

# A DESIGN FOR INTERCHANGABLE SIMULATION AND IMPLEMENTATION

---

Klaus Birkelund Jensen    Brian Vinter

August 25, 2015

Niels Bohr Institute

# OUTLINE

## 1. Introduction, background and motivation

Some context to understand why ISI was developed.

## 2. The current state of storage simulation

What techniques are we using today, and what are the advantages and disadvantages?

## 3. Our approach to interchangeability

What is interchangeability in simulation and implementation?

## 4. Scalability results

What makes the ISI approach viable for large scale (storage) simulation?

## 5. Summary

# INTRODUCTION AND MOTIVATION

---

# MOTIVATION

To understand how large scientific data sets can be stored efficiently.

Efficiency in

- Performance
- Resources usage
- Locality
- **Energy consumption**

We focus on energy consumption.

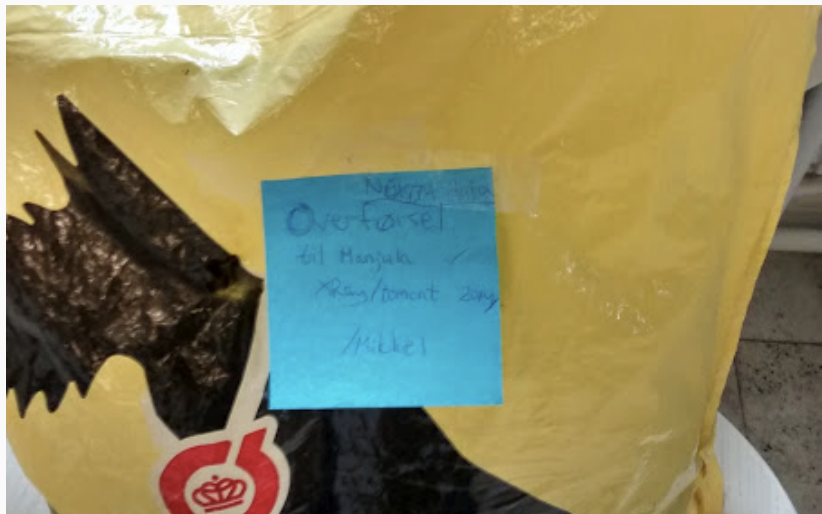
## ABOUT ME

Former systems operator at HPC/UCPH. Did storage and compute.

- Nordic T1 facility (storage & compute for ATLAS and ALICE)
- Multi PB disk, multi PB tape, thousands of compute cores.

Now, PhD student on the CINEMA project, working on storage techniques.

# MOTIVATION



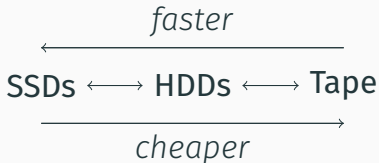
# THE PROBLEMS

The energy bill associated with storage is an ever larger part of the data center budget.

Most common technique to reduce energy consumption and maximize performance:

- **Hierarchical Storage Management (HSM)**

The notion of managing data according to popularity, age, size etc. Move *passive* data to cheaper lower tier storage (usually tape).



# THE PROBLEMS

HSM uses many reasonably good techniques including (but not limited to):

- LRU-caching and aging
- Manual tagging of data (i.e. “please do NOT move my data!”).
- Generally, *on-demand* retrieval. No prediction.



# THE PROBLEMS

HSM is too general to efficiently store what we define as *known data sets*.

We focus on scientific and industrial tomography imaging.

Imaging data exhibits known workloads and structure.

**We should acknowledge and exploit that.**

## THE PROBLEMS

In the data center, durability and reliability is most commonly provided by large RAID systems, but *erasure codes* are rapidly gaining traction.

In RAID, all drives must spin simultaneously. There are solutions to this in the literature, including:

- *Power-aware* RAID (or *gear shifting*).
- *Intelligent* data placement (e.g. locality optimized).

They are all general in nature.

## THE PROBLEMS

The principle of “optimizing for the common case” has always been a good strategy.

“This data was just used —  
let’s keep it around for a week... *or so*”

## THE PROBLEMS

The principle of “optimizing for the common case” has always been a good strategy.

“This data was just used —  
let’s keep it around for a week... *or so*”

But, the *common case* isn’t at all *common* when working with well-defined scientific data.

## THE PROBLEMS

The principle of “optimizing for the common case” has always been a good strategy.

“This data was just used —  
let’s keep it around for a week... *or so*”

But, the *common case* isn’t at all *common* when working with well-defined scientific data.

“This data has just been acquired —  
the physicist won’t use it for months... *if ever*”

What is possible if we exploit what is *known*?

- Raw data can be moved directly to tape
- Stream filtering

But how to quantify any possible benefits?

**Simulation** of storage hierarchies, workloads and data acquisition and consumption.

## BUILDING A STORAGE SYSTEM

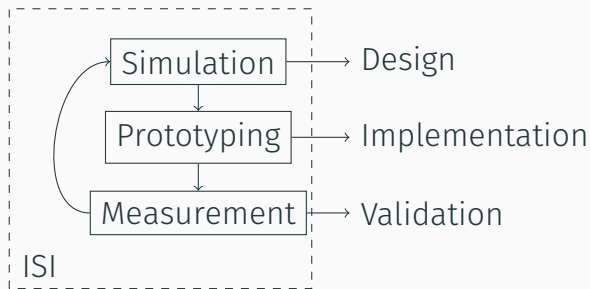
Developing a large-scale storage system where the design isn't exactly known in advance, could go something like this:

1. *Simulate* a model and identify a design.
2. *Implement* a prototype from the design.
3. *Measure* the prototype and *validate* it and the model against predictions.
4. *Repeat*. Feed the results of the validation back into the simulator and/or model and repeat from step 1.

The process is sound, but can we improve it?

# IMPROVEMENTS

Interchangeability of simulation and implementation  
Eliminating the *simulation-prototyping-measure* cycle.





Simulate the system model using *Discrete Event Simulation* (DES).

A DES is a priority queue of events, handled sequentially. Each event has a time stamp, updates the model and adds new events to the queue when handled.

Main loop of a DES.

---

## Algorithm 1 Discrete Event Simulation

---

```
1: procedure DES-LOOP( $Q$ )
2:   while  $Q \neq \emptyset$  do
3:      $e \leftarrow$  DEQUEUE( $Q$ )
4:      $T \leftarrow$  CLOCK( $e$ )  $\triangleright$  update world clock
5:     PROCESS( $e$ )  $\triangleright$  process event and add new
6:   end while
7: end procedure
```

---

An *event* is processed by a handler. Typically a huge function with a single switch-statement.

Parallel DES (PDES), generalizes this by allowing multiple processes to have a local priority queue.

ROSS (Rensselaer's Optimistic Simulation System) is an *optimistic* PDES.

- Extremely high performance
- Runs on millions of cores
- Relies on *Reverse Computation*

ROSS (Rensselaer's Optimistic Simulation System) is an *optimistic* PDES.

- Extremely high performance
- Runs on millions of cores
- Relies on *Reverse Computation*

**In summary: a savage beast**

# INTERCHANGEABLE SIMULATION AND IMPLEMENTATION

Model the system components as the individual processes they are.

The process logic directly implements a prototype.

Requires an environment supporting millions of independently communicating processes:

- Language based: Go, Erlang, occam- $\pi$
- Library based: ZeroMQ

Substantial reduction in time spent going from modeling to prototype.

Do measurement at the same points that does simulation.

No (explicit) priority queues. Communication is done directly between interacting entities.

Communicate instead of dictating events.

## DISCRETE VS. REAL-TIME

Simulated durations are calculated in the processes that does the work.

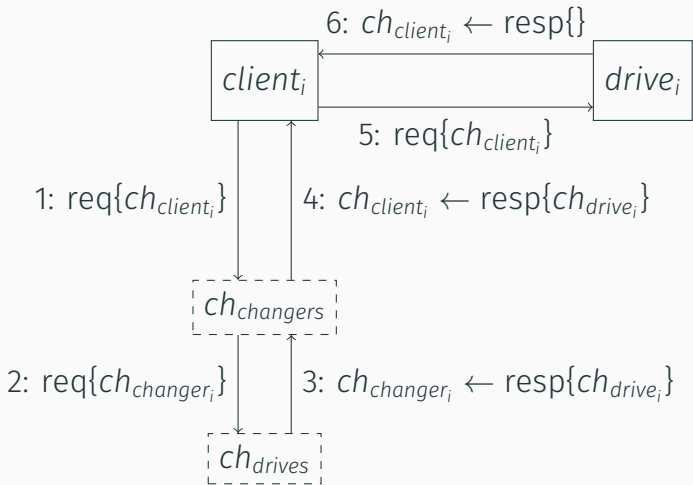
Interchangeability allows components to be swapped around and possibly mixing discrete time for some components with real-time for other components.



## SIMULATING A RATHER HUGE TAPE LIBRARY

- 90 days of constant I/O
- Three types of entities: clients, tape drives and changers
- Fixed ratio of 16 : 8 : 1
- Up to 250,000 processes simulated.

# I/O COMMUNICATION PATH



# GO

Open source concurrent programming language, created and primarily developed by Google.

Designed to be highly productive and easy to learn.

Follows the *principle of least surprise*.

Key features:

- CSP and  $\pi$ -calculus style channels and processes as low level language features.
- Garbage-collected
- Compiled
- Statically typed

```
func client(lib *library) {
    ch := make(chan response, *chanBufSize)

    for {
        lib.changers <- request{mount, ch, clock}

        resp = <-ch
        clock = clock.Add(resp.t)

        waitTime += resp.t

        t += resp.t

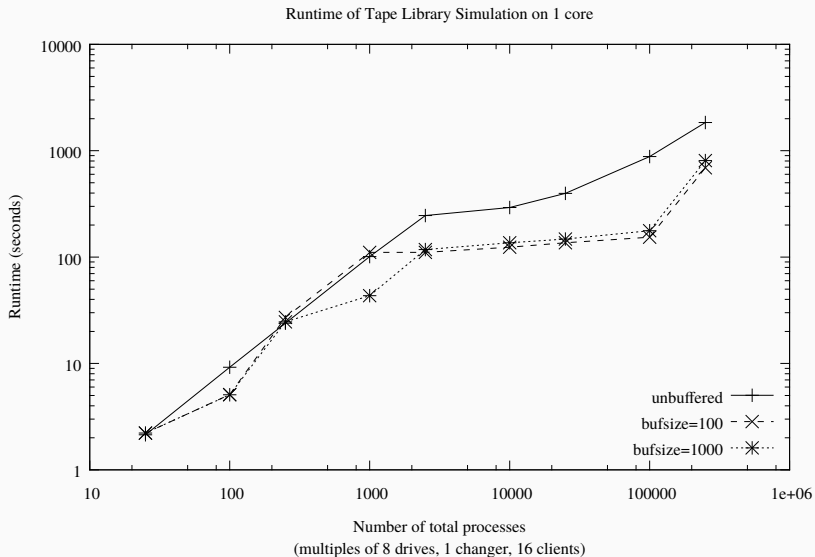
        resp.ch <- request{read, ch, clock}

        resp = <-ch
        clock = clock.Add(resp.t)
        t += resp.t
        ioTime += resp.t
    }
}
```

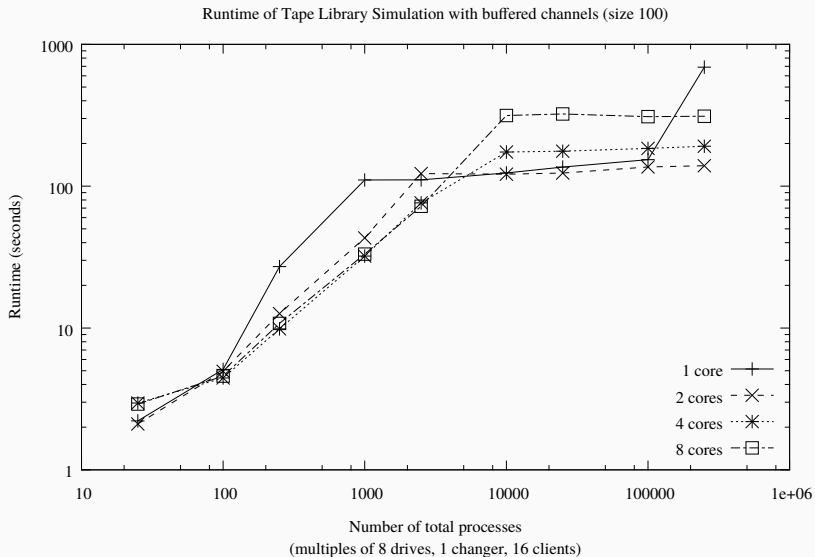
## SCALABILITY RESULTS

---

# RESULTS (SEQUENTIAL)



# RESULTS (PARALLEL)



## UNBUFFERED CHANNELS

Processes	1 core	2 cores	4 cores	8 cores
25	<b>2.14</b>	4.14	4.04	4.00
100	9.22	<b>4.96</b>	5.82	6.15
250	23.90	10.77	<b>10.71</b>	13.39
1000	101.09	37.75	<b>32.00</b>	37.77
2500	245.64	80.37	<b>70.10</b>	75.45
10000	292.40	365.42	585.33	<b>243.83</b>
25000	<b>397.34</b>	419.40	652.45	528.45
100000	881.00	<b>726.77</b>	902.13	1788.21
250000	1839.43	<b>1307.85</b>	1392.19	3671.10



## BUFFERED CHANNELS

Processes	1 core	2 cores	4 cores	8 cores
25	2.22	<b>2.11</b>	2.96	2.91
100	5.11	4.96	<b>4.44</b>	4.58
250	27.09	12.63	<b>9.86</b>	10.78
1000	110.72	43.19	<b>32.16</b>	33.19
2500	110.83	122.83	76.30	<b>72.19</b>
10000	123.91	<b>121.59</b>	174.01	315.17
25000	136.47	<b>123.72</b>	176.50	322.98
100000	153.77	<b>136.59</b>	184.69	309.25
250000	691.15	<b>139.34</b>	191.30	311.50

## SUMMARY AND FUTURE WORK

- Rapid transition from simulation/modeling to prototype
- Communicate instead of dictating events
- No reverse computation
- Scales well with at least Go
  
- Further refinement and packaging of the ISI patterns.
- Look into locality management of Goroutines.

THANK YOU

QUESTIONS?