# SUMMIT

## JBoss WORLD

## PRESENTED BY RED HAT

# LEARN. NETWORK.
# EXPERIENCE OPEN SOURCE.

www.theredhatsummit.com

# Simplifying Parallel Programming
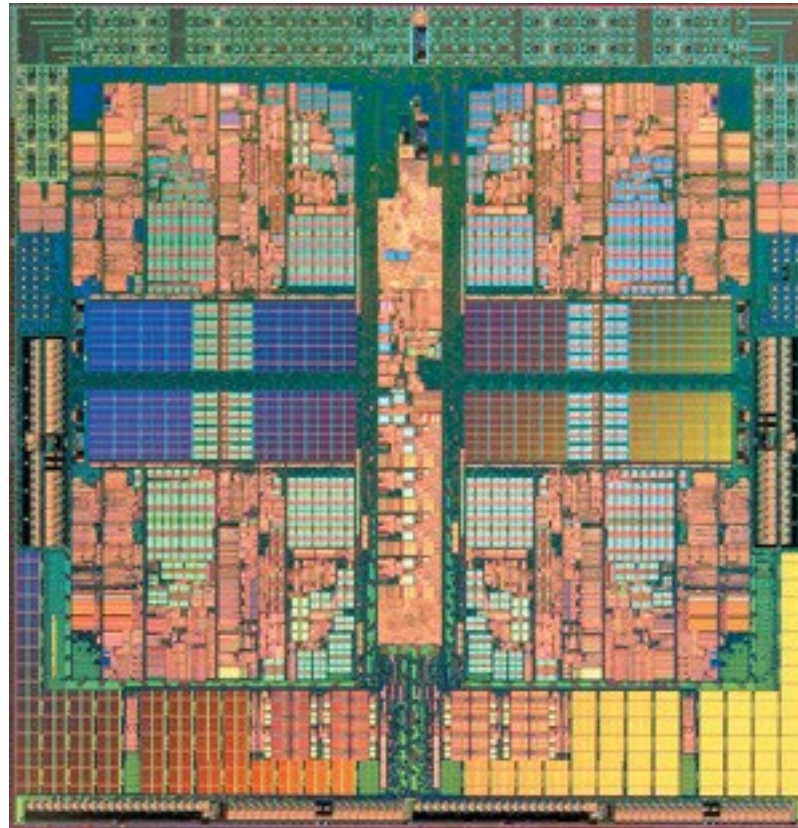
Ulrich Drepper
Consulting Engineer, Red Hat
2010-6-25

# The Problem

# The Problem

# The Problem

# The Problem

# The Problem

# The Reason

$$E = C \times V^2 \times f$$

# More Correctly

$$E = C \times V(f)^2 \times f$$

# Use of Transistors

- Increasing frequency is out

- Two uses

  - More complex architecture

    - Handle existing instructions faster

    - More specialized instructions

  - Horizontal growth

    - More execution cores; or          **Requires Parallelism!**

    - Only more execution contexts

# Cost of Too Little Parallelism

- Idealized Amdahl's Law

$$S \; = \; \cfrac{1}{(1-P) \; + \; \cfrac{P}{N}}$$

- Problems
  - $P$ too small
  - $N$ is steadily growing
- Formula is unrealistic though…

# A More Realistic Formula

- Extended Amdahl's Law with Overhead

$$S = \frac{1}{(1-P)\ (1+O_S)\ +\ \dfrac{P}{N}\ (1+O_P)}$$

- Parallelization is not free
  - Most of the time not even for serial code
- The results are not *that* bad…

# Even with Overhead P=0.6



Legend:
- 0%
- 20%
- 40%
- 90%
- 1000%

- Even with 40% overhead not that much slower

- Speed-up from two threads on

  - Eleven threads for 10x slowdown

# Programming Goals

$$S = \frac{1}{(1-P)\ (1+O_S)\ +\ \dfrac{P}{N}\ (1+O_P)}$$

- Two goals: 1. ease parallel programming to increase $P$

  2. reduce $O_S$ and $O_P$

# Getting Parallelism

- Multi-process Pipeline

| | | |
|---|---|---|
| **Process 1** | **Process 2** | **Process 3** |

**Unix Pipeline**

# Problems with Pipelines

- Marshalling needed for transmission

- Protocol standardization required

- Limited buffer sizes

  - Lots of scheduling needed

- Program need to be designed for pipeline

  - Extending an existing program not easy

  - Major code restructuring needed

# Problems with Pipelines

- Marshalling needed for transmission
- Protocol standardization required
- Limited buffer sizes
  - Lots of scheduling needed
- Program need to be designed for pipeline
  - Extending an existing program not easy
  - **Major code restructuring needed**
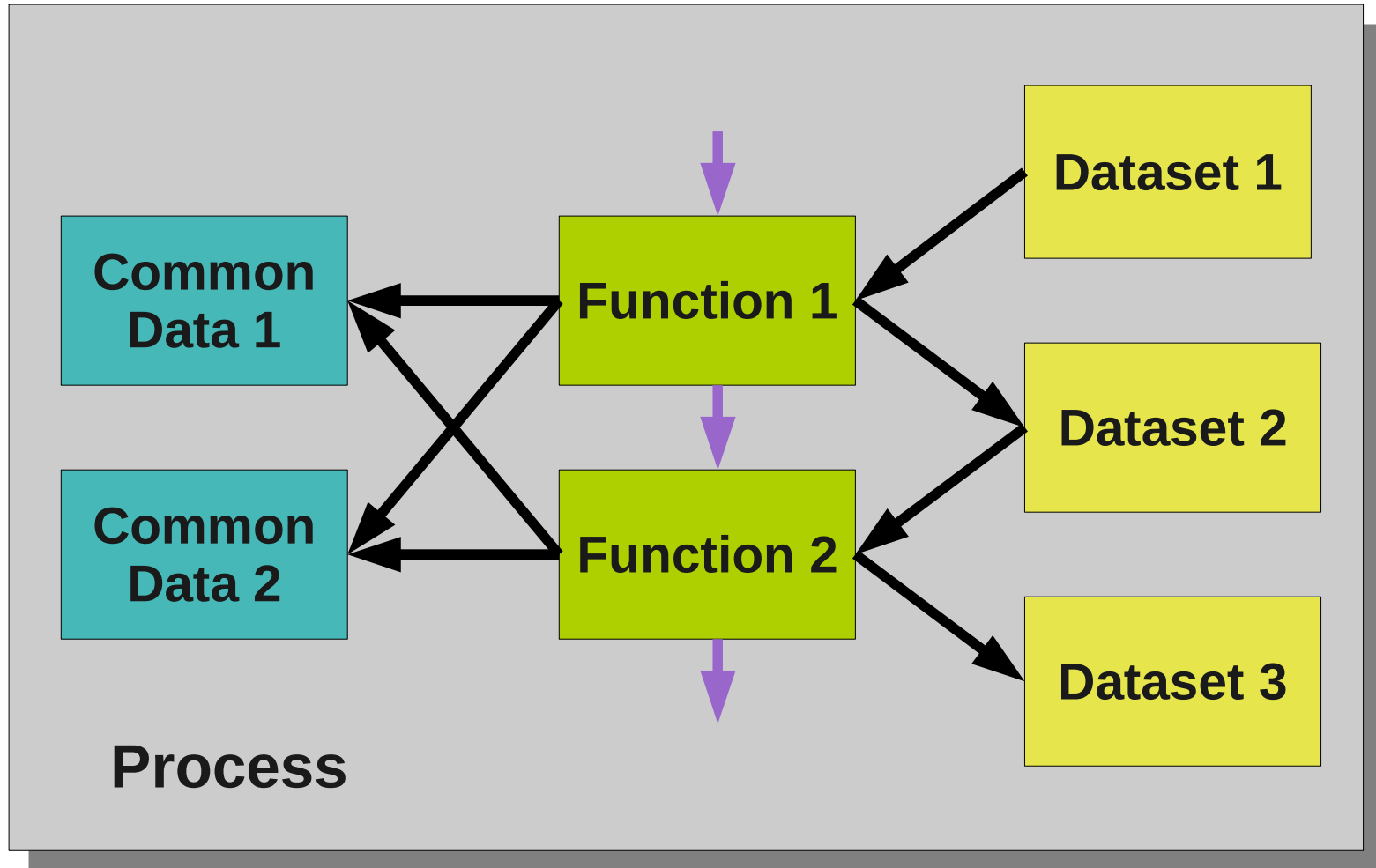
# Simple Program Structure

# "Easy" Fix

# "Easy" Fix



Process

# It seems easy…

# Explicit Multi-Threading

- Ill-conceived solution
  - Yes
    - Existing code can be reused, easier to set up
    - High-bandwidth inter-thread communication
    - On some OSes context switching faster
  - But:
    - Fragile programming model (one thread dies, the process dies)
    - Memory handling mistakes have global effects
    - Unix model initially not designed for multiple threads

**Hard to write correct code! High Cost!**

# Measures



| | | | |
|---|---|---|---|
| Common Data 1 | Mutex | Thread 1 | Dataset 1 |
| Common Data 2 | | Thread 2 | Mutex Dataset 2 |
| Process | | | Dataset 3 |

| Reuse | Fragile |
|---|---|
| Bandwidth | Overwrites |
| Context Cost | Unix model |
| Ease Program | Error Prone |

# Alternative 1: fork and Shared Memory

- All in POSIX:

```
int fd = shm_open(name, O_RDWR|O_CREAT);
ftruncate(fd, size);
p = mmap(NULL, size, PROT_READ|PROT_WRITE,
         MAP_SHARED, fd, 0);
if (fork() == 0)

    ...
```

# fork and Shared Memory



| | |
|---|---|
| Reuse | Fragile |
| Bandwidth | Overwrites |
| Context Cost | Unix model |
| Ease Program | Error Prone |

State Data

Process 1

Dataset 1

Mutex

Mutex

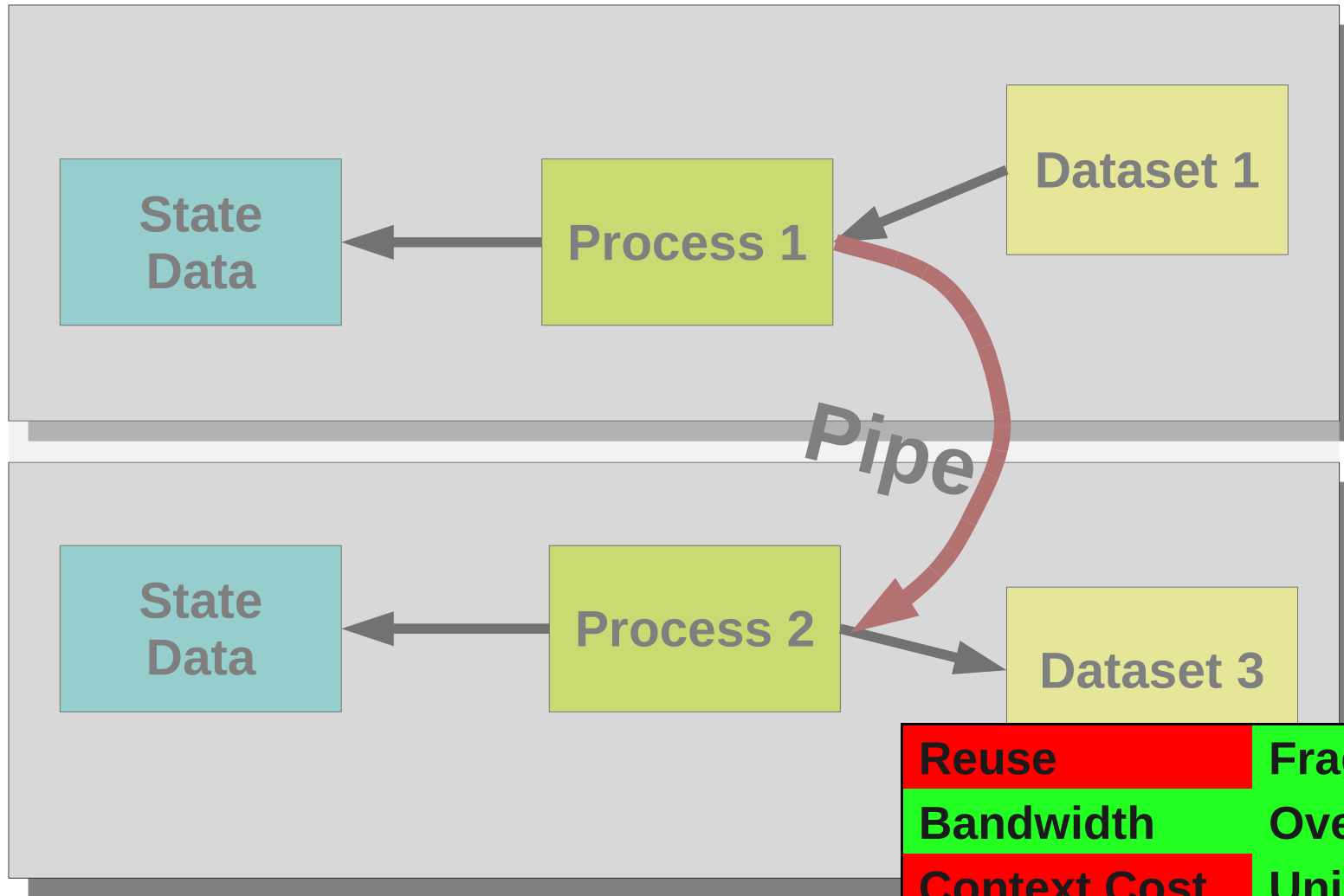Dataset 2

State Data

Process 2

Dataset 3

# Alternative 2: `fork` and Linux Pipes

- Linux extensions, not POSIX (yet ☺ )

- Can be zero-copy

- Use if just transferring data without inspection

- splice: transfer from file descriptor to pipe

- tee: transfer between pipes and keep data usable

- vmsplice: transfer from memory to pipe

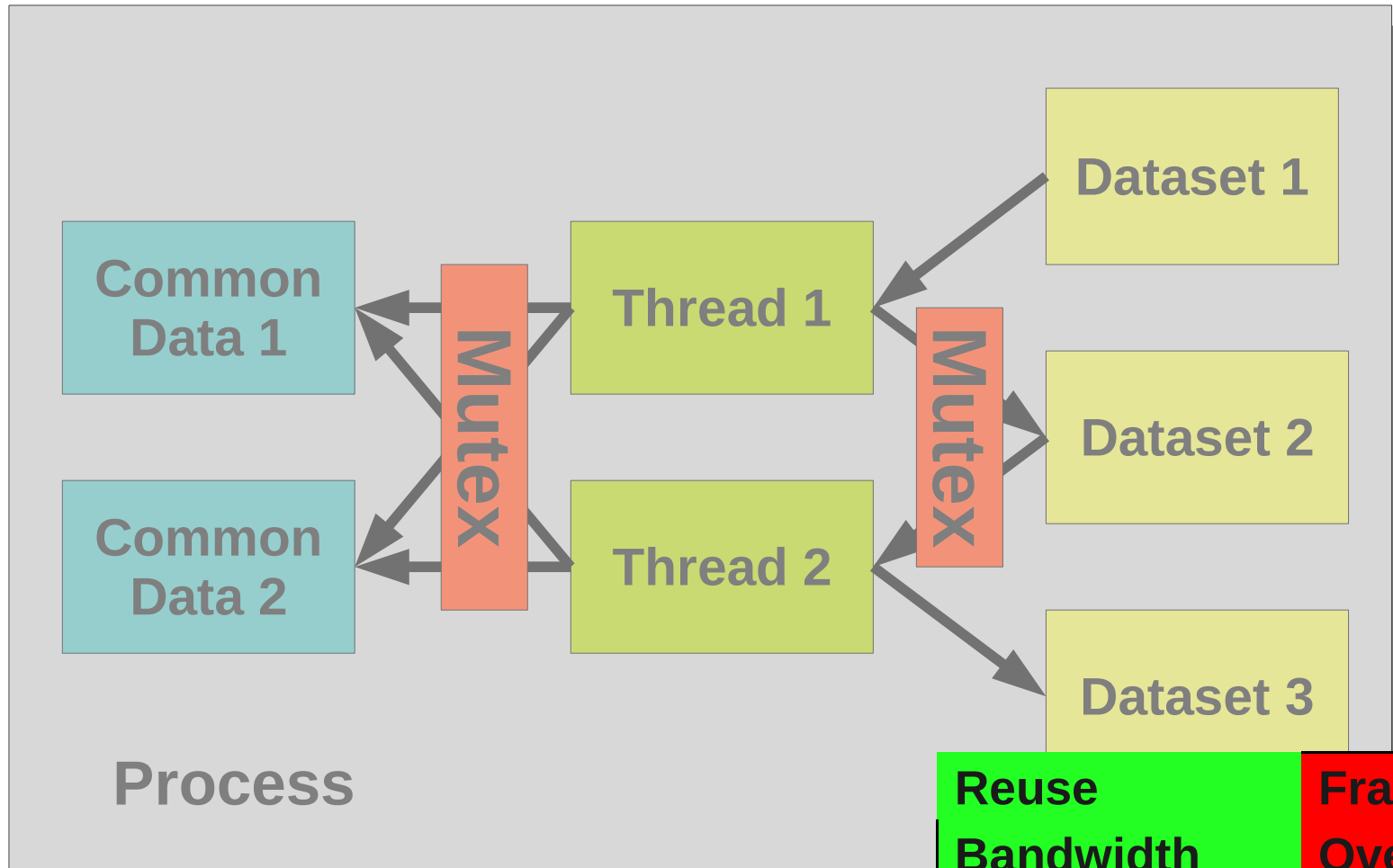# fork and Linux Pipes

# Alternative 3: Thread Local Storage

- Use thread-local storage
    - Very much simplifies use of static variables
    - No more false sharing of cache lines

```
__thread struct foo var;
```

# Thread Local Storage



Process

Common Data 1

Common Data 2

Mutex

Thread 1

Thread 2

Mutex

Dataset 1

Dataset 2

Dataset 3

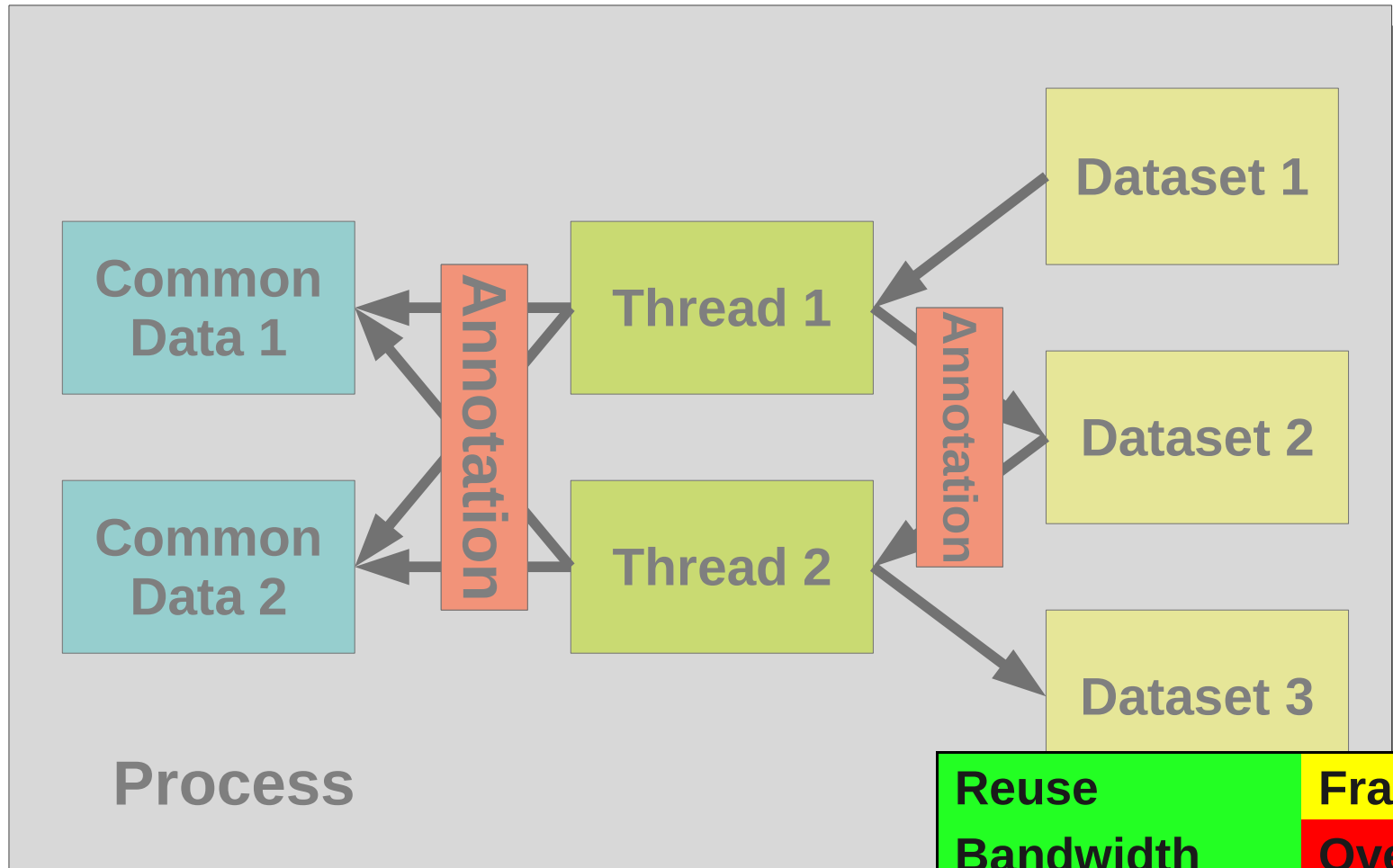| Reuse | Fragile |
|---|---|
| Bandwidth | Overwrites |
| Context Cost | Unix model |
| Ease Program | Error Prone |

# Alternative 4: OpenMP

- Language extension to C, C++, Fortran languages
- Implements many thread functions with very simple interface for
  - Thread creation (controlled)
  - Exclusion
  - Thread-local Data

# OpenMP



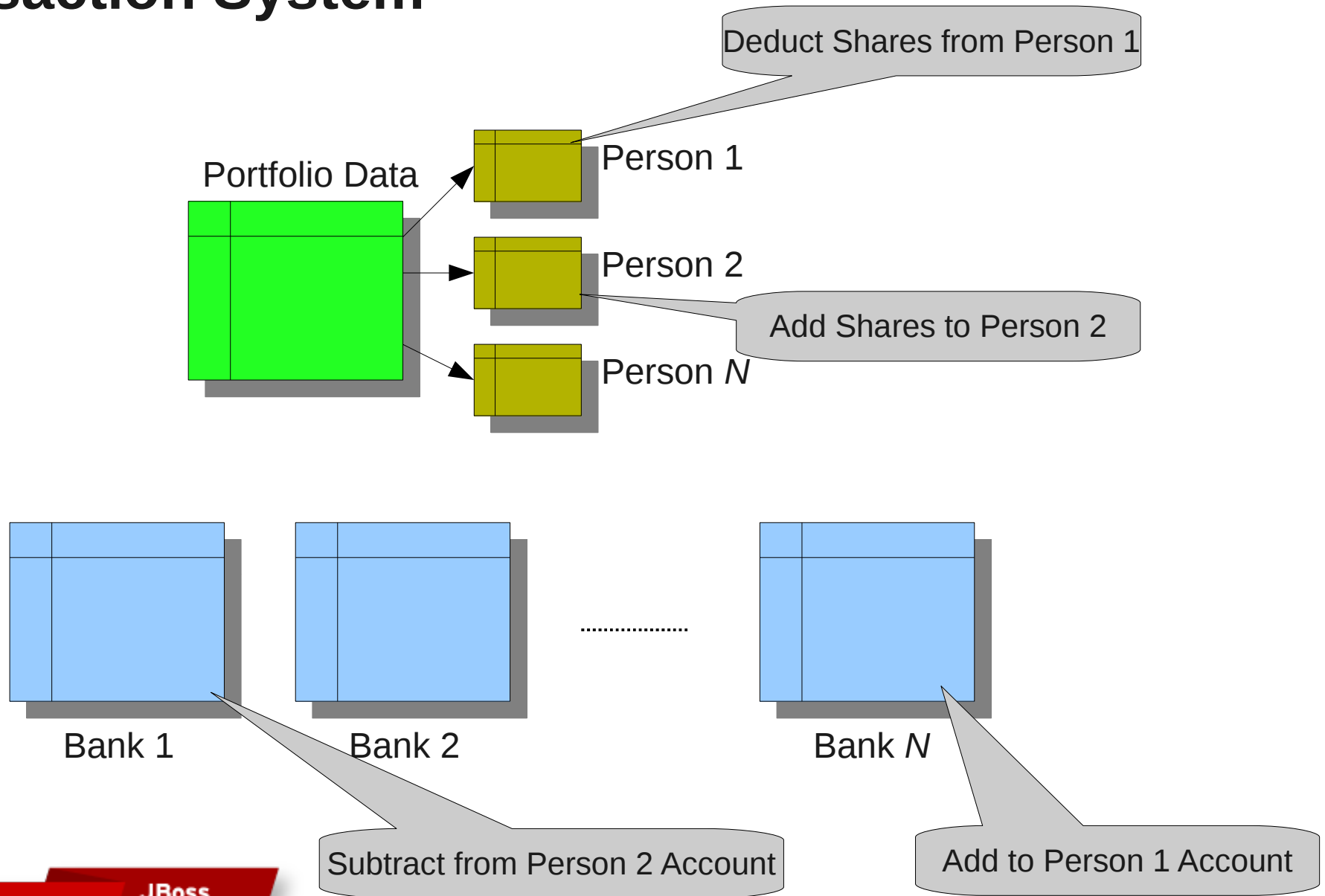| | | | |
|---|---|---|---|
| Reuse | Fragile | | |
| Bandwidth | Overwrites | | |
| Context Cost | Unix model | | |
| Ease Program | Error Prone | | |

# Alternative 5: Transactional Memory

- Extensions to C and C++ languages
- Can help to avoid using mutexes
    - Just source code annotations
    - No more deadlocks!!
    - Fine-grained locking without the problems
- Slow as pure software solutions
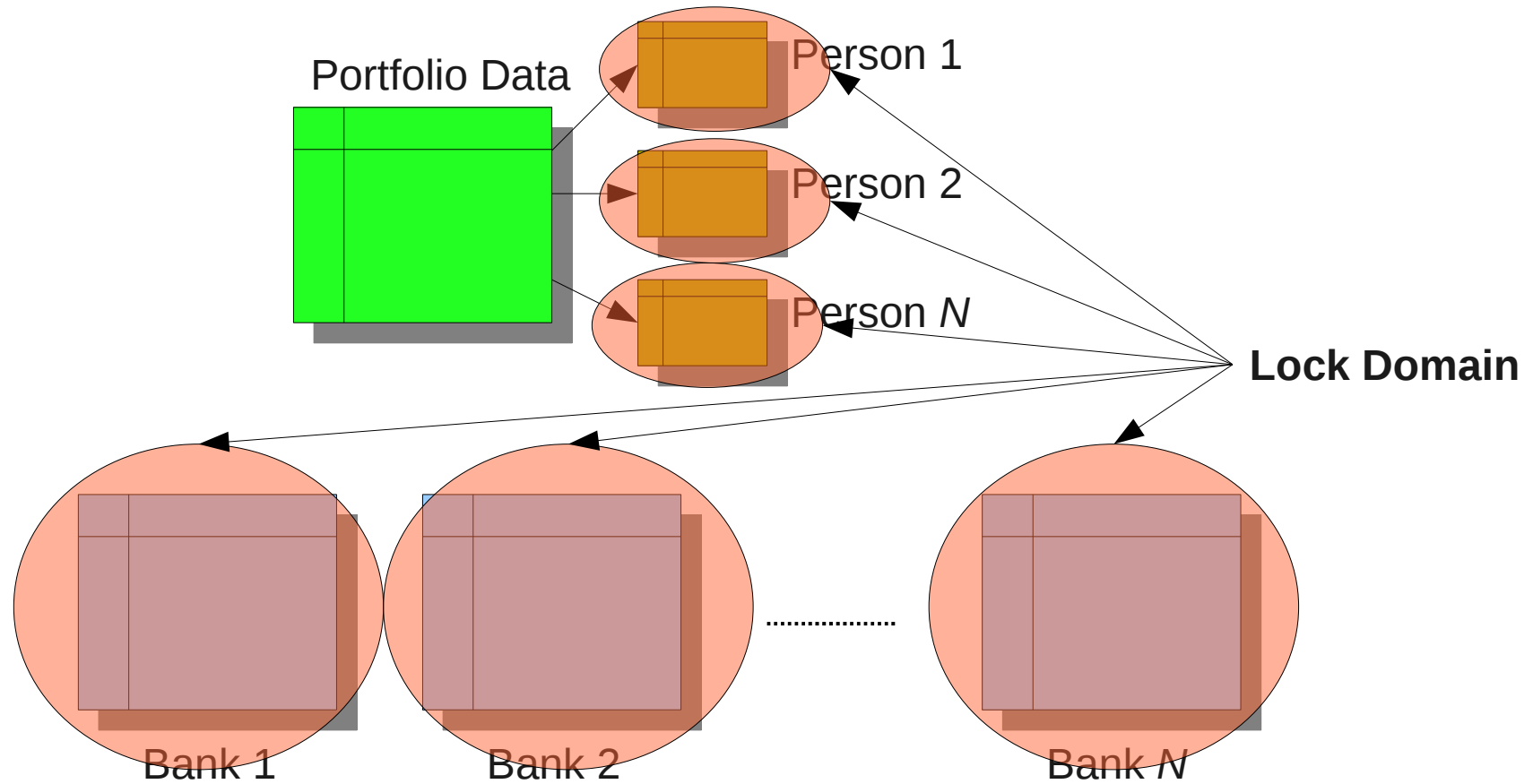    - Hardware support on the horizon

# Transaction System

# Trying to Parallelize

Portfolio Data

Person 1

Person 2

Person *N*

**Lock Domain**

Bank 1

Bank 2

..................

Bank *N*

# Not What We Want



Single Core i7



Opteron NUMA

SUMIT **JBoss WORLD**

PRESENTED BY RED HAT

# Trying to Parallelize

# Somewhat Better But…



Single Core i7



Opteron NUMA

# Transactional Memory



Process

| | |
|---|---|
| **Reuse** | **Fragile** |
| **Bandwidth** | **Overwrites** |
| **Context Cost** | **Unix model** |
| **Ease Program** | **Error Prone** |

# Conclusion

- Abilities to exploit hardware are there
    - Explicit threading only for experts
- But there is a lot of help
    - Use processes, not threads; or
    - If threads are used combine
        - Thread-local storage
        - Implicit thread creation
            - OpenMP
            - Futures
        - Transactional memory

# Questions?

# FOLLOW US ON TWITTER

www.twitter.com/redhatsummit

# TWEET ABOUT IT

#summitjbw

# READ THE BLOG

http://summitblog.redhat.com/

SUMMIT JBoss WORLD

PRESENTED BY RED HAT