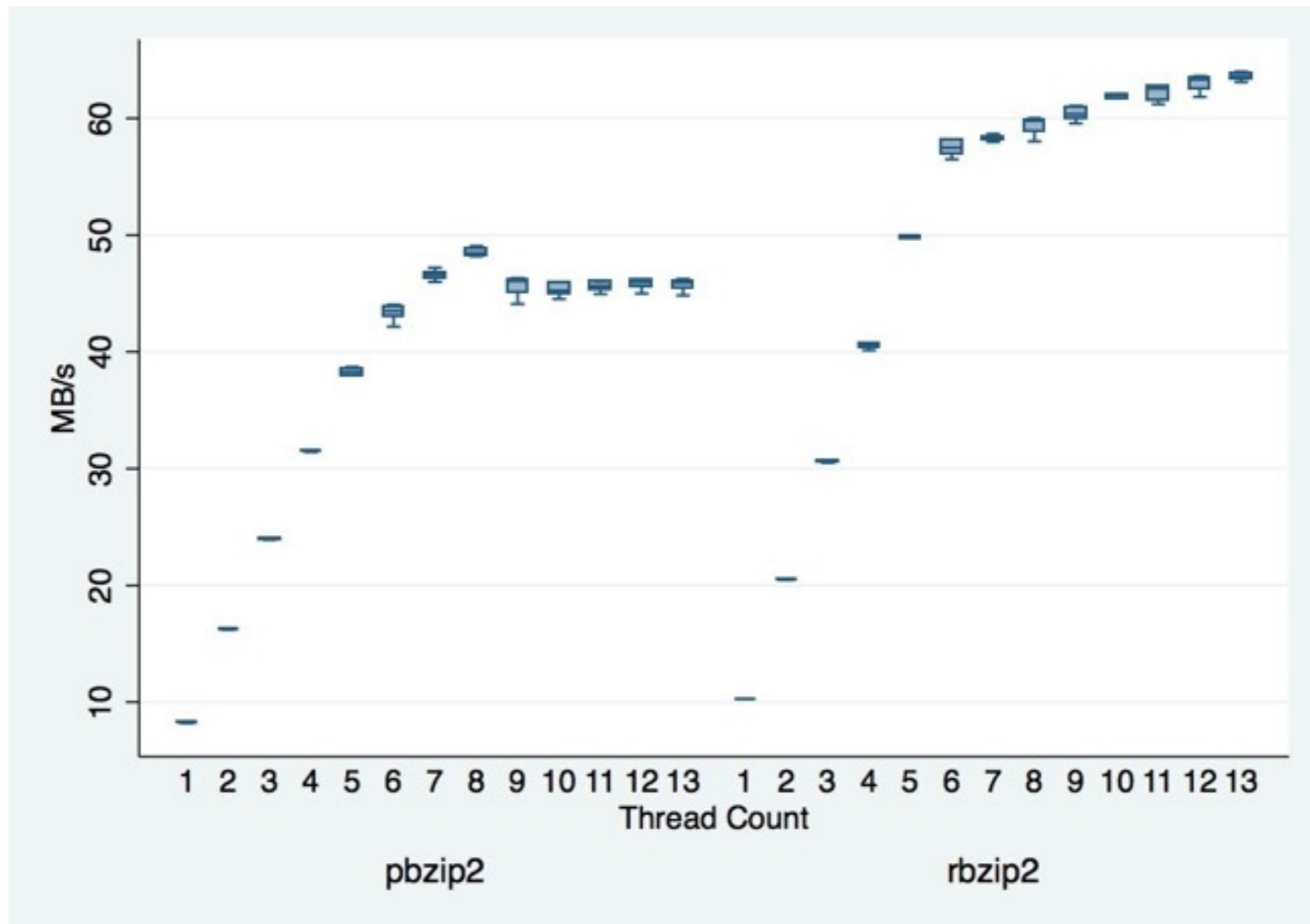
- Find fast SVM library to integrate (research code used LibSVM) for buffer model selection
- Integrate more production-capable network flow model for run-time re-partitioning choices
- Performant TCP links....
- RDMA on wish list
- QThreads Integration (see pool scheduler)
- *hwloc* and *netloc* integration (see `partition_scoth`)
- Perf data caching (useful for initial partition)

CHOOSE YOUR ADVENTURE



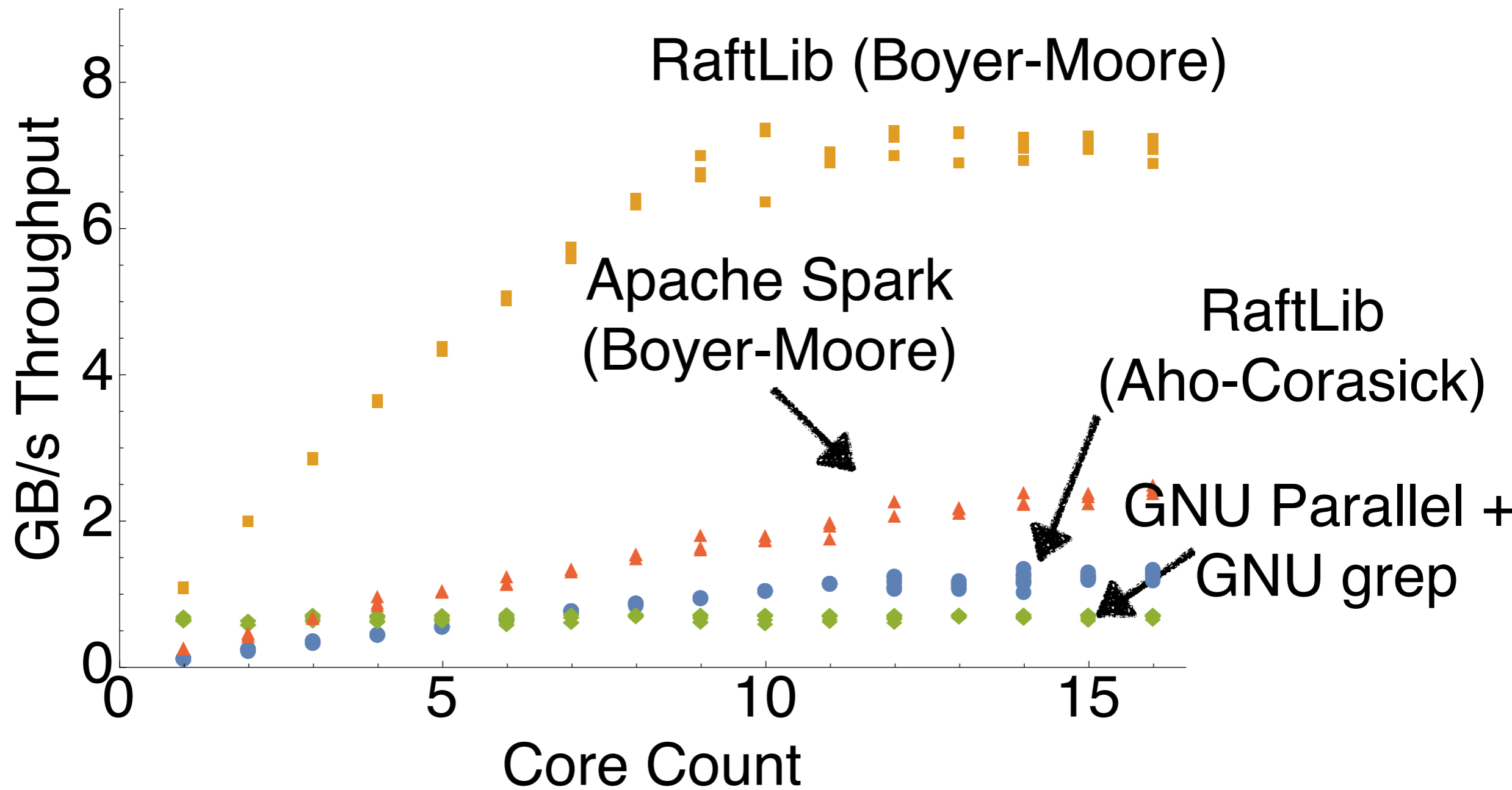
PERFORMANCE

- Decent compared to pthread stock implementation of pbzip2
- Parallel Bzip2 Example: <https://goo.gl/xyQAhm>



PERFORMANCE

- Fixed string search compared to Apache Spark, GNU Parallel + GNU Grep





ABOUT ME

my website

<http://www.jonathanbeard.io>



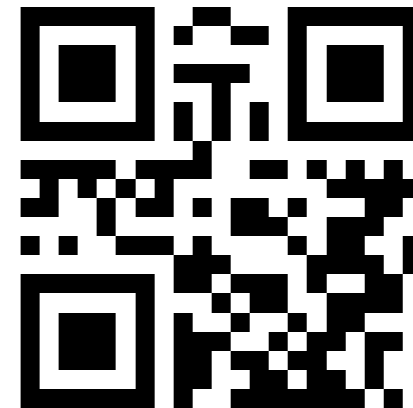
slides at

<http://goo.gl/cwT5UB>



project page

raftlib.io





video of talk given at #CppNow2016

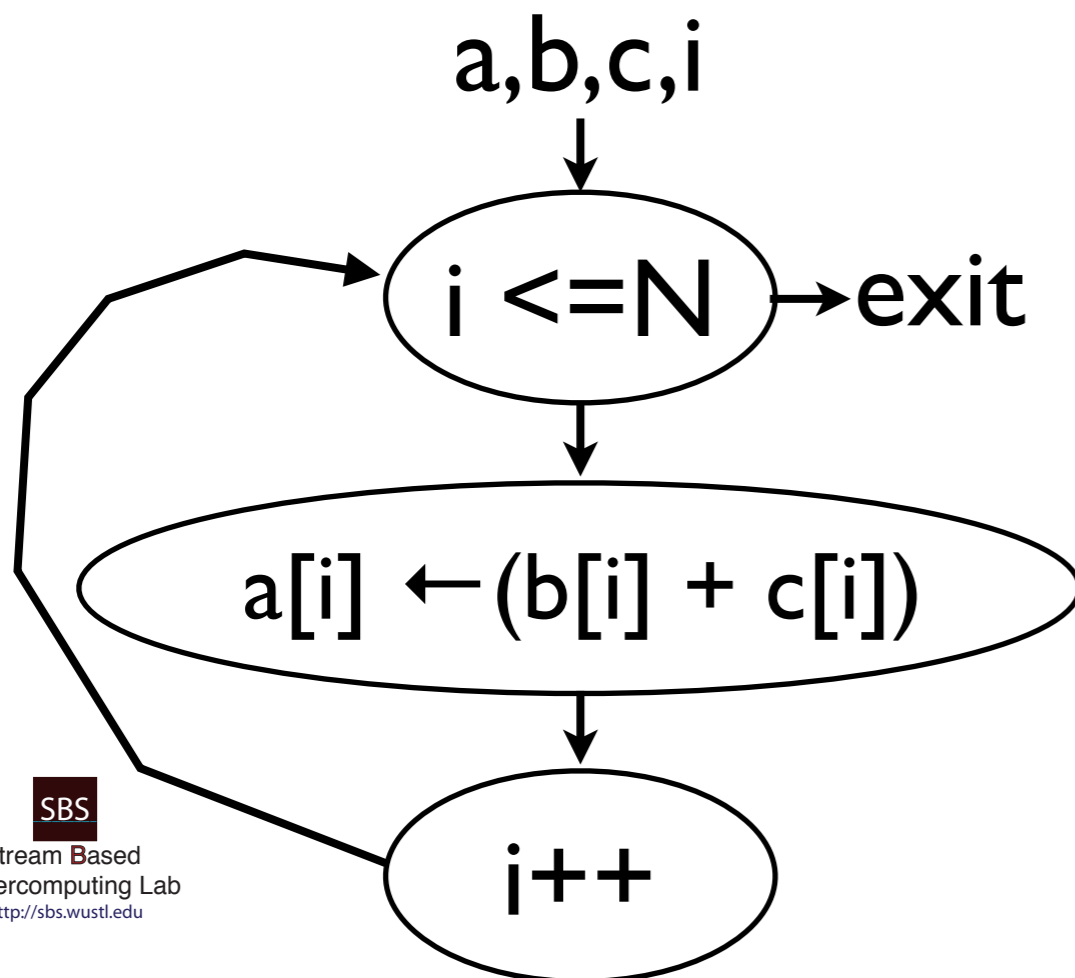
<http://goo.gl/mbxAwK>



Stream Processing

➔ for $i \leftarrow 0$ through N do
 $a[i] \leftarrow (b[i] + c[i])$
 $i++$
end do

Traditional Control Flow



Streaming

