# VANDERBILT/JPL COLLABORATION ON A SCALABLE DATA PROCESSING PIPELINE FOR PROTEOMICS
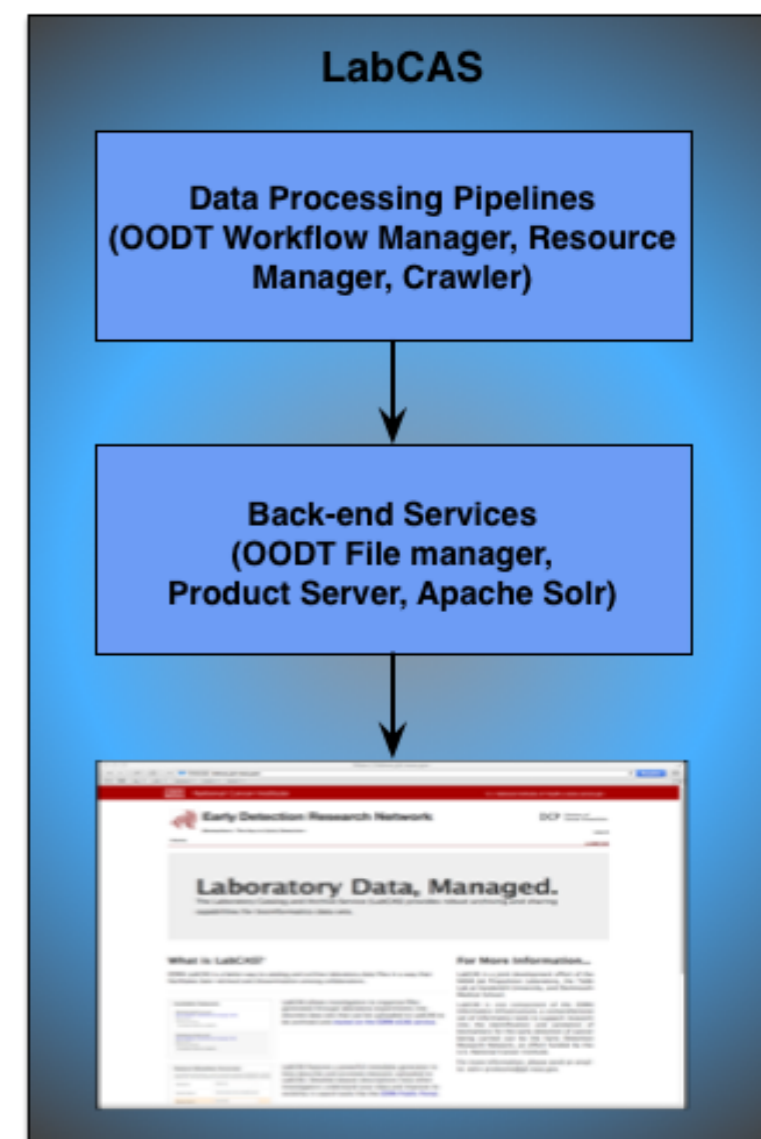
**Kirsten Anton, Luca Cinquini, Maureen Colbert, Dan Crichton (PI), Thomas Fuchs, Sean Kelly, Heather Kincaid, Ashish Mahabal, Chris Mattman, Rishi Verma, Paul Zimdars**
California Institute of Technology & Jet Propulsion Laboratory (NASA)

**Matthew Chambers, David Tabb (PI)**
Vanderbilt University

- <u>LabCAS is an IT integrated environment for managing biomedical laboratory data</u>: generation, publication, documentation, search, discovery and access
- Overall goal: support research for identification and validation of cancer biomarkers for early detection of cancer - focus on early sharing of results before publication
- Developed by JPL in collaboration with EDRN partners (Vanderbilt, Dartmouth, BU…)

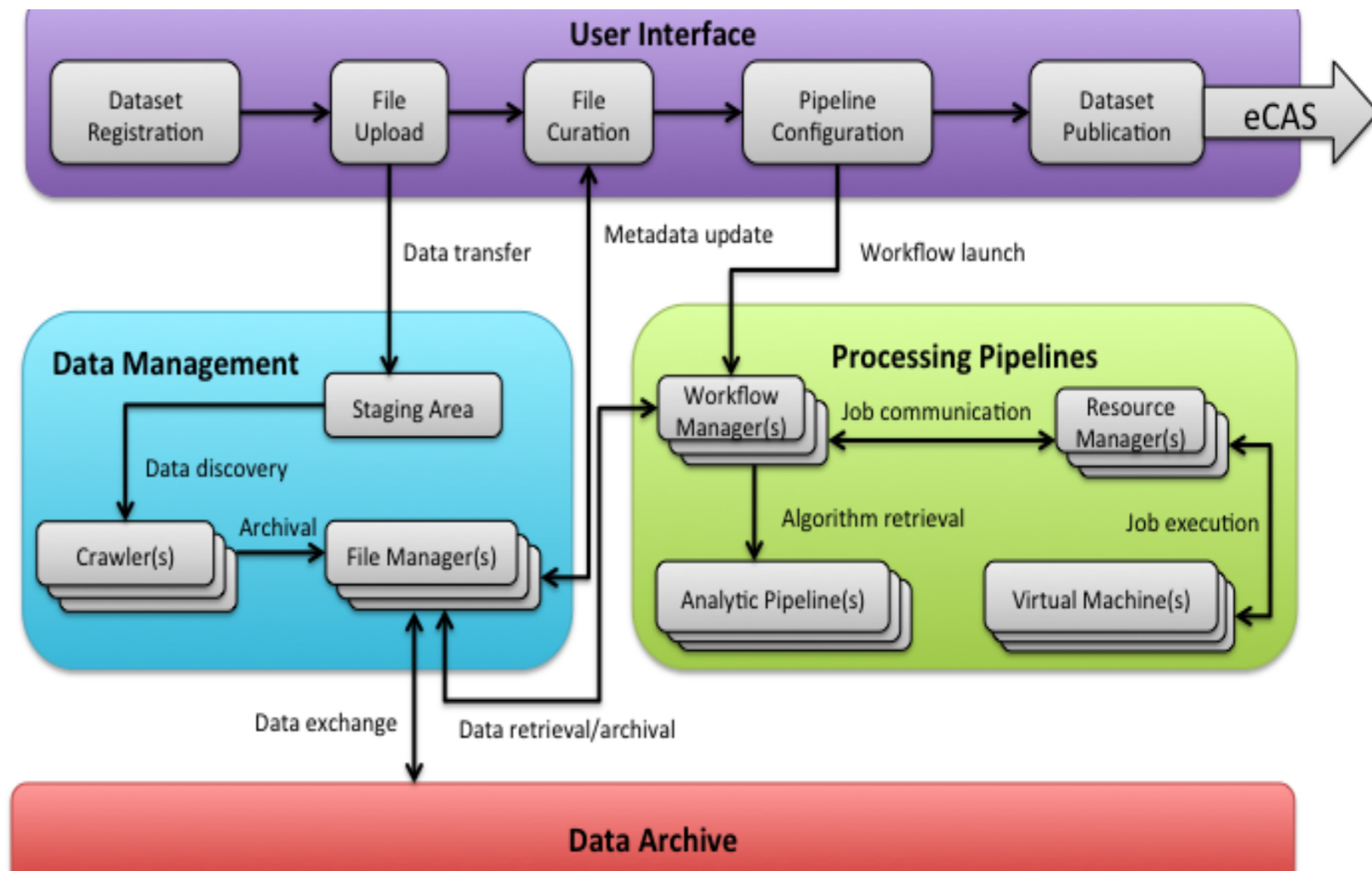**Architecture Software Layers**

- <u>Data Processing Pipelines</u>:
  ‣ framework for processing laboratory data as integrated workflows, generate data products, and publish them to the archive
- <u>Back-end Services</u>:
  ‣ components for publishing, searching and downloading data products
- <u>Front-end Web Portal</u>:
  ‣ public site for authorized access of data products, as well as UI for execution of data processing pipelines



**LabCAS**

Data Processing Pipelines
(OODT Workflow Manager, Resource Manager, Crawler)

Back-end Services
(OODT File manager,
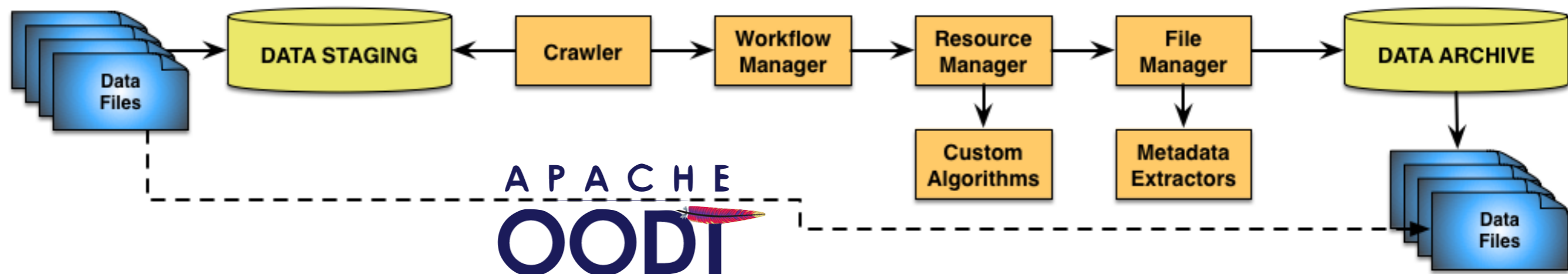Product Server, Apache Solr)

Laboratory Data, Managed.

Our group is working at establishing a state-of-the-art computing environment at JPL for execution of biology data processing pipelines for the EDRN program.
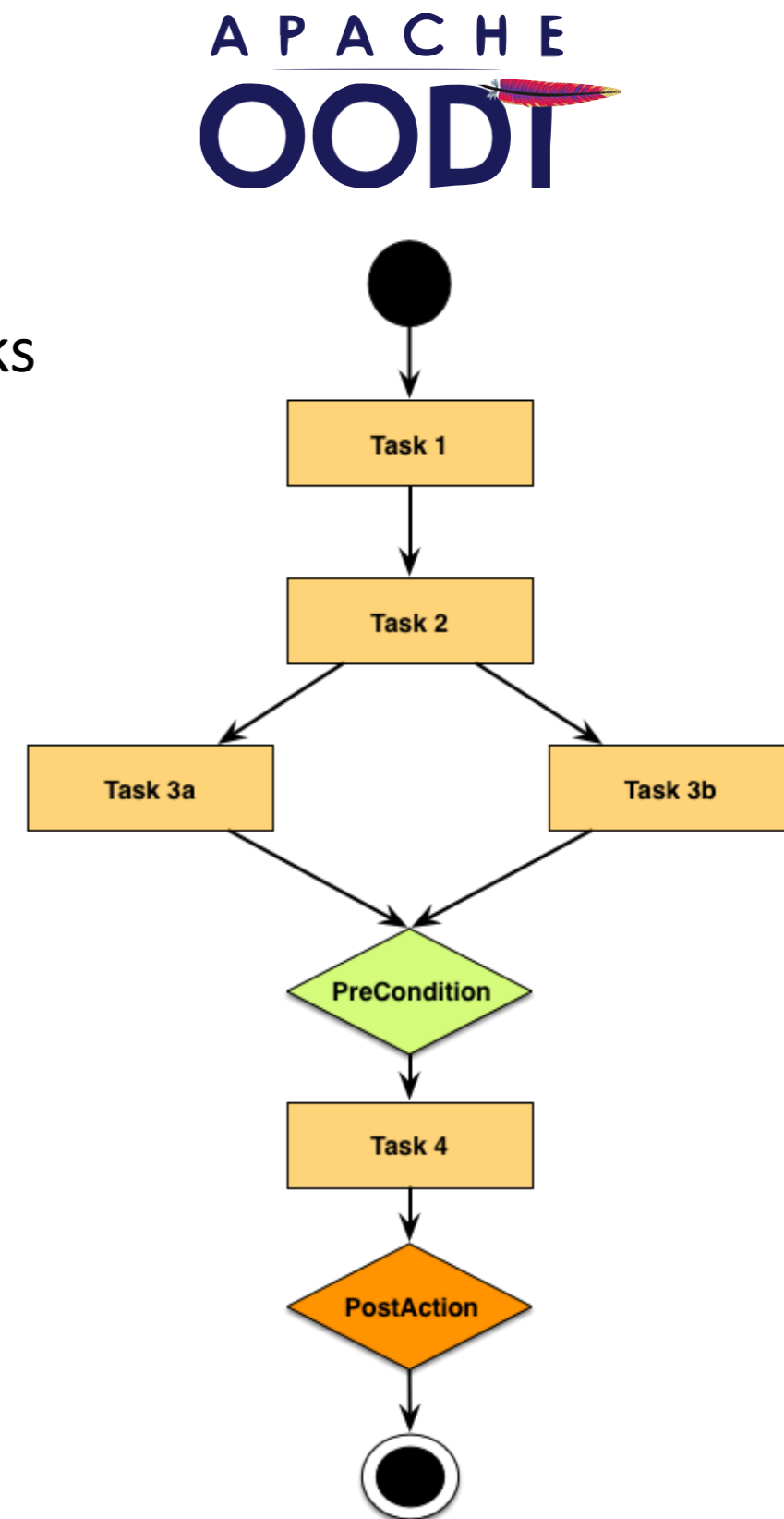
### Goals

- <u>Performance</u>: easy configuration, distributed processing, support multiple pipelines
- <u>Reproducibility</u>: capture detailed metadata about workflow execution ("provenance") so that other investigators can reproduce and validate the results
- <u>Sharing</u>: publish data and metadata to the web portal for authorized access by other investigators and the general public

- Apache OODT is an Open Source framework for management, processing, discovery and access of scientific data
- Modular architecture allows to instantiate and combine different components to realize the most appropriate architecture for a specific data processing environment
  - ▸ File Manager: data access server and metadata catalog. May be backed up by Apache Solr - web-enabled, high performance search engine
  - ▸ Workflow Manager: general workflow engine for execution of pipelines composed of sequential or parallel tasks
  - ▸ Resource Manager: allocates computing resources for task execution
  - ▸ Crawler: service for monitoring data spaces and to trigger ingestion of files into the File Manager or submission of jobs to Workflow Manager

- Each Workflow Manager server can be configured to execute one or more workflows
- A client starts a workflow by sending an "event" with optional configuration metadata
- Each workflow is composed of an arbitrary number of tasks (sequential or parallel)
- For scalability, tasks can be run on the local Resource Manager, or sent to Resource Managers on other servers
- Optional pre-conditions cause task execution to wait until they are satisfied
- Optional post-actions trigger operations when a workflow terminates (on success or failure)
  ‣ Example: start crawling for products to ingest
- Workflow products are categorized according to custom types:
  ‣ Specific metadata elements
  ‣ Specific archive location and versioner
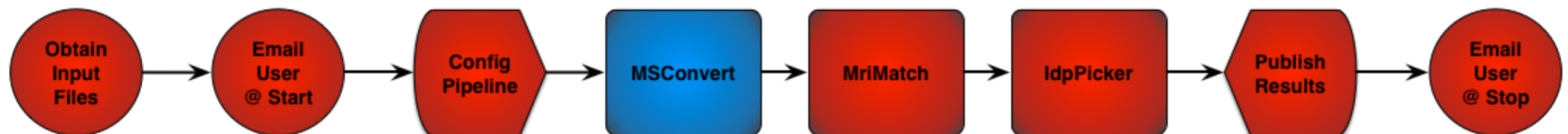  ‣ Specific metadata parsers

# LabCAS Pipelines Collaborations

JPL is collaborating with several EDRN partners to enable their data processing pipelines to be executed within the LabCAS environment:

- First, understand and run the different stages of a pipeline as standalone tasks
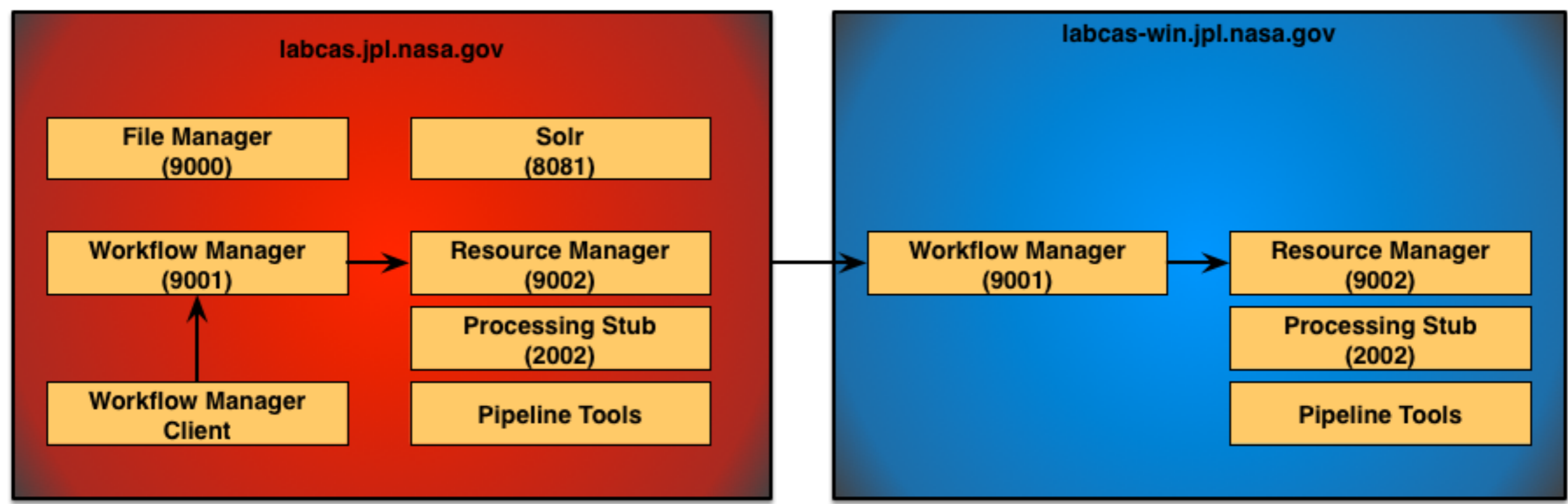- Then, instrument all stages as single runnable OODT workflow

- <u>Vanderbilt University</u> (PI: Dave Tabb)
  - ▸ Proteomics
- <u>Boston University (PI: Marc Lenburg)</u>
  - ▸ Biomarker Discovery
  - ▸ RNA Sequencing
  - ▸ Microarray pre-processing
- <u>University of Washington (PI: Alvin Liu)</u>
  - ▸ Microarray data
- <u>Memorial Sloan Kettering Cancer Center (PI: Gunnar Ratsch)</u>
  - ▸ Genomics
- <u>Cedar-Sinai Medical Center (PI: Beatrice Knudsen)</u>
  - ▸ Pathology
- <u>Cedar-Sinai Medical Center (PI: Michael Freeman)</u>
  - ▸ Proteomics

# Proteomics Pipeline

- Data processing pipeline for identification and analysis of protein cancer biomarkers in body fluids
- Developed by David Tabb's group at Vanderbilt, part of CPTAC (Clinical Proteomic Tumor Analysis Consortium) activities
- Composed of several programs (part of ProteoWizard suite of open source tools):
  ▸ MSConvert (pre-processing): conversion of RAW files to mzML format
  ▸ MyriMatch (database search): searches sample data for peptides (.pepXML)
  ▸ IdpPicker (filtering: IdPQonvert+IdPAssemble+IdPQuery): generates protein identification reports
- Instrumented as single OODT workflow executing a sequence of 13 tasks
- First test case at JPL consisted in running the pipeline on a medium-size pancreatic dataset composed of 675 input files of approx. 200MB each (total: 1.35TB)
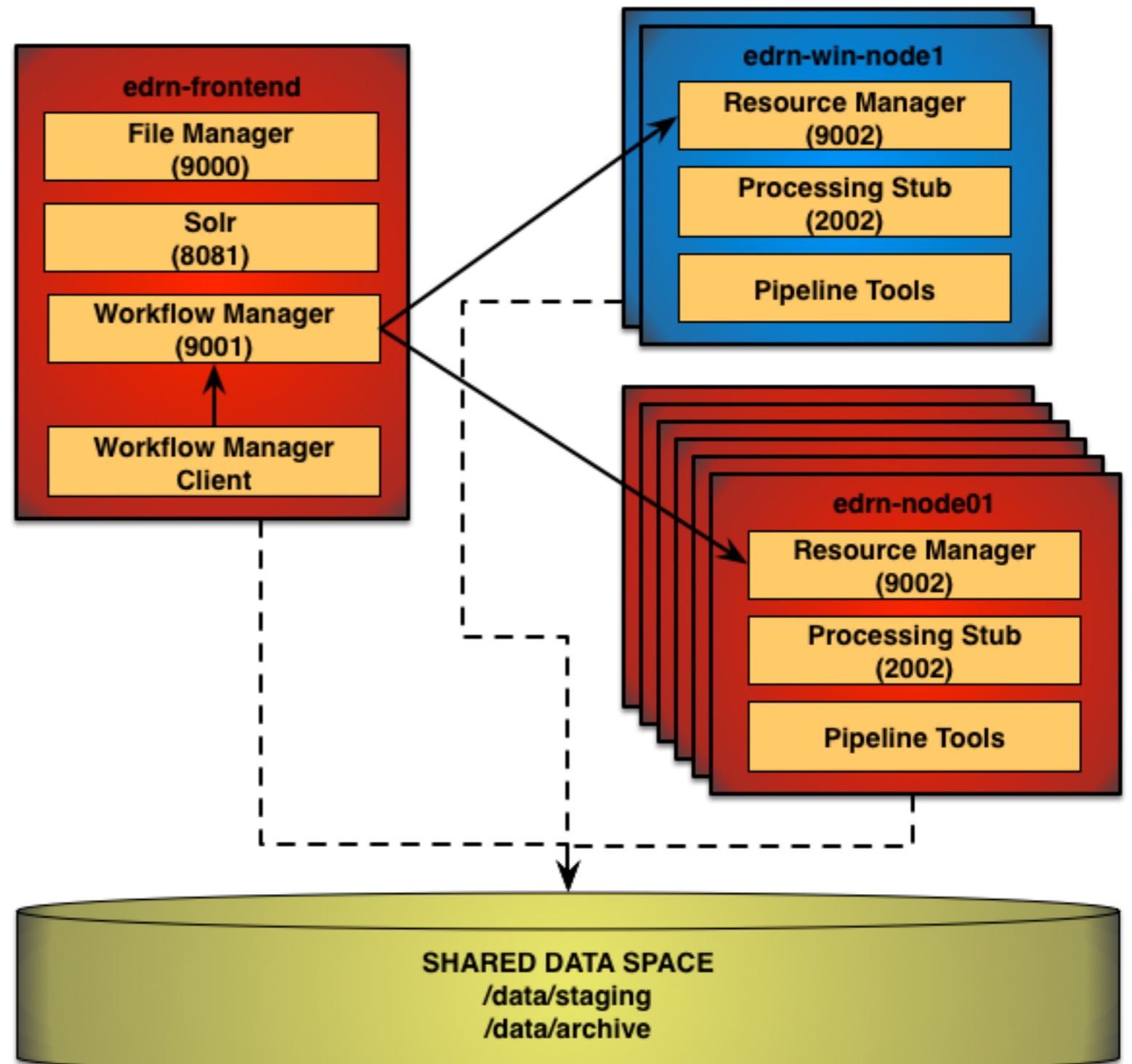
- Use system of 1 Linux + 1 Windows servers (2 cores each)
- Each task setup to run over all files sequentially
- Tasks 1-7 (up to MSConvert) completed in about 1 week on Windows server
- Task 8 (MyriMatch) started on Linux server, stopped after estimating 40+ days before completion (approximately 1.5 hours for each file)
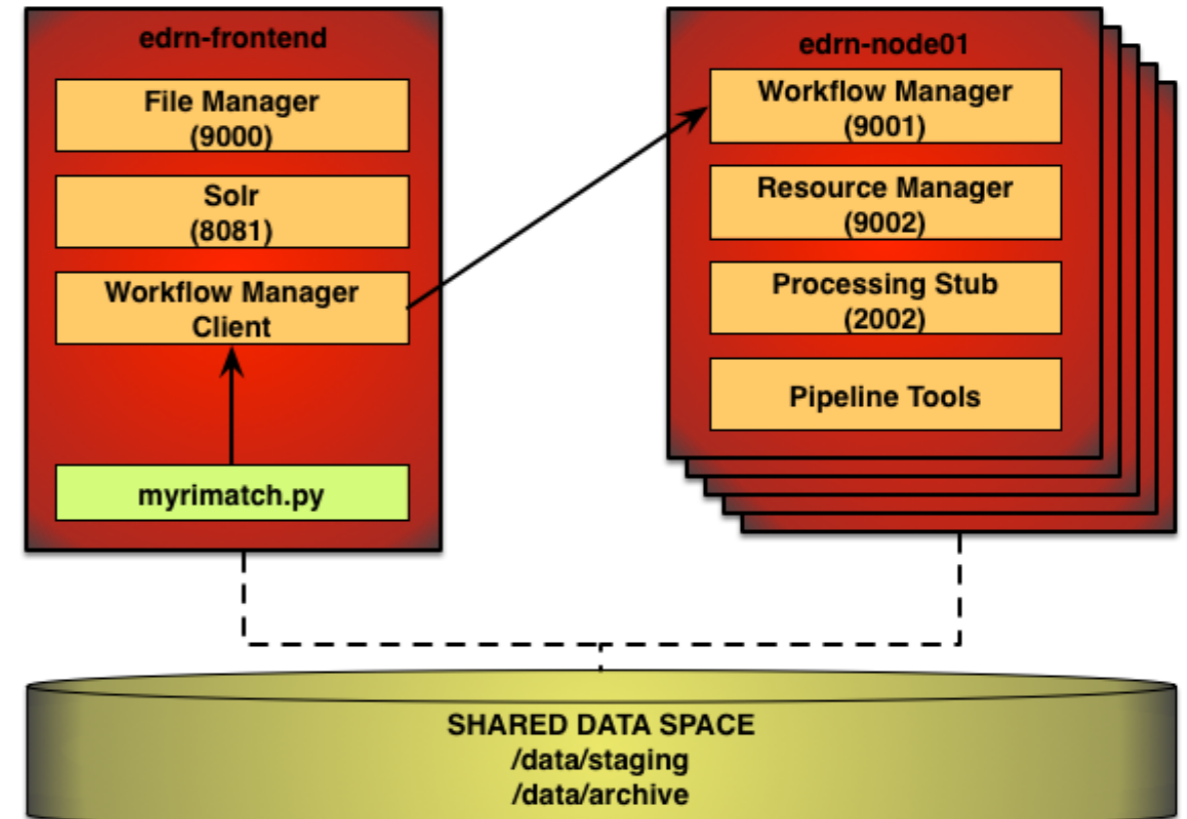
- System of 21 Virtual Machines (6 cores each)
  - ‣ 1 front-end server hosting common services
  - ‣ 18 Linux back-end processing nodes
  - ‣ 2 Windows back-end processing nodes
  - ‣ 10TB shared storage
- Automatic software replication from 1 processing node to all the others
- Execution of commands from front-end node to all the back-end nodes
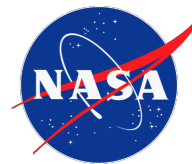  - ‣ Start/stop OODT services on all node simultaneously

- Use new JPL/EDRN cluster
- MyriMatch step was identified as bottle-neck in Proteomics pipeline
- Because MyriMatch only processes 1 mzML file at a time, this step can be fully parallelized!
- Using mzML files already produced by the old pipeline, sub-workflow was setup to run MyriMatch in parallel over all available Linux nodes

## Results

- Single MyriMatch job completes in approx. 30min (instead of 90+ min)
- Full MyriMatch sub-workflow completed in approx. 1 day (for all 675 files)!
- Performance improvements:
  - ▸ Factor of 3 from using more powerful hardware
  - ▸ Factor of 18 from parallelizing a single step onto multiple processing nodes

# JPL/EDRN Cluster Usage (1 day)



Show Hosts Scaled: ● Auto ○ Same ○ None | EDRN-Cluster **load_one** last **day** sorted **by name** | Size [ small ⇅ ] Columns [ 4 ⇅ ] (0 = metric + reports)

(Nodes colored by 1-minute load) | Legend

# JPL/EDRN Cluster Usage (1 day)

DAWN (Distributed Analytics, Workflows and Numeric): model to simulate and optimize Big Data computational pipelines, developed at JPL under the Data Science initiative

Analysis of Proteomics distributed workflow
(using rough benchmarking):
- Clear reduction in overall elapsed time when MyriMatch step is executed on multiple nodes
- Efficiency gain levels off around 15 nodes
- Another gain is obtained by using 2 Windows nodes instead of 1
- Full pipeline should complete in 4-5 days



EDRN CPTAC Proteome Pipeline
Number of Files: 675

Number of Windows Nodes: 1
Number of Windows Nodes: 2
Number of Windows Nodes: 3
Number of Windows Nodes: 4

## Proteomics Pipeline

- Re-run and benchmark full pipeline on new hardware for pancreatic test case
  - ▶ Submitted "report" to Vanderbilt to fix IdpPicker bug on Linux OS
- Run pipeline on other datasets
- Finalize pipeline product metadata
- Publish generated products to LabCAS

## LabCAS

- Total redesign of Front-end Web Portal
  - ▶ Streamlined uploads of data files
  - ▶ Flexible UI for execution of data processing pipelines
  - ▶ Enhanced searching of datasets (by keywords or facets)
  - ▶ Improved access to documentation
- Back-end Services and Pipelines
  - ▶ Ability to package specific workflows for deployment at other institutions
  - ▶ Develop test suite
  - ▶ Upgrade to latest OODT release