# MeDiCI: UQ's Metropolitan Data Caching Infrastructure

### David Abramson

Director, Research Computing Centre

Professor of Computer Science

University of Queensland

david.abramson@uq.edu.au

# ~~Turtles~~ Caches all the way down

"a jocular expression of the infinite regress problem in cosmology posed by the "unmoved mover" paradox.

The metaphor in the anecdote represents a popular notion of the theory that Earth is actually flat and is supported on the back of a World Turtle, which itself is propped up by a chain of larger and larger turtles.

Questioning what the final turtle might be standing on, the anecdote humorously concludes that it is turtles all the way down""
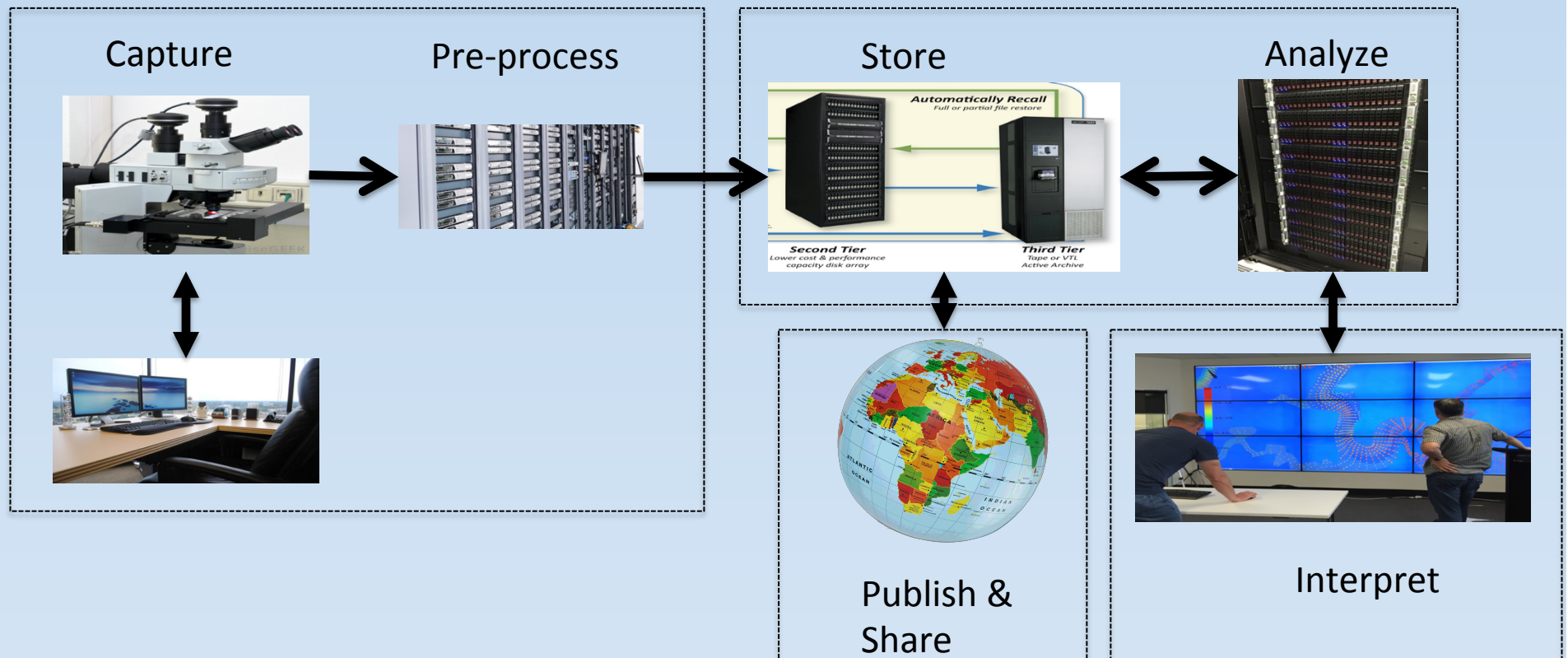


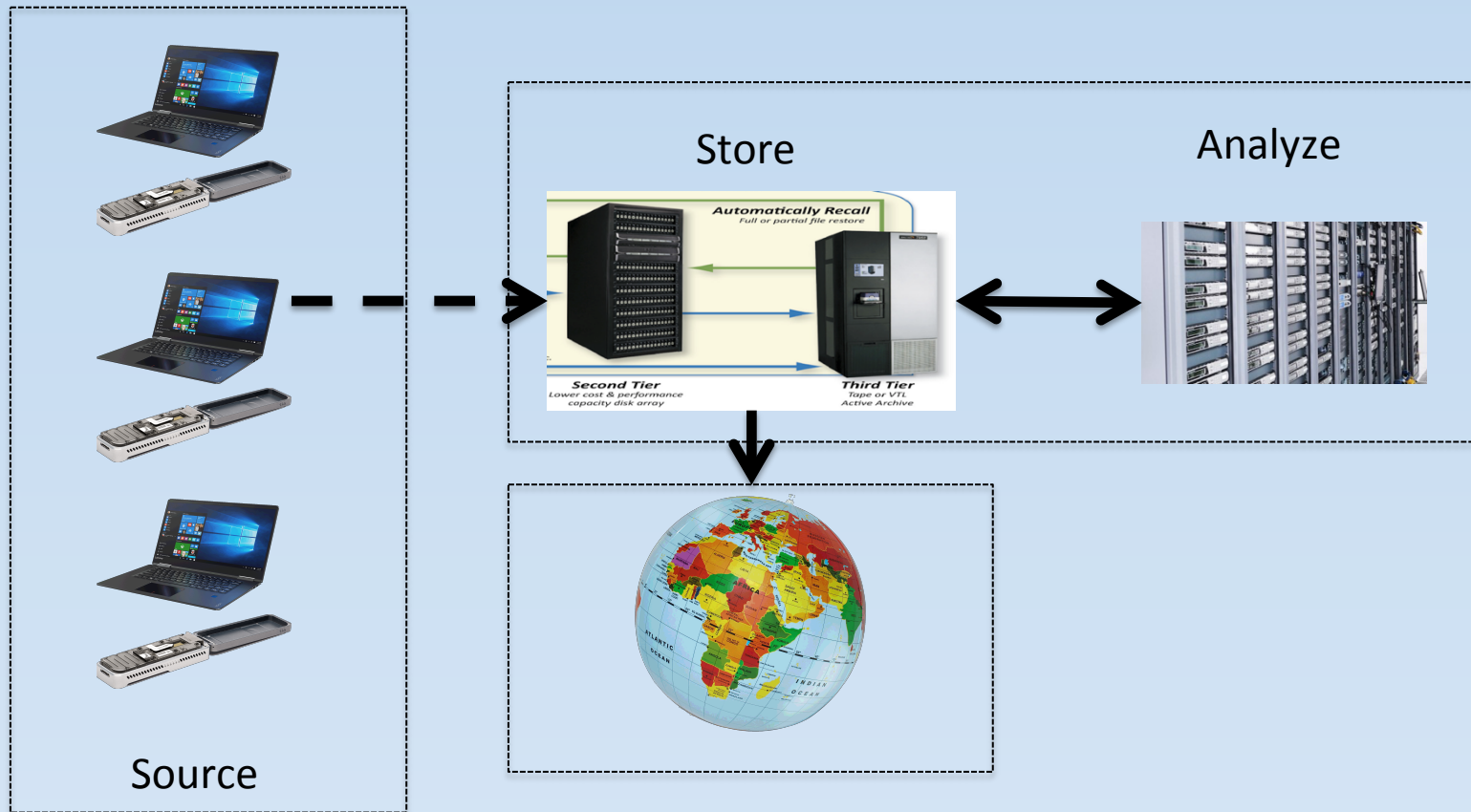https://en.m.wikipedia.org/wiki/Turtles_all_the_way_down

Why do we need to do anything special?
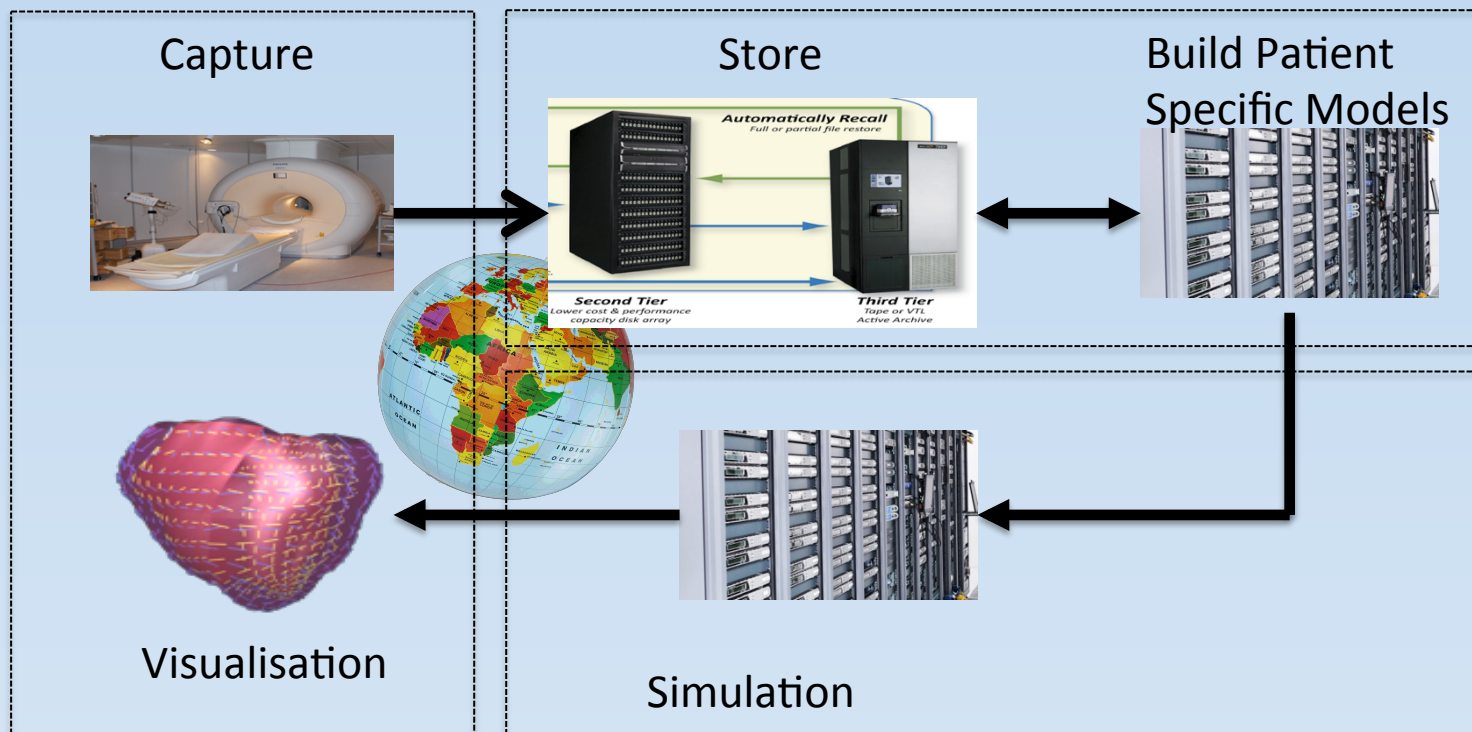
# Data Intensive Computing

# Use Case: Microscopy



Capture

Pre-process

Store

Analyze

Publish & Share

Interpret

# Use Case: Personal Genomics



Source

Store

Analyze

**Automatically Recall**
Full or partial file restore

**Second Tier**
Lower cost & performance
capacity disk array

**Third Tier**
Tape or VTL
Active Archive

# Use Case: Cardiac Science



Capture

Store

Build Patient Specific Models

Visualisation

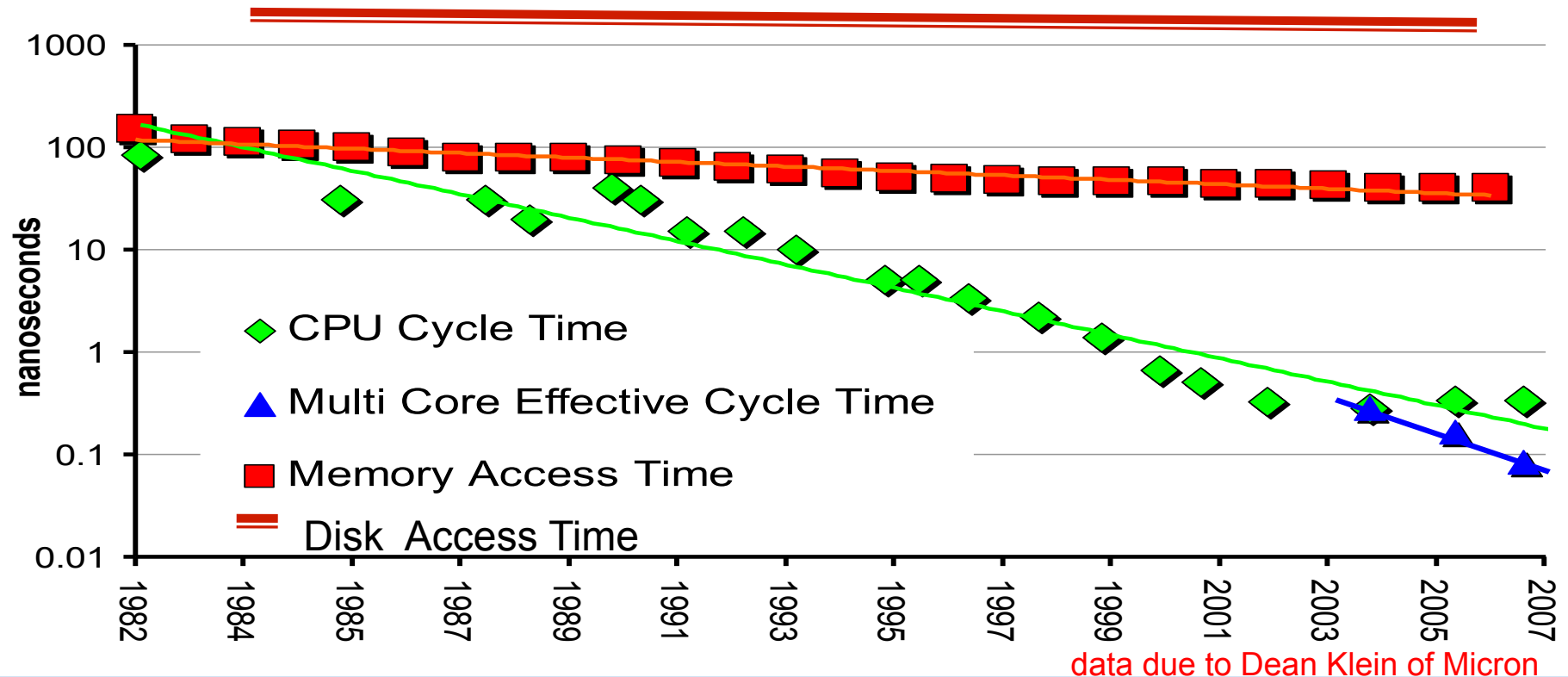Simulation

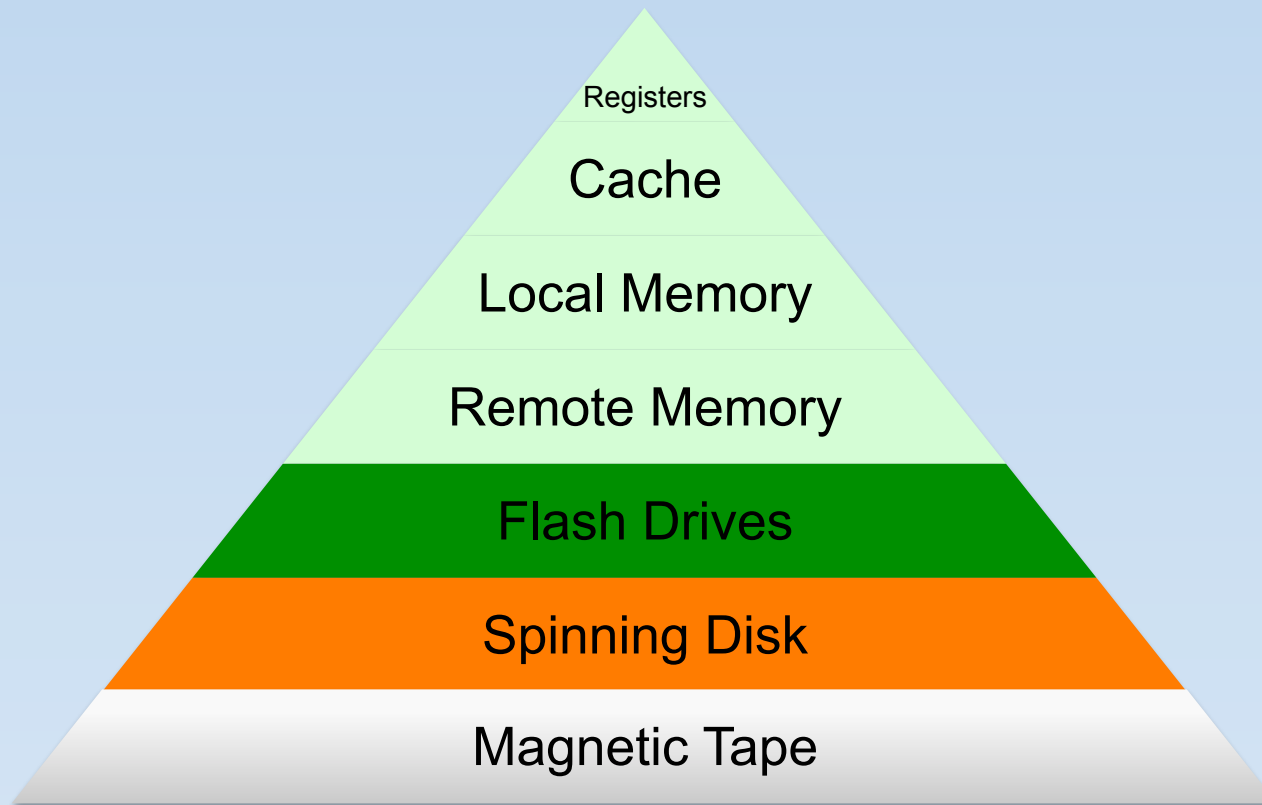# Infrastructure Challenges of Big Data

Red Shift: Data keeps moving further away from the CPU with every turn of Moore's Law

- ◆ CPU Cycle Time
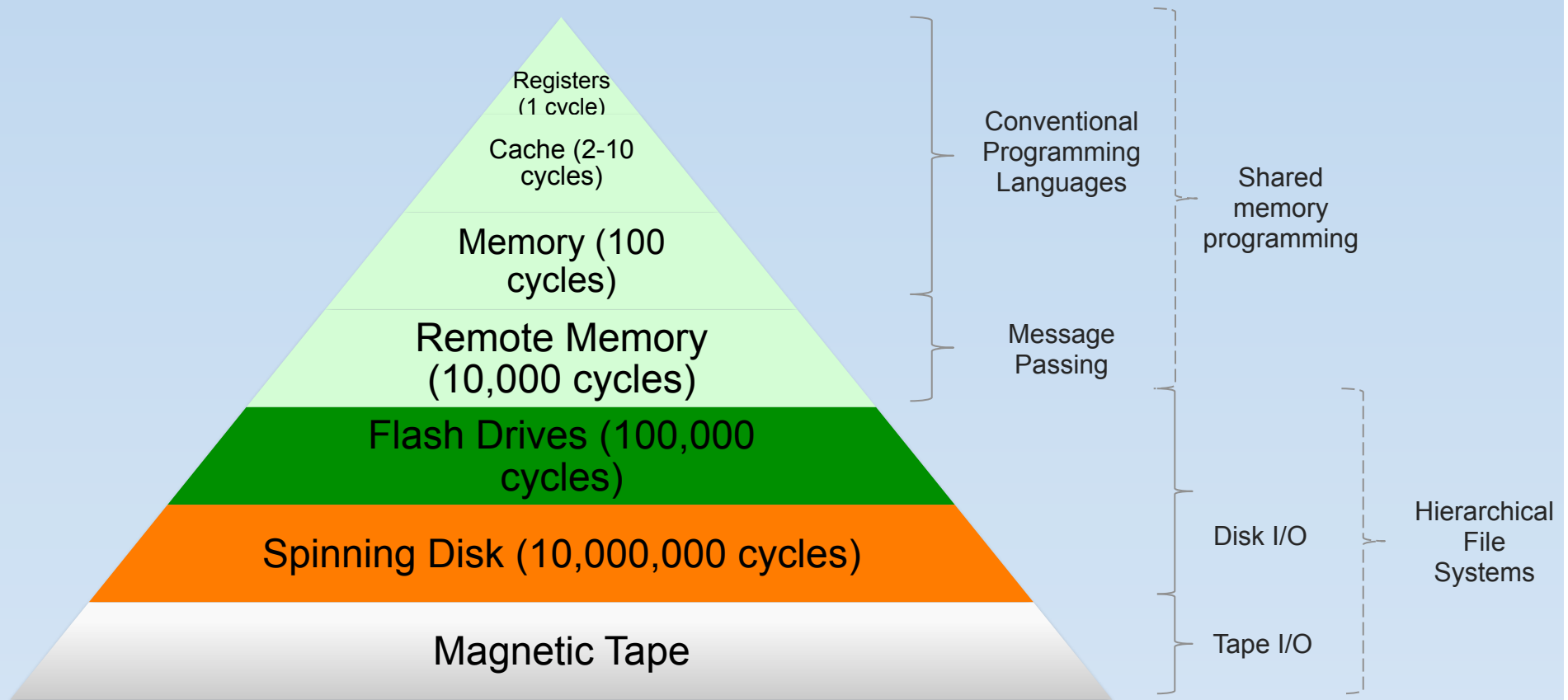- ▲ Multi Core Effective Cycle Time
- ■ Memory Access Time
- ▬ Disk Access Time

data due to Dean Klein of Micron

Slide courtesy Mike Norman, SDSC

# It's always been caches all the way down

Registers

Cache

Local Memory

Remote Memory

Flash Drives

Spinning Disk

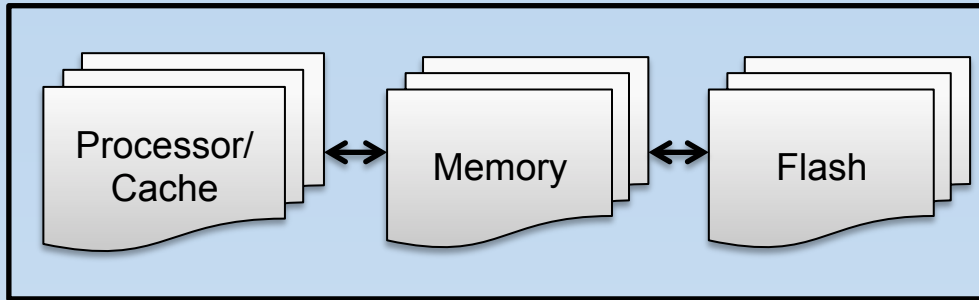Magnetic Tape

Explicit vs Implicit management

# Infrastructure for Data Intensive Computing

- Computation
  - Large amounts of main memory
  - Parallel processors
  - Smooth out memory pyramid

- Storage
  - Significant long term storage
  - Smooth out the memory pyramid
  - Many views of same data
    - Parallel File System
    - Local access (POSIX)
    - Remote collaboration and sharing (Object store)
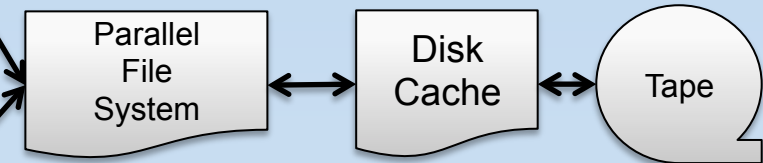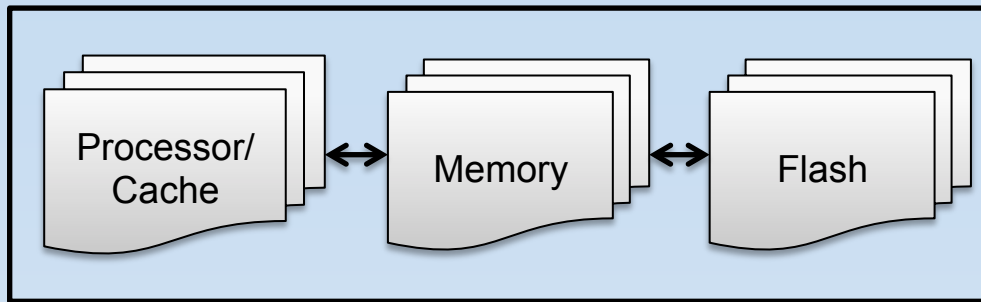    - Sync-and-share
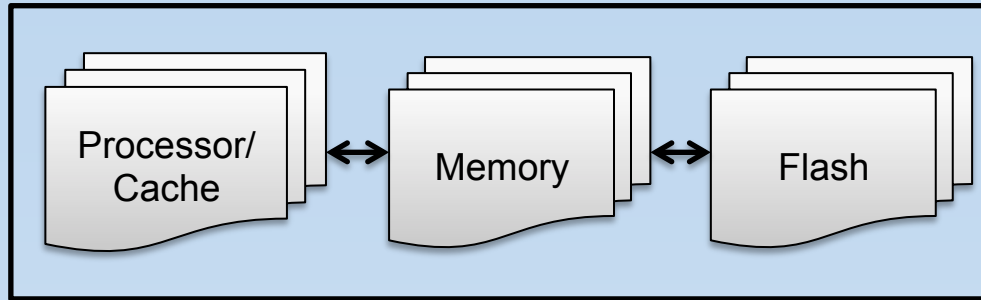    - Web
    - Cloud

Reference Architecture

Cluster B

Processor/
Cache ↔ Memory ↔ Flash

Cluster A

Processor/
Cache ↔ Memory ↔ Flash

Parallel
File
System ↔ Disk
Cache ↔ Tape

Shared Memory
Programming

Hierarchical File
System

Reference Architecture

Cluster B
- Processor/Cache
- Memory
- Flash

FlashLite
- Processor/Cache
- Memory
- Flash

Parallel File System
Disk Cache
Tape

Shared Memory Programming

Hierarchical File System

# Data Intensive Computation Engine

- Parallel
  - High performance network
  - Good numeric performance

- Massive memory
  - Ability to hold whole data sets or data bases in memory
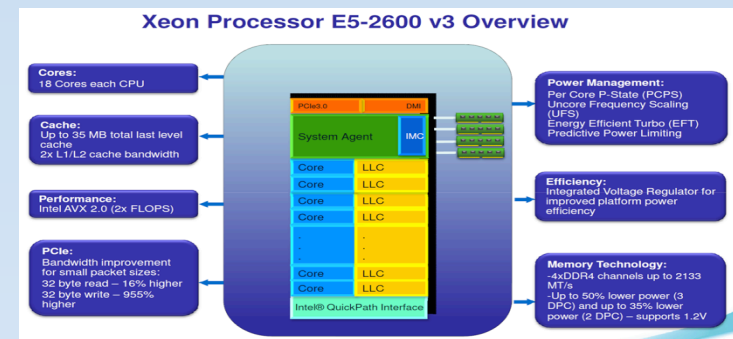
- High IO throughput

# FlashLite

- High throughput solid state disk

- Large amounts of main memory

- Software shared memory

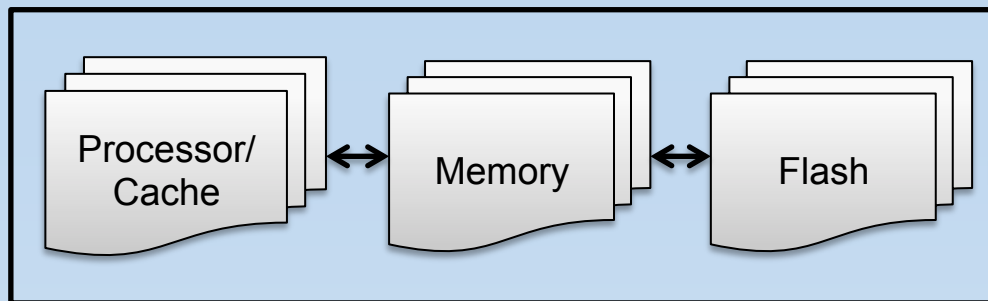- Inspired by SDSC Gordon
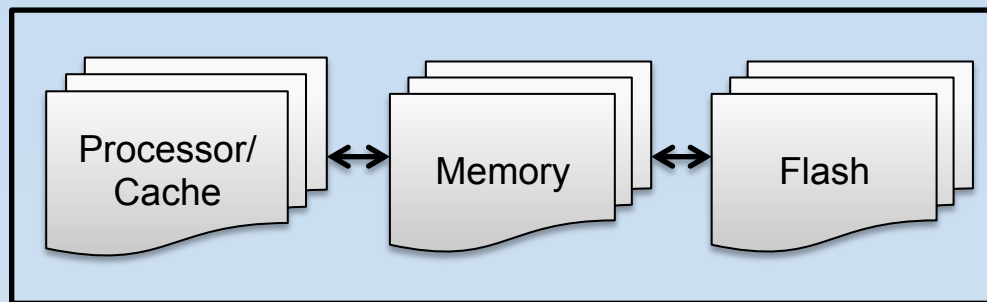
# What is FlashLite?

- FlashLite
  - ~ 70 compute nodes (~1600 cores)
    - Dual socket Intel E5-2680v3 2.5GHz (Haswell)
    - 512 GB DDR-2
    - 4.8 TB NVMe SSD
  - ScaleMP vSMP virtual shared memory
    - 4TB RAM aggregate(s)
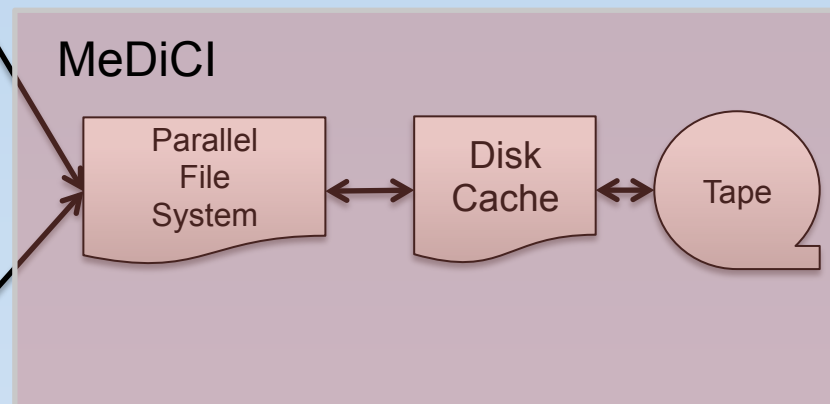
Reference Architecture

Cluster B

Processor/Cache ↔ Memory ↔ Flash

FlashLite

Processor/Cache ↔ Memory ↔ Flash

MeDiCI

Parallel File System ↔ Disk Cache ↔ Tape

Shared Memory Programming
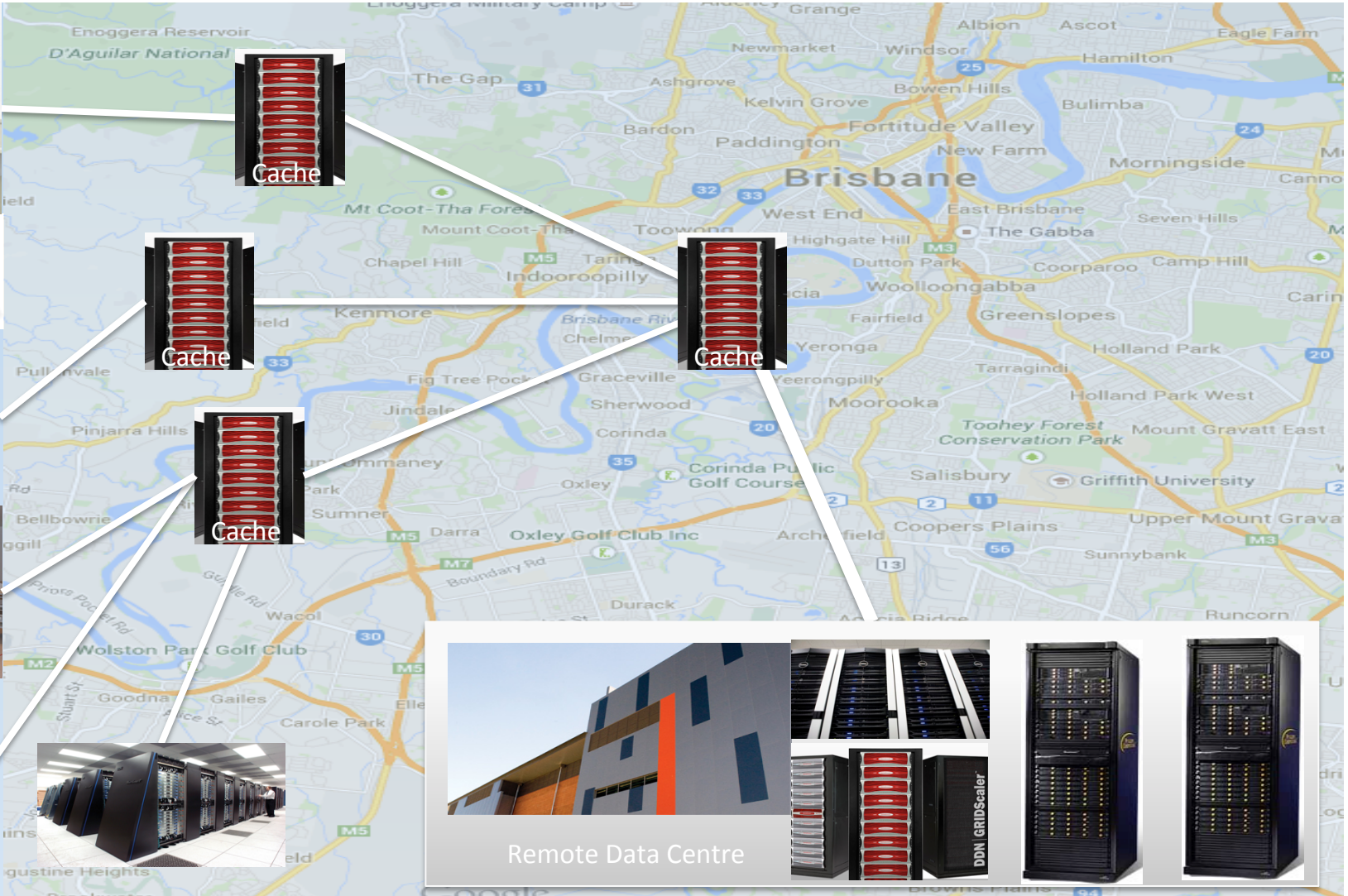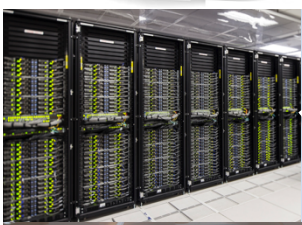
Hierarchical File System
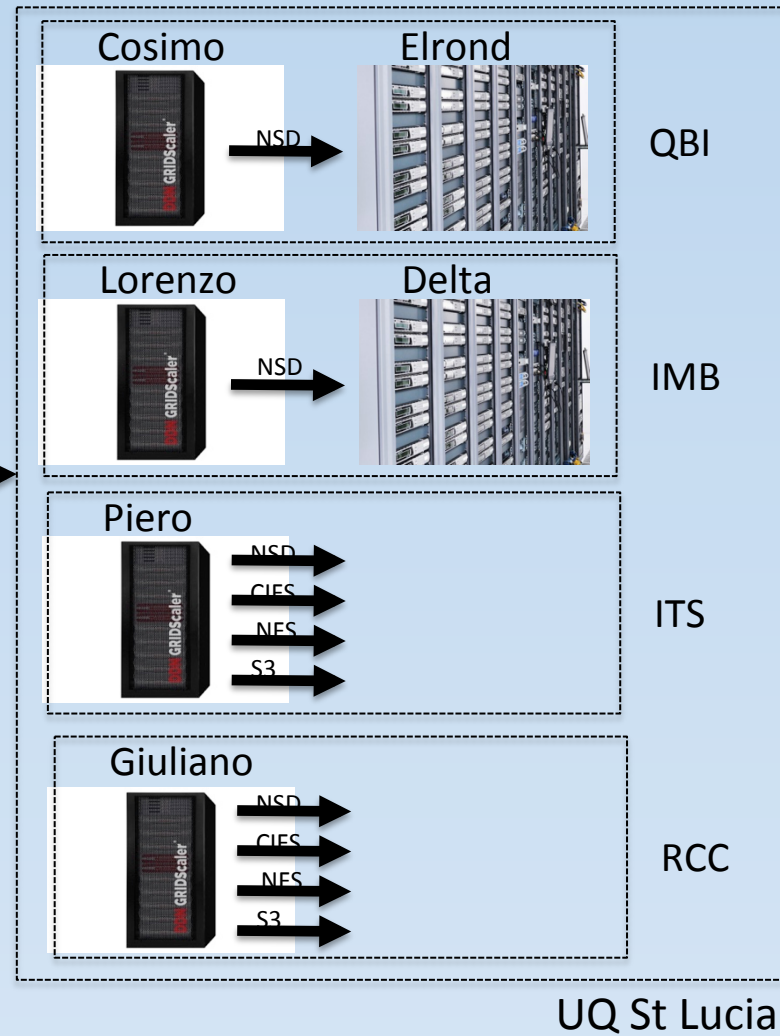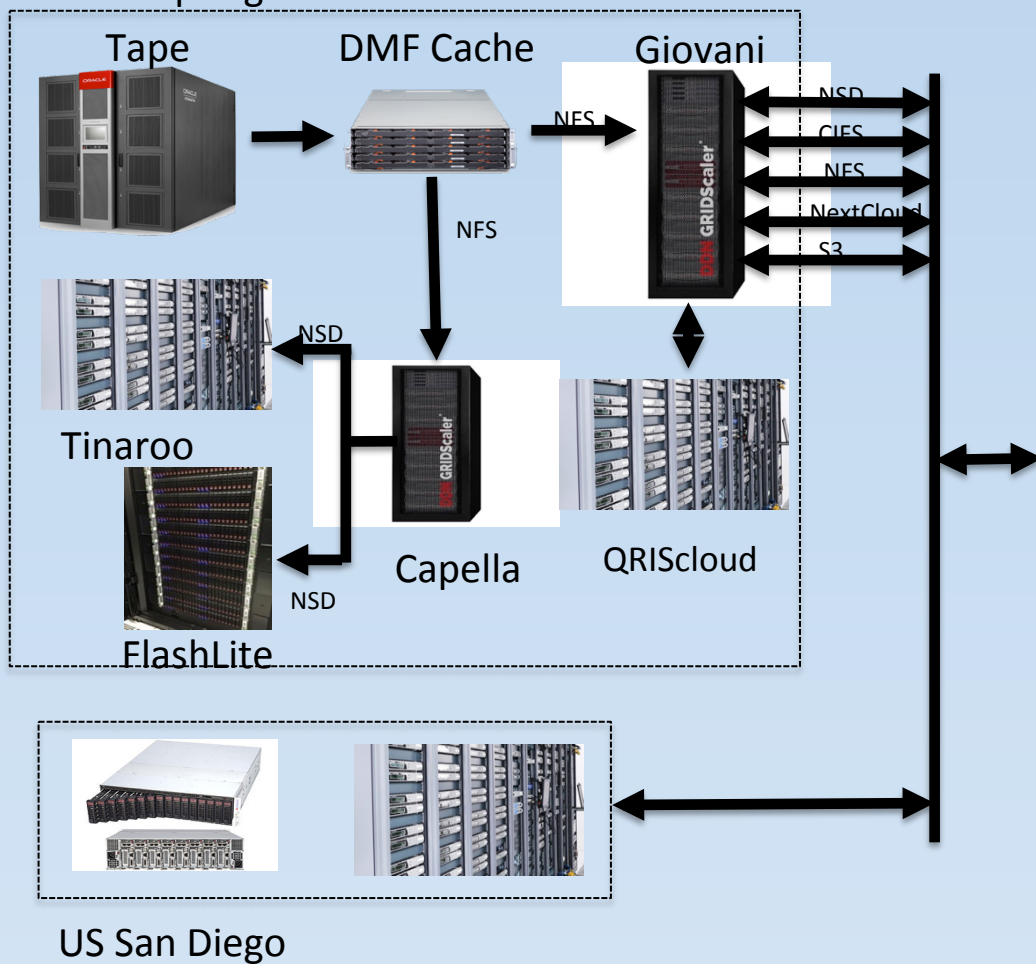
The caches continue …

# MeDiCI

# MeDiCI

- Centralising research data storage and computation
- Distributed data is further from both the instruments that generate it, some of the computers that process it, and the researchers that interpret it.
- Existing mechanisms manually move data
- MeDiCI solves this by
  - Augmenting the existing infrastructure,
  - Implementing on campus caching
  - Automatic data movement
- Current implementation based on IBM Spectrum Scale (GPFS)

Cache

Cache

Cache

Cache

Remote Data Centre

DDN GRIDScaler

# FlashLite in the Data Centre



GPFS

GPFS

GPFS

GPFS

DDN SFA12KXE

FlashLite

Parallel file system

# FlashLite in the Data Centre



GPFS

GPFS

GPFS

GPFS

NFS

**Automatically Recall**
Full or partial file restore

**Second Tier**
Lower cost & performance
capacity disk array

**Third Tier**
Tape or VTL
Active Archive

DDN SFA12KXE

SGI DMF Disk/Tape

FlashLite

Parallel file system

Long term data collections

# MeDiCI Wide Area Architecture

SGI DMF Disk/Tape

GPFS

GPFS

GPFS

GPFS

NFS

Automatically Recall
Full or partial file restore

Second Tier
Lower cost & performance
capacity disk array

Third Tier
Tape or VTL
Active Archive

Machine
Room
Network

Parallel
File
System

Compute

Long term data collections

Cache

# MeDiCI Wide Area Architecture



SGI DMF Disk/Tape

GPFS

GPFS

GPFS

GPFS

GPFS

GPFS

GPFS

GPFS

NFS

Machine Room Network

Wide Area Network

Long term data collections

Compute

Parallel File System

Cache

Cache

# MeDiCI Wide Area Architecture

SGI DMF Disk/Tape

GPFS GPFS GPFS NFS

GPFS GPFS

GPFS GPFS

GPFS GPFS

Automatically Recall
Full or partial file restore

Second Tier
Lower cost & performance
capacity disk array

Third Tier
Tape or VTL
Active Archive

BitsPerSecond - bes-1.router.uq.edu.au  Ethernet1/48  AARNet 10GE to St Lucia
Fri Jul 22 00:00 2016 to  Fri Jul 22 14:45 2016    Resolution: 1 min   Speed: 10G

Receive  Transmit

10.0G
8.0G
6.0G
4.0G
2.0G
0

22 Jul    1am    2am    3am    4am    5am    6am    7am    8am    9am    10am    11am    noon    1pm    2pm
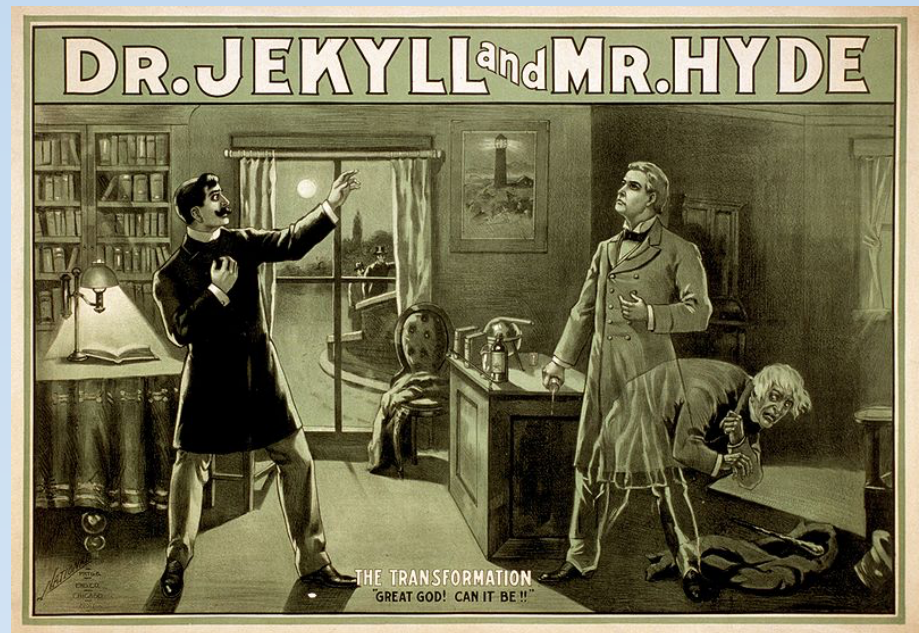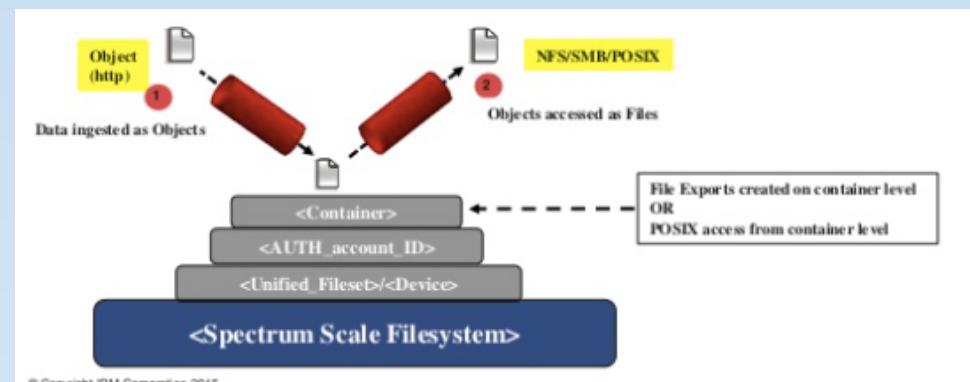
# Identity!

- No single UID space across UQ/QCIF users
- Need to map UID space between UQ and Polaris
- GPFS 4.2
  - mmname2uid/mmuid2name

# Object Storage

- S3 style objects becoming defacto standard for distributing data

- http put/get protocol

- Swift over GPFS
  - Unified Object/file interfaces
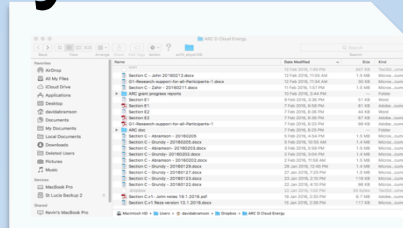
# Data Data everywhere anytime



ImageTrove    myTardis    OMERO

Managed Data

MeDiCI    ownCloud
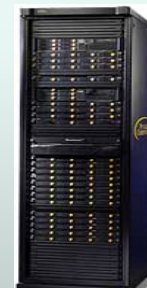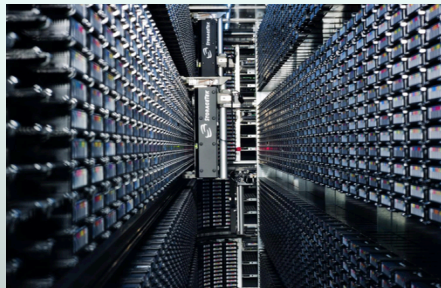
Synchronous    Asynchronous

Unmanaged Data

OpenClinica

Clinical Data

S3, Swift

Cloud Access

MeDiCI

QRIScloud Compute and Storage Fabric

# Building on basic architecture

- A Declarative Machine Room

- Alternative backends

- Leveraging Cloud Storage

- Very Very Wide Area File Systems

- Supporting repository stacks

- Orchestrating Workflows

# Conclusions

- FlashLite
  - Parallel computer
  - Very large amounts of local memory and Flash disk
  - Still learning what works
  - Need Burst Buffer s/w
- MeDiCI
  - Caches all the way down
  - IBM Spectrum Scale
  - AFM semantics

# Acknowledgments