# Collaborative Data Life-cycle Management (CDLM) for Petascale Projects

## Arun Jagatheesan

## iRODS.org, DICE, SDSC/UCSD

# Agenda

- **Introductions**
- **LSST as use case**
- **CDLM**
- **Attributes of CDLM**

# **History behind the story**

- **MDAS (Massive Data Analysis System)**
  - Support data-intensive applications that manipulate very large data sets by building upon object-relational database technology and archival storage technology
  - 1995 by DARPA
- **SDSC SRB (Storage Resource Broker)**
- **iRODS**
  - Flexible license for our community
  - Flexible rules for users
  - Flexible data management

# My role in iRODS Community

- Large-scale usage and adoption of iRODS
  - Research and Analysis of large-scale use-cases
  - Design requirements for large-scale users
  - Consult on iRODS-based storage infrastructure
- Community Growth
  - Tutorials, dissemination
  - iROD-Chat (2006), SRB-Chat (2003)
  - Academic and Industrial users

# Large Scale Synoptic Survey

- Survey entire sky every 3 nights
- Dark Energy, Dark Matter, Near Earth Asteroids, and more
- World's largest digital camera (3 billion pixels)
- Images 3000 times wider than Hubble
- Data from Chile to US and rest of the world
- 15 TB/night, over hundred(s) petabytes

- www.youtube.com/watch?v=LtMJ_WwvBb8

# Data Products

QuickTime™ and a
TIFF (Uncompressed) decompressor
are needed to see this picture.

- **Releases**
- **Cataloged database**
- **Provenance Info**
- **Metadata**
- **Processed Data Sets**
- **Raw Images**

# LSST Data Infrastructure Layout

# LSST Data Train and iRODS

/file1..10.fits
/nobel.event



/file1..10.fits

/file1..10.fits
/catalog1.db

/catalog1.db

UK or IN2P3

/file1..10.fits
/catalog1.db

/file1..10.fits
/catalog1.db

# LSST CDLM Problem Statement

- **LSST data-lifecycle management infrastructure for:**
  - Performance oriented data storage sub-systems
  - Capacity oriented data storage sub-systems
  - Data (usage oriented) distribution networks
  - [Provenance and archive storage systems]
- **Confluence of three major storage dimensions**
  - HPC data processing  (pipelines to produce our data)
  - Datacenter  sharing   (data centers that host our data)
  - Data delivery and distribution (usage of our data)

# CDLM

- **Collaborative Data Lifecycle Management**
  - Multiplexing of a single data life-cycle amongst more than one autonomous partner
  - Attributes of data-lifecycle is shared
  - Varying levels of autonomy and inter-dependence

# Multiplexing a Data Life-cycle

- **Data Creation (Raw data)**
- **Data Processing (Derived data)**
- **Data Analysis (Data warehouse, ..)**
- **Data Namespace**
- **Data Dissemination**
- **Data Provenance**
- **Data Archival**

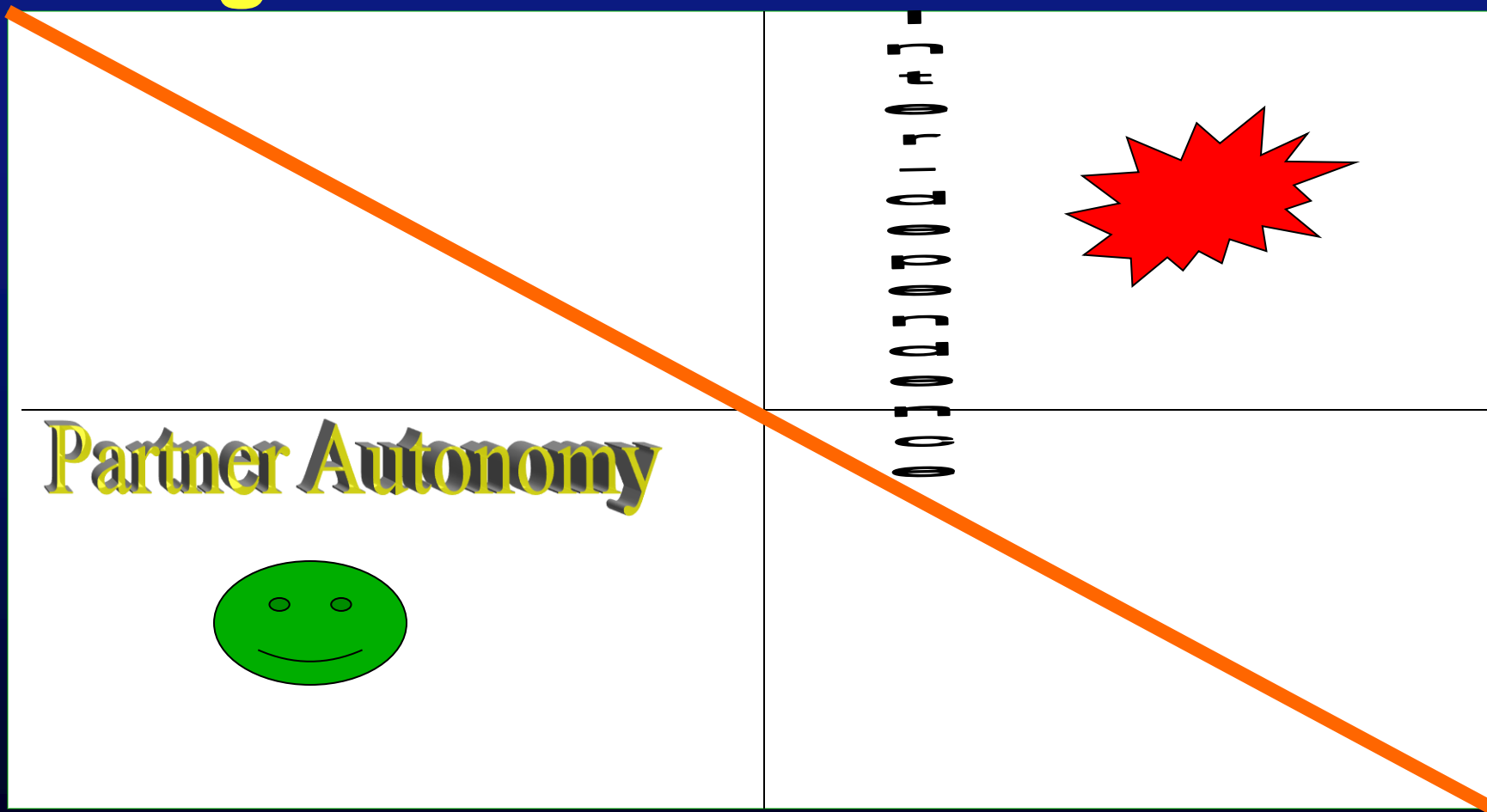# Levels of Collaboration

- **Collaboration on Data Life-cycle not necessarily mean collaboration of businesses**
- **Some types of CDLM**
    - Symbiotic - All partner businesses benefit from CDLM
    - Neutral - No effect on businesses due to CDLM
    - Competitive - partners of CDLM are actually competitors of the resulting business process (forced to have a common platform to compete)
    - Hybrid - Multiple or transient partner relationships

# Autonomy & Inter-dependence at right levels for CDLM to work

# LSST Data Layout

# ALMA data flow

# LSST SC-2008 Prototype

# CDLM Infrastructure Design

- **Requirements, Expectations and Performance Management**

- **Minimize dependencies (without affecting cost)**

- **Reduce individual autonomy into hierarchical groups (that can remain autonomous)**

- **Hierarchical rules and community rules**

# iRODS enabling CDLM

- **Global Namespace**
- **Resource allocation and service levels as policies/rules**
- **Hierarchical rules and access controls**
- **Highly Flexible System**

# Similar projects? Let's talk

- **The power of the community**
- **Not necessarily "large" scale**
- **Symbiotic**
- **arun@diceresearch.org**