# Storage as a First Class Citizen in HPC Environments.

James S. Plank
University of Tennessee

CCGSC
September 9, 2010

# A Personal Historical Perspective

Me – Erasure codes

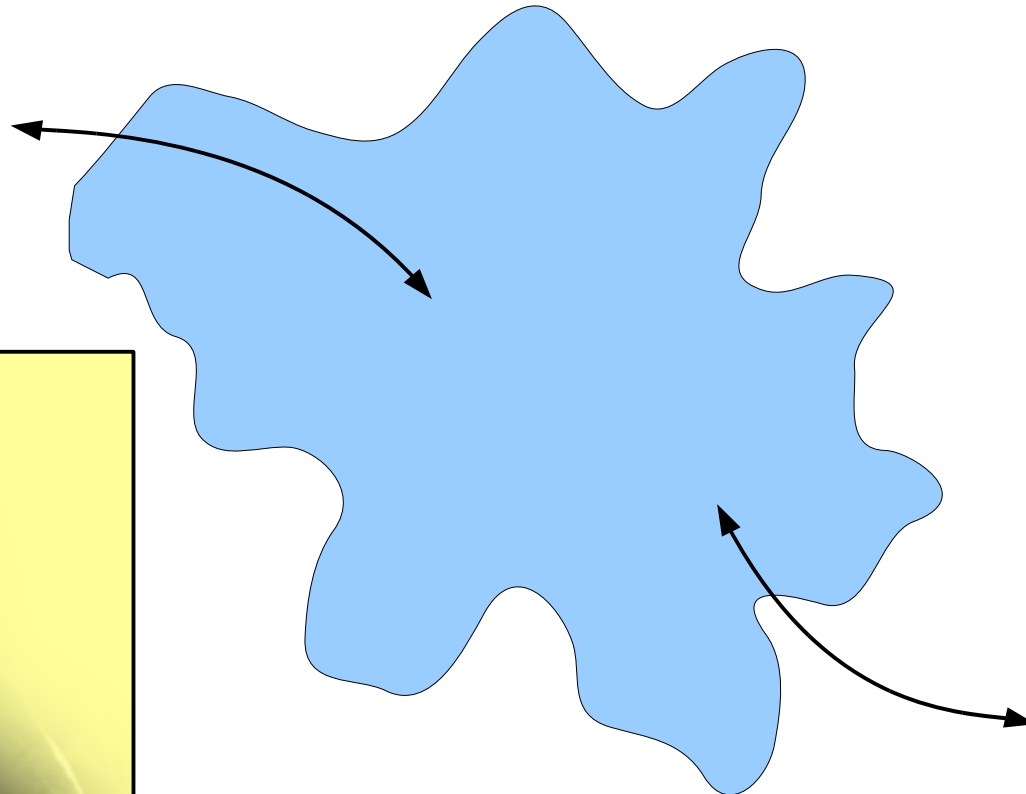Y'all –
HPC

# A Personal Historical Perspective



Jim - 1987

# A Personal Historical Perspective

LINDA: Parallel computing with a "tuple space."



Data tuples

Jim - 1987

Gelernter

Processing tuples

**Tuple Space**

# A Personal Historical Perspective

LINDA: Parallel computing with a "tuple space."

- "Linda processes aspire to know as little about each other as possible.

Gelernter

Jim - 1987

- They never interact directly with each other;

- they only deal with tuple space."

# A Personal Historical Perspective
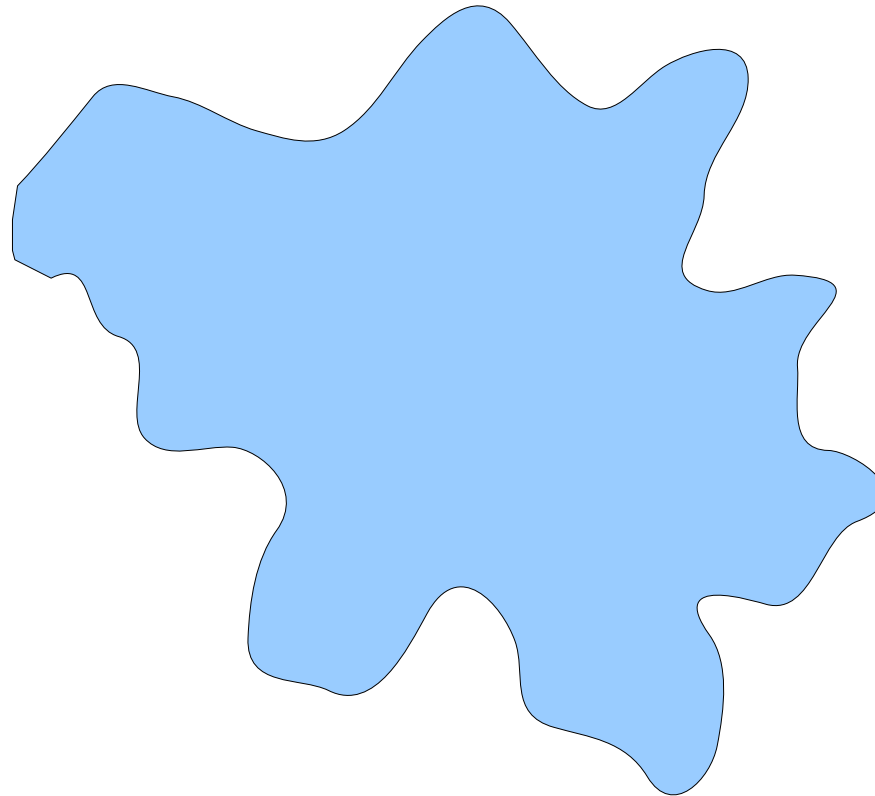


Jim - 1988

# A Personal Historical Perspective

SSLS: Shared Single Level Store



Jim - 1988

Naughton

Gigantic shared, persistent address space

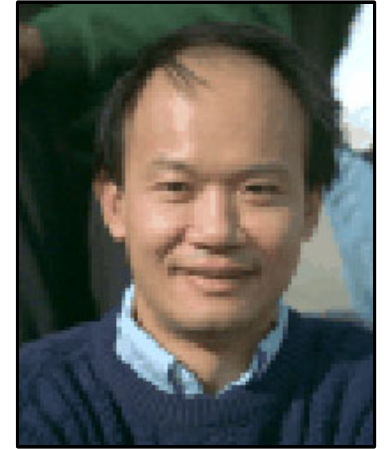# A Personal Historical Perspective
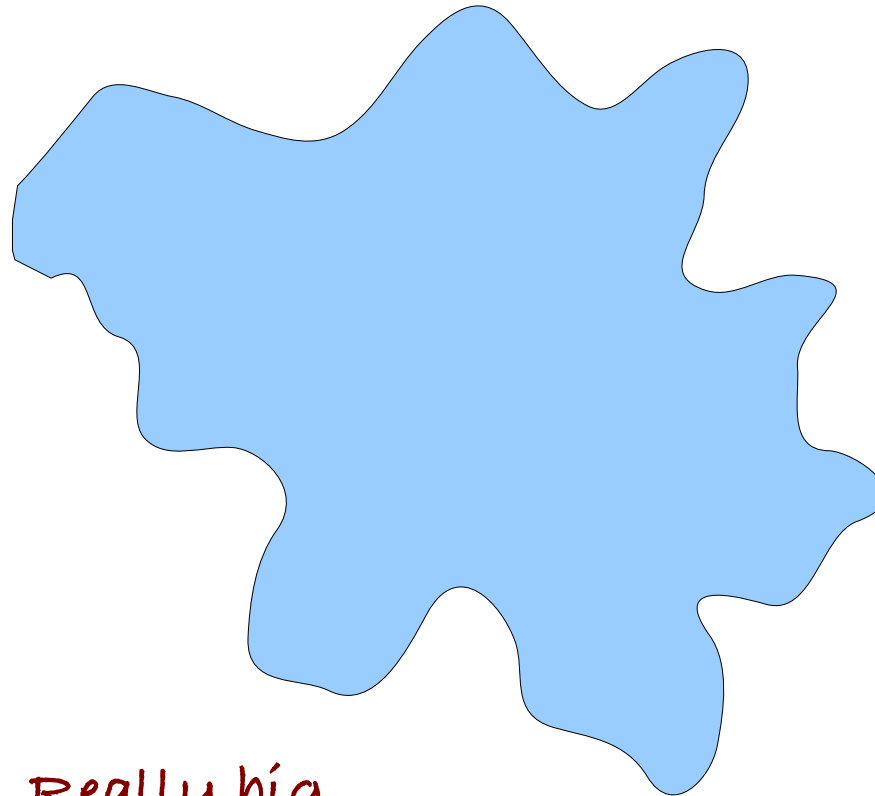
SSLS – Shared Single Level Store

Jim - 1988

Naughton

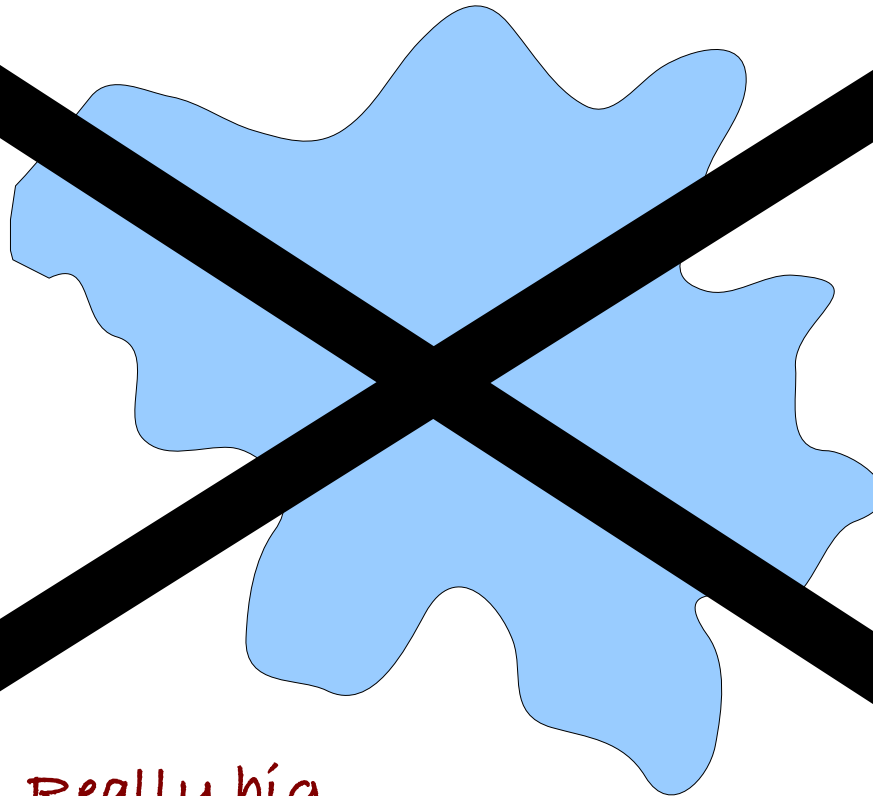Gigantic shared, persistent address space

# A Personal Historical Perspective
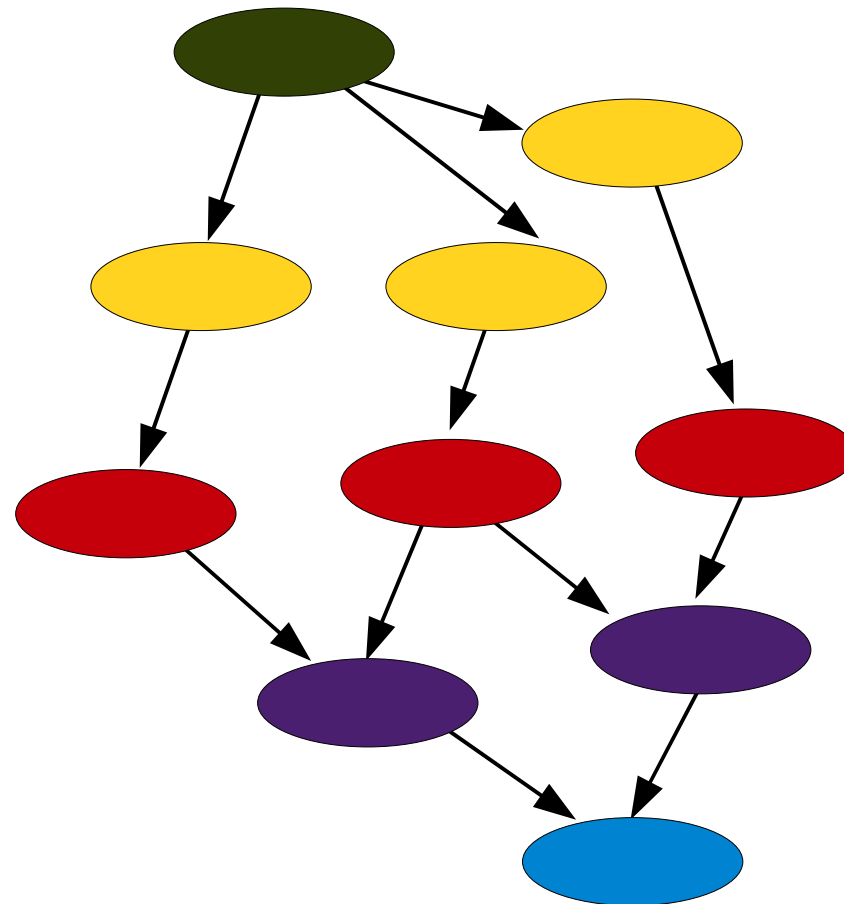
SVM: Shared Virtual Memory

Li

Jim - 1989

~~Really big~~

~~Gigantic~~ shared, ~~persistent~~ address space

# A Personal Historical Perspective

# A Personal Historical Perspective

HeNCE: Heterogeneous Network Computing Environment.

Jim - 1990

Grand
Fromage

Functional
Dataflow
DAG Processing
System

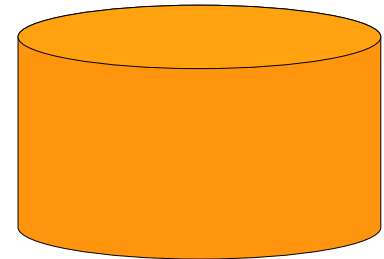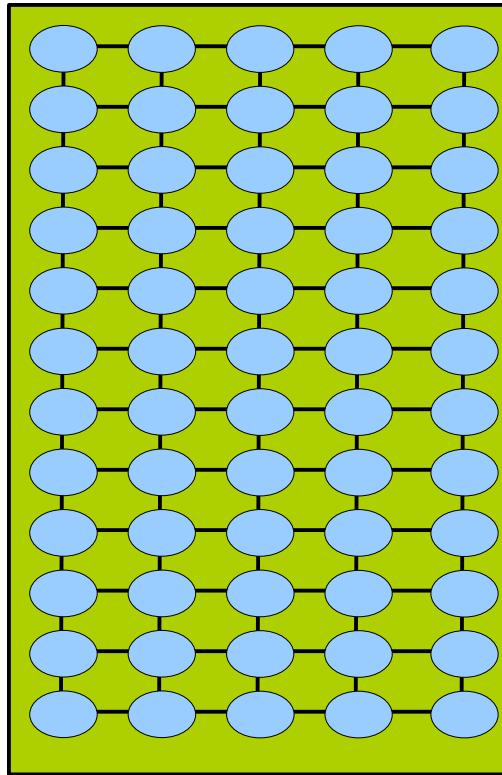# A Personal Historical Perspective



Jim - 1991-98

# A Personal Historical Perspective

Mr. Checkpointing:

Jim - 1991-98

# A Personal Historical Perspective

There are two major difficulties with checkpointing:

## What did I learn:

Jim - 1991-98

1. Fighting the OS / Getting it to work.

2. Mitigating the overhead of getting
   all those bytes to disk.

Everything else (synchronization,
consistency, Lamport time, etc, etc)
is in the noise.
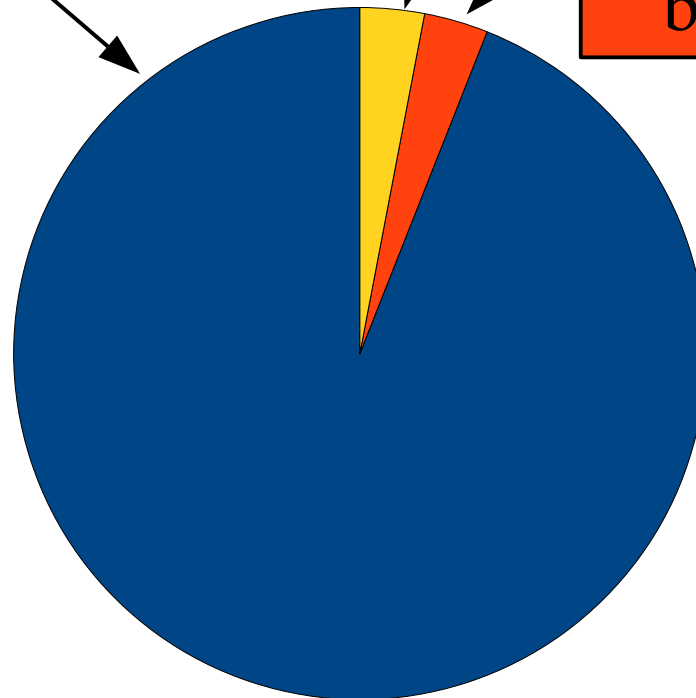
# A Personal Historical Perspective

## Where's the research?

Getting it to work.

Synchronization, consistency, Lamport time, etc, etc.

Mitigating the overhead of getting all those bytes to disk.
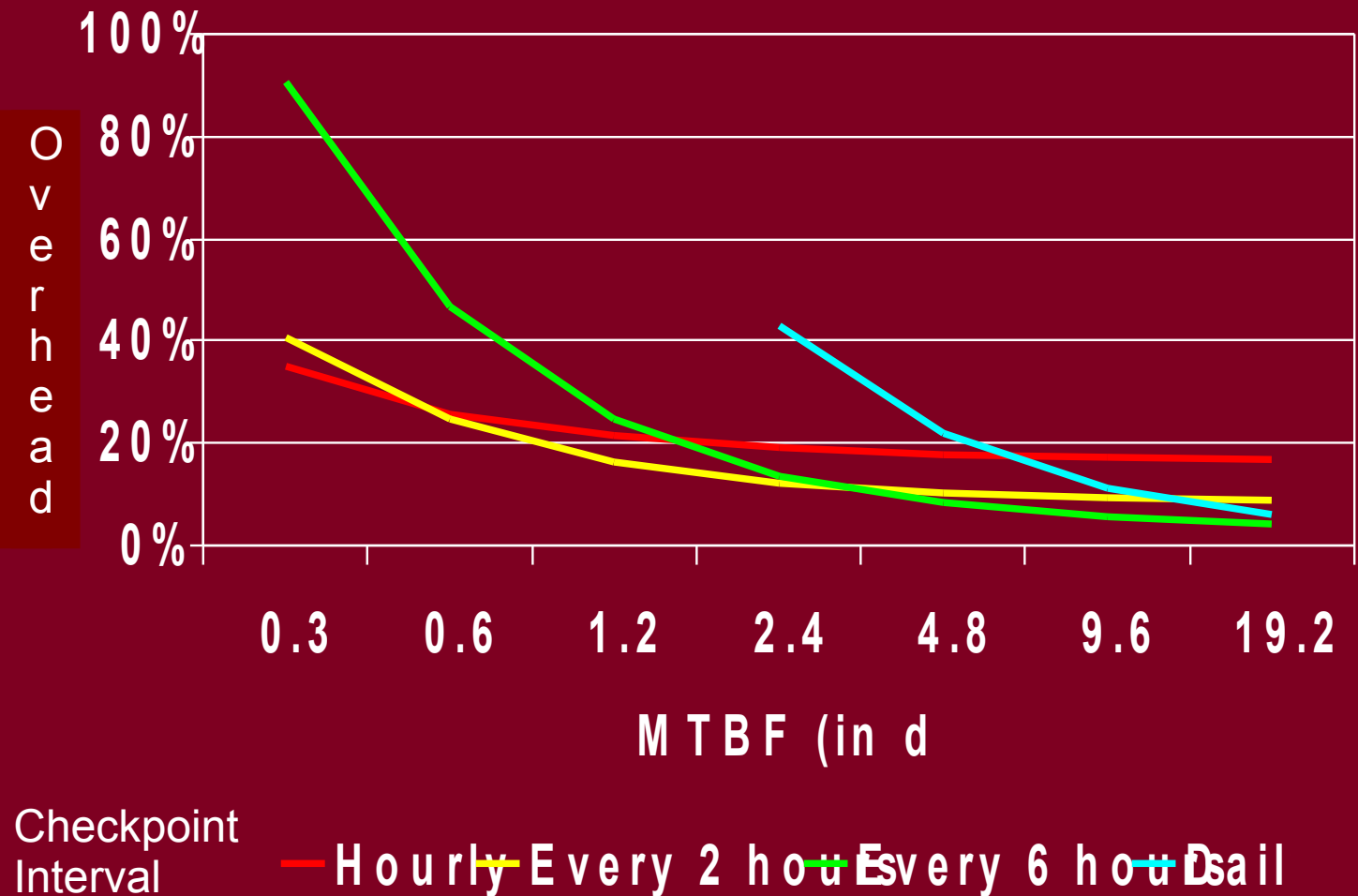
Jim - 1991-98

# A Personal Historical Perspective

[Elnozahy/Plank 2004]

Jim - 1991-98

**Cost of Reliability (at 1**

Overhead

100%
80%
60%
40%
20%
0%

0.3    0.6    1.2    2.4    4.8    9.6    19.2

**MTBF (in d**

Checkpoint
Interval

—Hourly —Every 2 hours —Every 6 hours —Dail
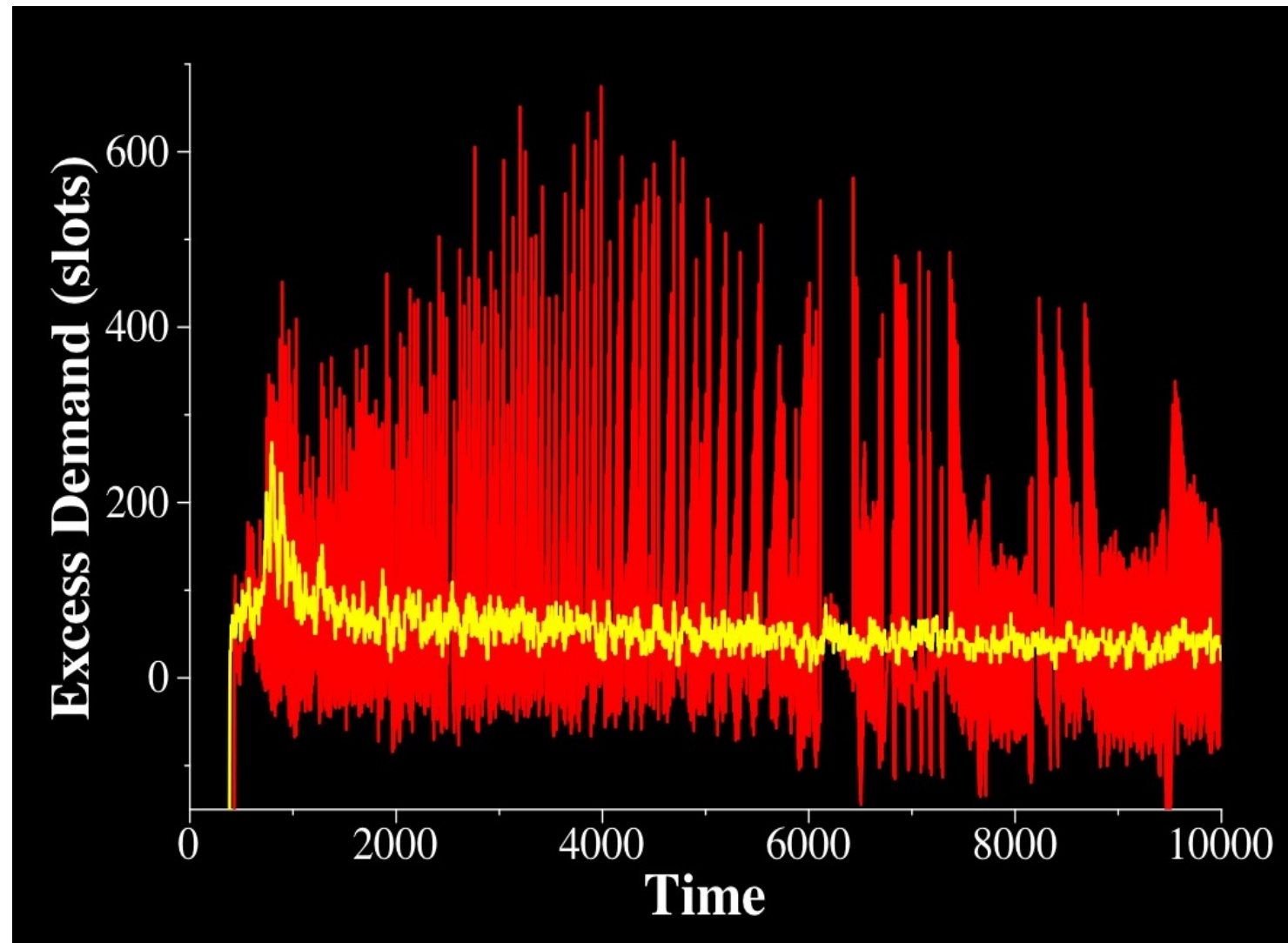
# A Personal Historical Perspective



Jim – 1999

# A Personal Historical Perspective

## G-Commerce: Brief Foray into Grid Computing

Jim – 1999

# A Personal Historical Perspective

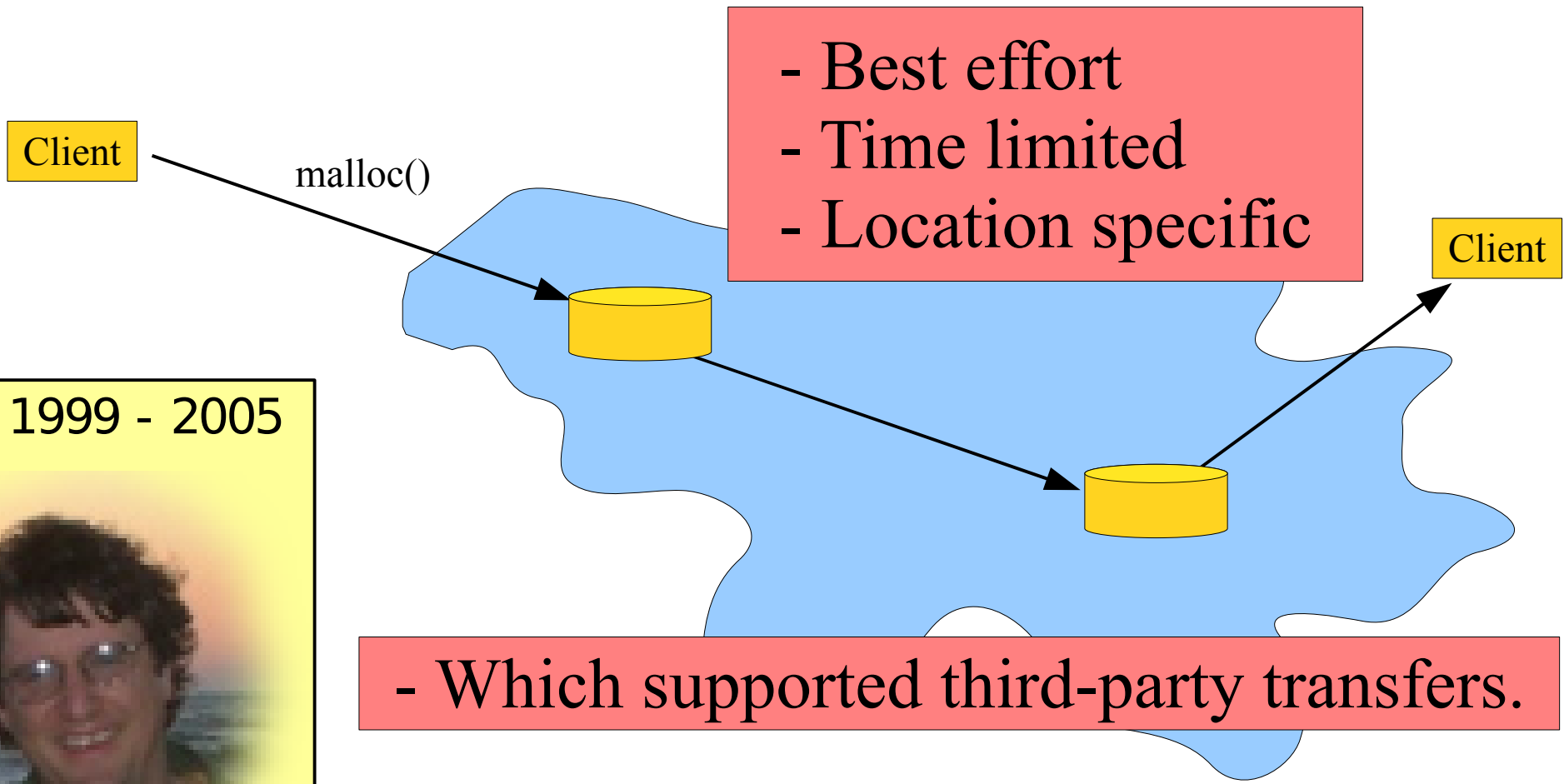## IBP: Internet Backplane Protocol (Logistical Networking)
w/ Micah Beck

Client

malloc()

Jim – 1999 - 2005

# A Personal Historical Perspective

## IBP: Internet Backplane Protocol (Logistical Networking)
w/ Micah Beck

Client
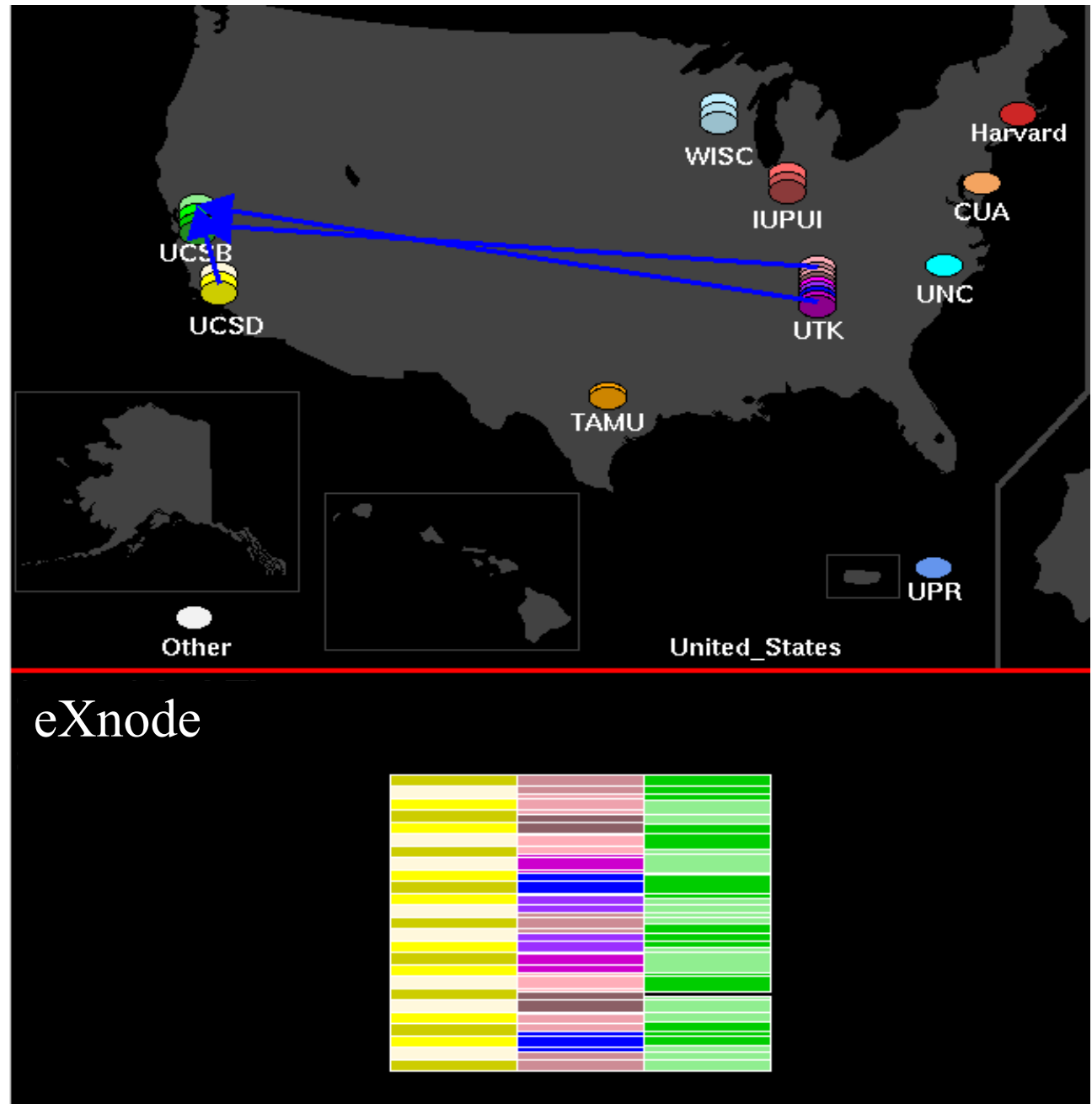
malloc()

- Best effort
- Time limited
- Location specific

Client

Jim – 1999 - 2005

- Which supported third-party transfers.

# A Personal Historical Perspective

IBP gave data a place to "live" on the network, perhaps moving from site to site.

Jim – 1999 - 2005

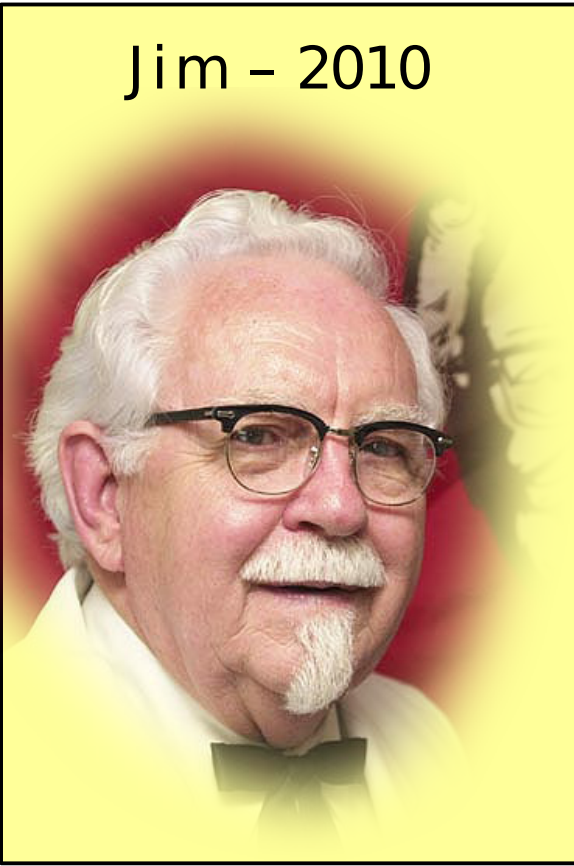# A Personal Historical Perspective

Into the land of erasure coding.

Jim – 2005 - ???
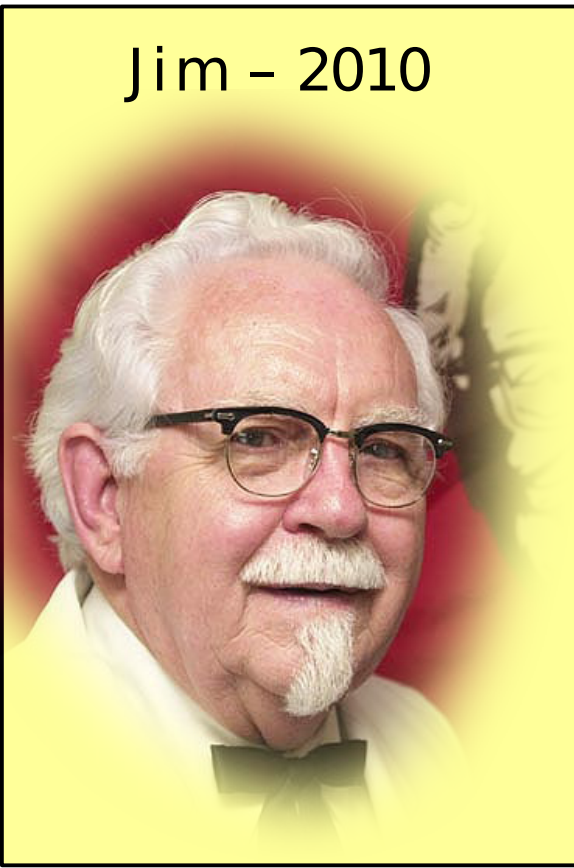
I won't bore you with it.
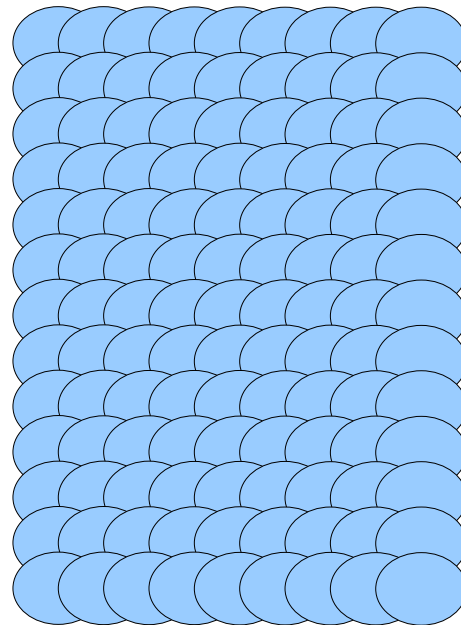
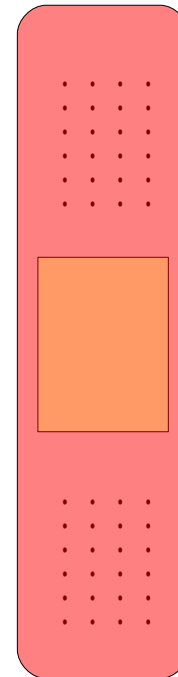But there's more...



Jim – 2010

# A Personal Historical Perspective

2010 Meeting on
Staging for HPC

Jim – 2010

The Big Iron

"Staging"

The Disks

# A Personal Historical Perspective
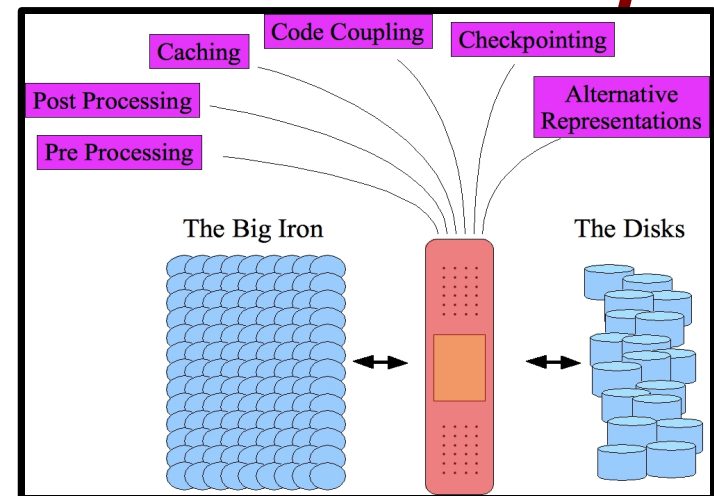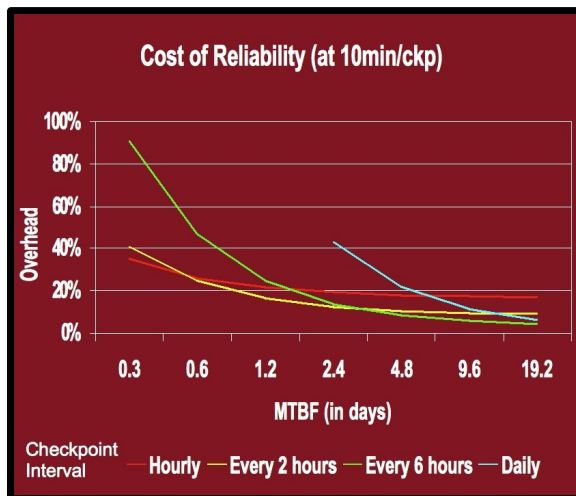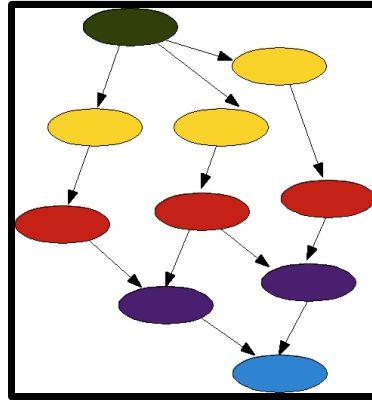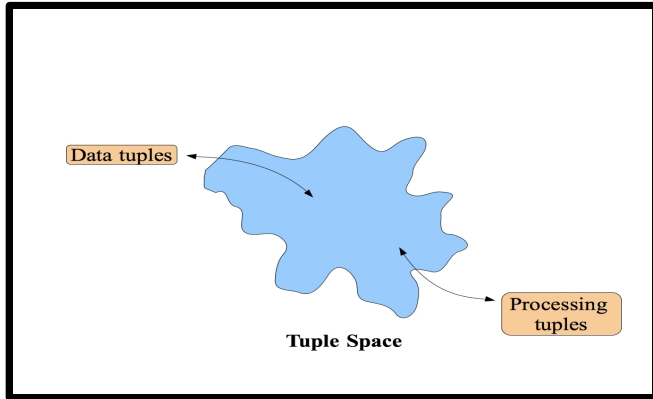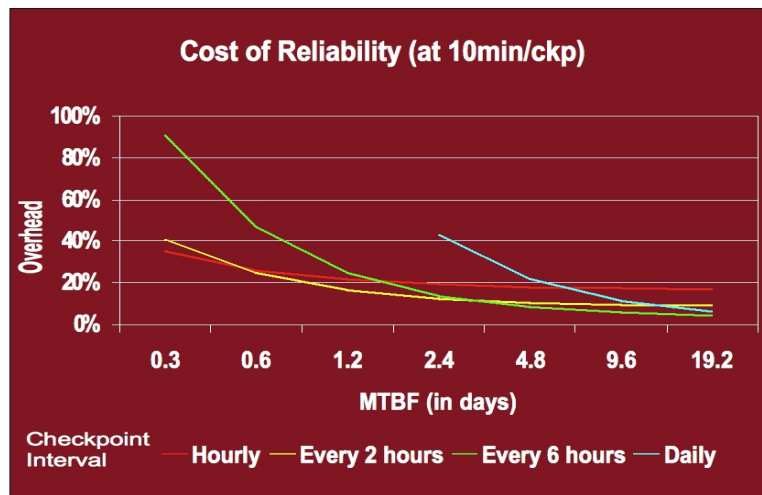
# What do we make of all this?

# What do we make of all this?
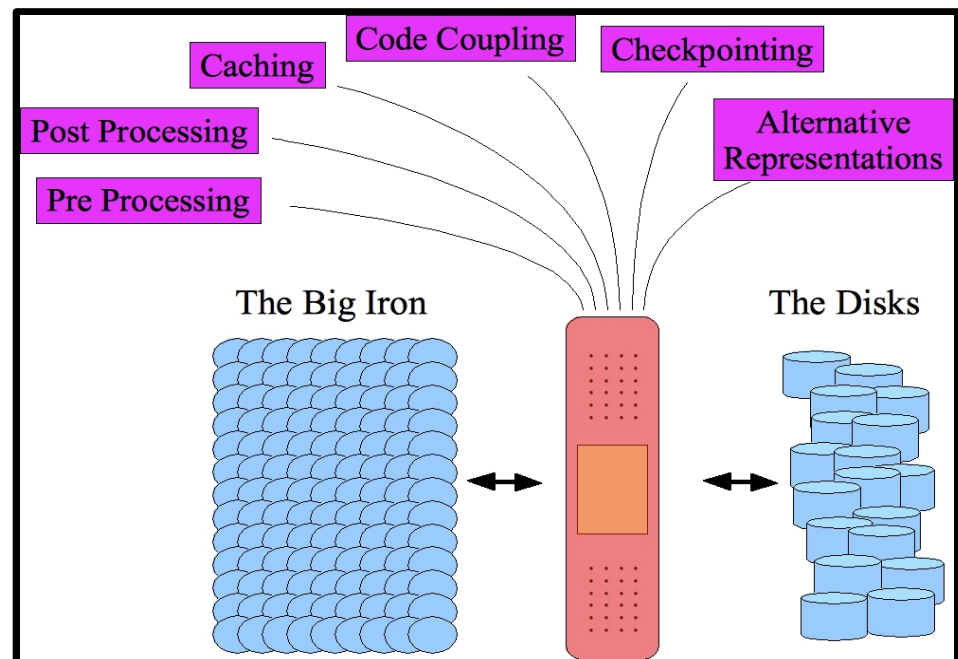
## 1. Checkpointing Sucks.

- Slow
- Inelegant
- Swamps disks and networks to store gigantic files that are almost never read.
- Enables you to perform "bad fault-tolerance."
- Is a manifestation that something is wrong.
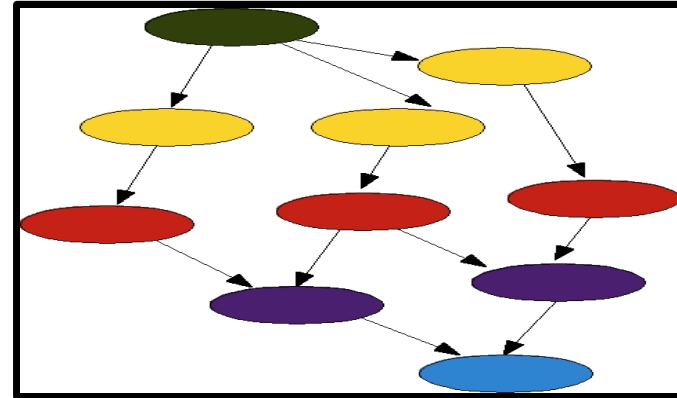


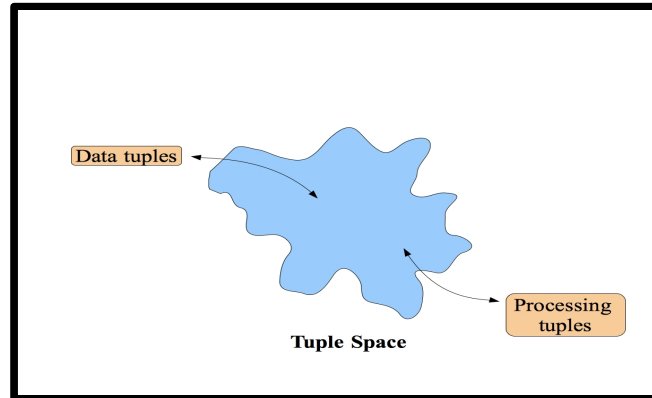Cost of Reliability (at 10min/ckp)

# What do we make of all this?

## 2. Band-Aids Are Only Temporary Solutions

- Non-reusable
- Cover the wounds but don't address the root cause
- Are a manifestation that something is wrong.
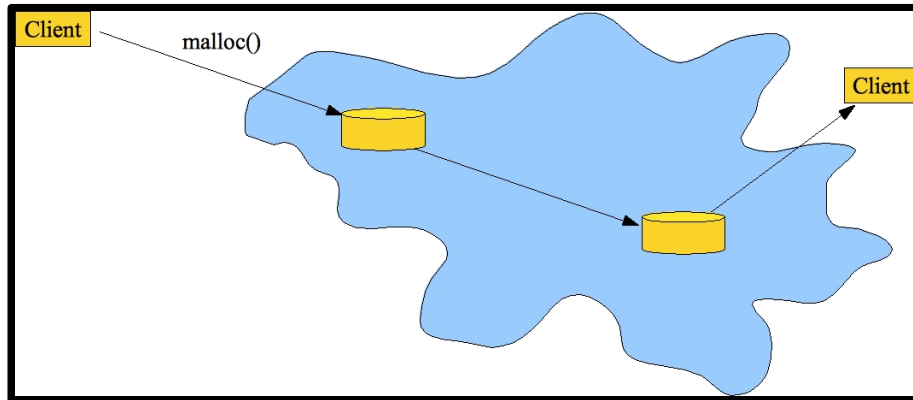
# What do we make of all this?



## 3. Saving State Sure is Attractive

- Lets you reason about programs
- (In theory) lets balance load
- Allows fault tolerance to fall out naturally

- However, it's really difficult to do.
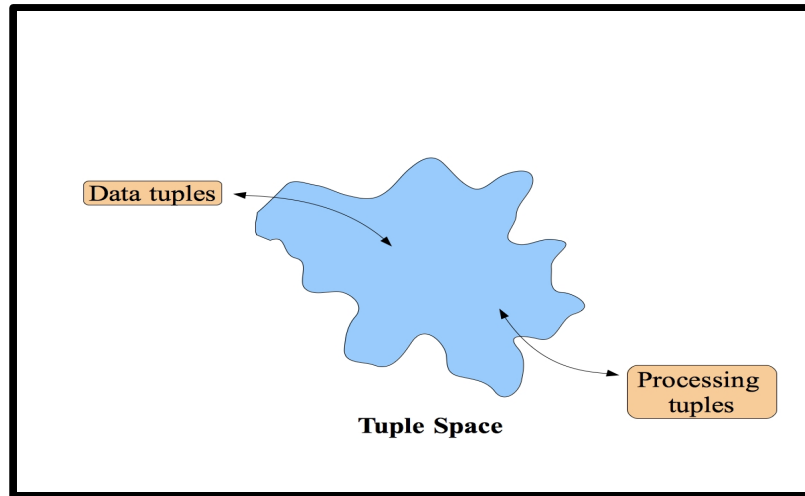- This is why the MPI model throws it in the trash can.
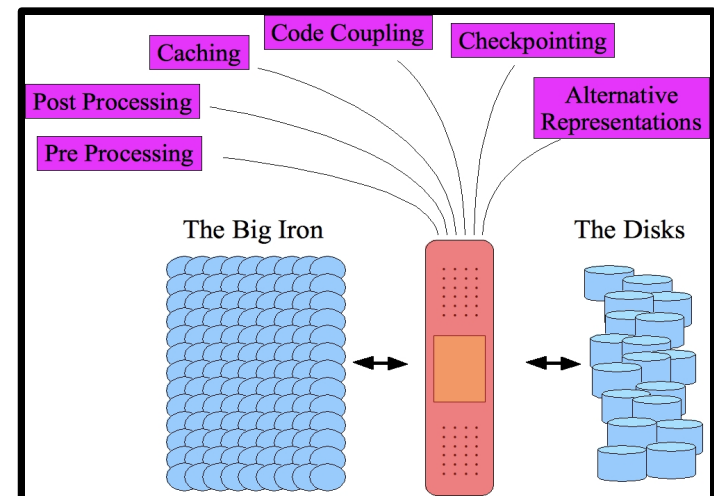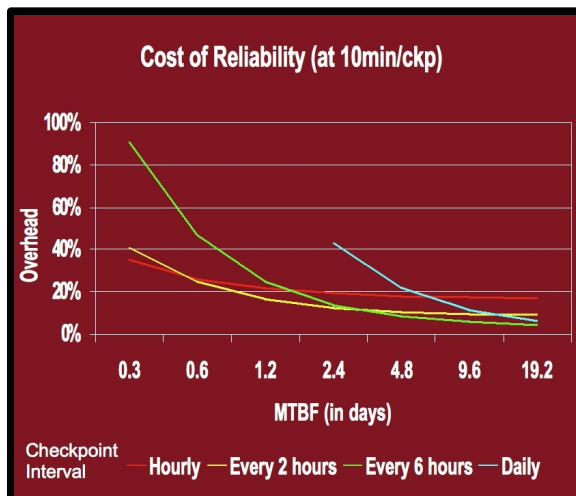
# What do we make of all this?



## 4. I Still Think IBP is Pretty Cool & That There Are Lessons To Be Learned From It

- Why do we constrain our view of storage as either the file or the memory segment?
- Why is storage either permanent or limited by program lifetime?
- Why do we jettison best-effort storage resources?
- Why don't we manage the location of storage?

# What do we make of all this?



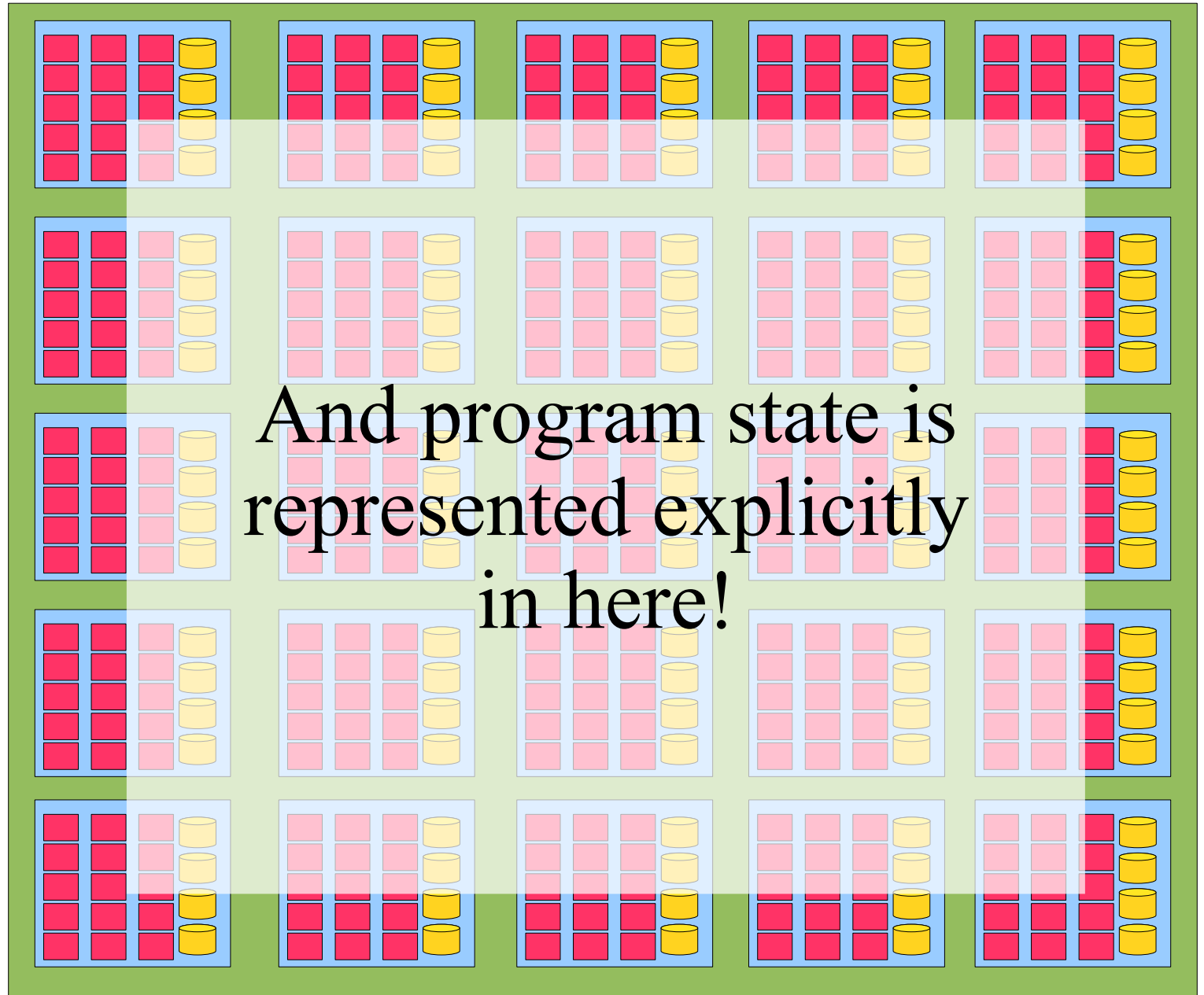## Why are storage and processing not equal first-class citizens in HPC?

# When I Close My Eyes and Dream ……

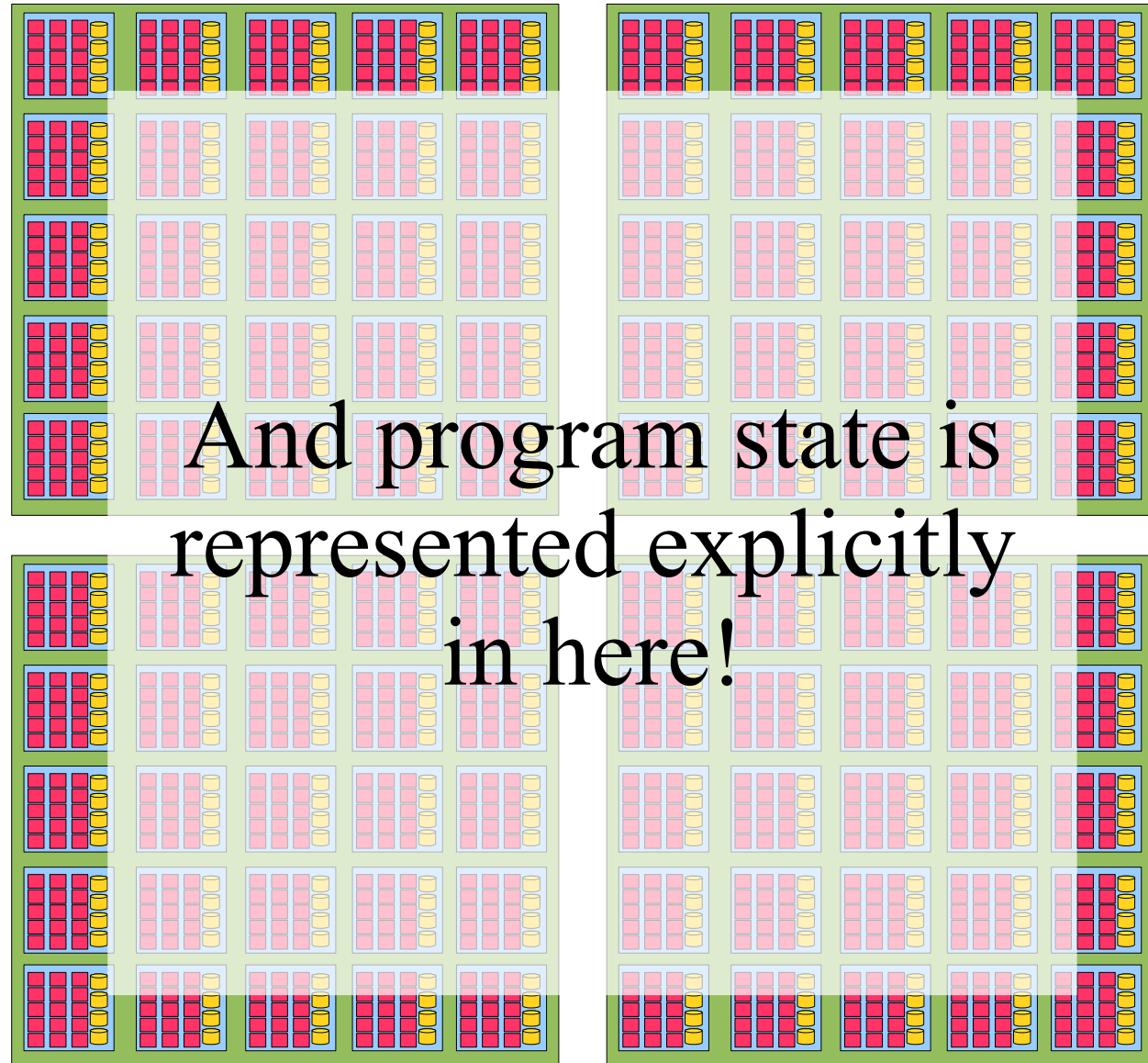The Big Iron looks like this.

And these guys:

are promoted to first class citizens.

And program state is represented explicitly in here!

And these guys compose seamlessly.

Over extremely wide areas ....

And program state is represented explicitly in here!

# When I Close My Eyes and Dream ......

And the Eagles win the Super Bowl...      Every Year...



And I retire to that mansion in Capri...

# Storage as a First Class Citizen in HPC Environments.

James S. Plank
University of Tennessee

CCGSC
September 9, 2010