# NCSA

# Medici in a nutshell

# Medici in a nutshell

- A multimedia content management system based on

  - Web 2.0 interfaces

  - Semantic web technologies (RDF)

  - Cloud-based processing and preprocessing

# Motivations

- Address research and education data collection and analytic needs
  - Manage large collections of heterogeneous data
  - Organize data with metadata and provenance information
  - Facilitate collaborations and data sharing
  - Enable curation and data preservation
- Support community collections of heterogeneous data (documents, images, video, sensor, modeling, etc.)
- Enable automated data extraction, analytic and preprocessing services on local and remote systems
- Provide data preview capabilities specific to different data types

# Why Not Flickr or YouTube

- Maybe?
  - Web-accessible tools are relatively generic
  - Users like not having to manage storage
  - Metadata, tagging, linking, etc. are effective means of organizing information (i.e., no need for "folders")

- Maybe not?
  - No individual or community ownership
  - No control of resources
  - Inadequate privacy (e.g., for unpublished work)
  - Limits on format, volume, throughput, resolution
  - No domain-specific processing
  - No provenance (everything is a stream of "posts")
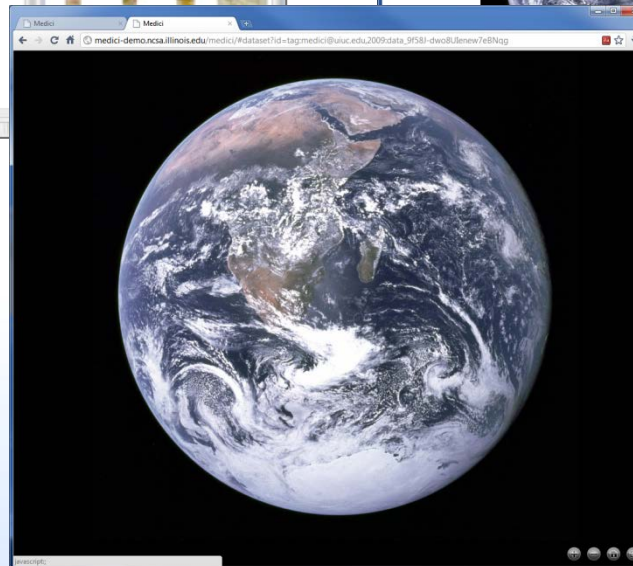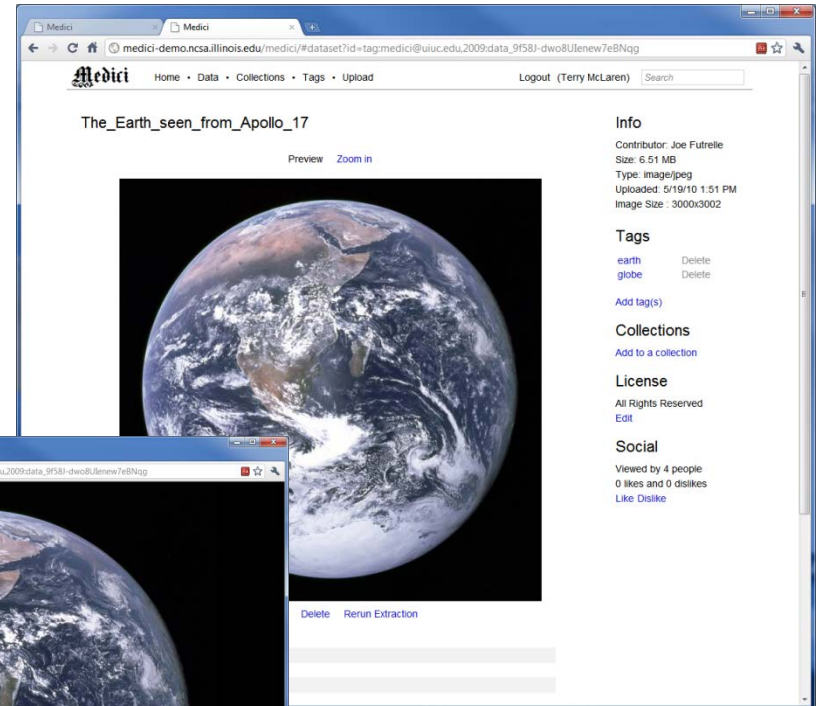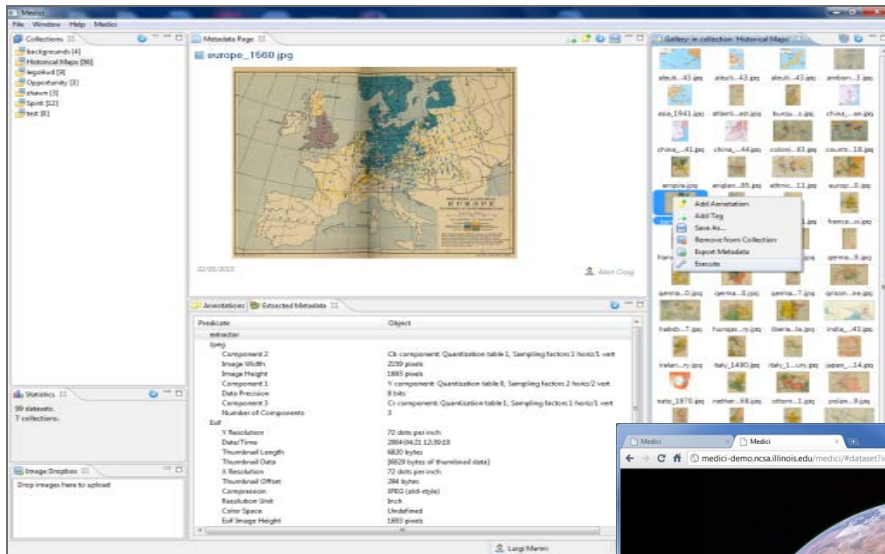
# Community Drivers

- Use cases and requirements driven by
  - US Office of Naval Research (ONR)
  - US National Archives and Records Administration (NARA)
  - US National Endowment for the Humanities (NEH)
  - US National Institute of Health (NIH)
  - US National Science Foundation (NSF)
  - Institute for Advanced Computing Applications and Technologies (IACAT)
  - Seagrant/EPA
  - EU LinkSCEEM-2 (Cyprus Institute)

# Why Medici?

- Provides a customizable turn key solution to store, organize, analyze, view, share, and preserve research content

- Supports heterogeneous files
  - Single file or directory upload via click-n-drag
  - RESTful web service for batch or script based uploading
  - Owner defined copyright and download permissions

- Standards based (RDF) semantic content model

- Conforms to open data and metadata standards
  - Supports OPM (Open Provenance Model)
  - Tags, comments, ratings

- Supports customizable automated extraction services
  - E.g.: Image pyramid creation, OCR for scanned text, movie frame extraction (.mpeg), file transformation, etc…

- Leverages proven technologies
  - Lucene indexing, MySQL, any command-line tool for extraction/analytic services

- Open source, public APIs
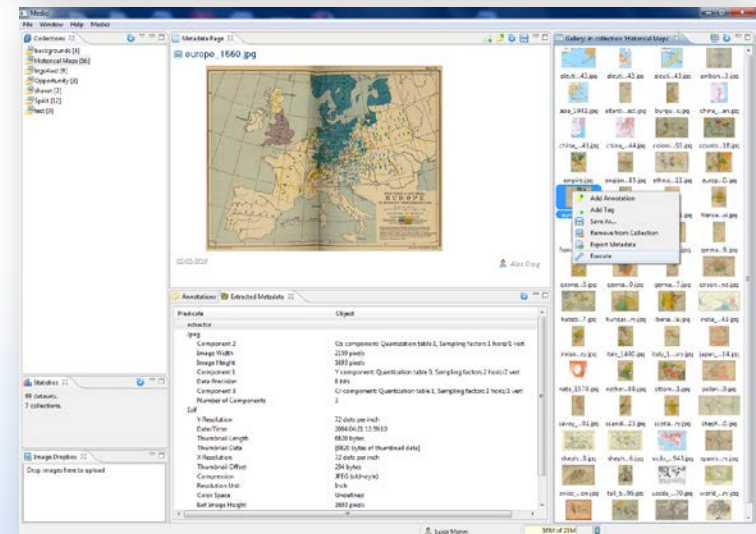
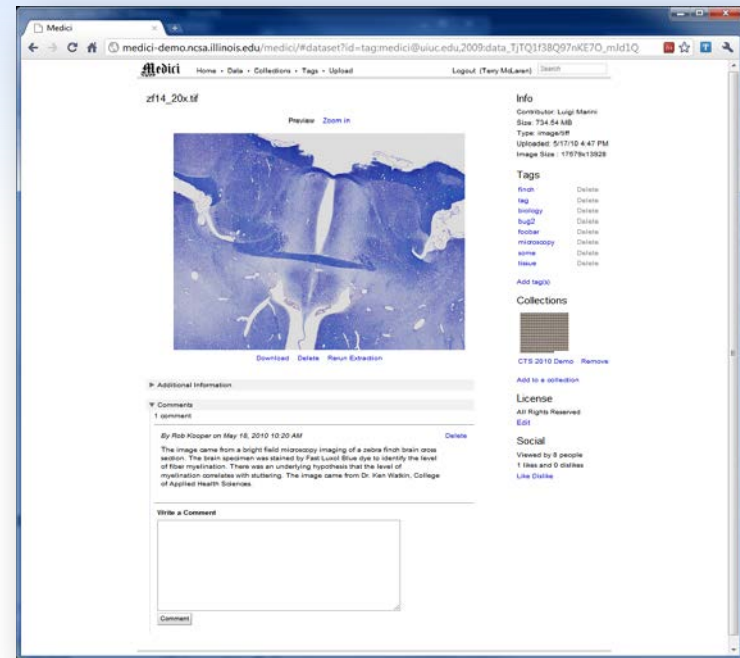# Medici – Semantic Data Repository

- Web and Desktop access to a semantic content repository.
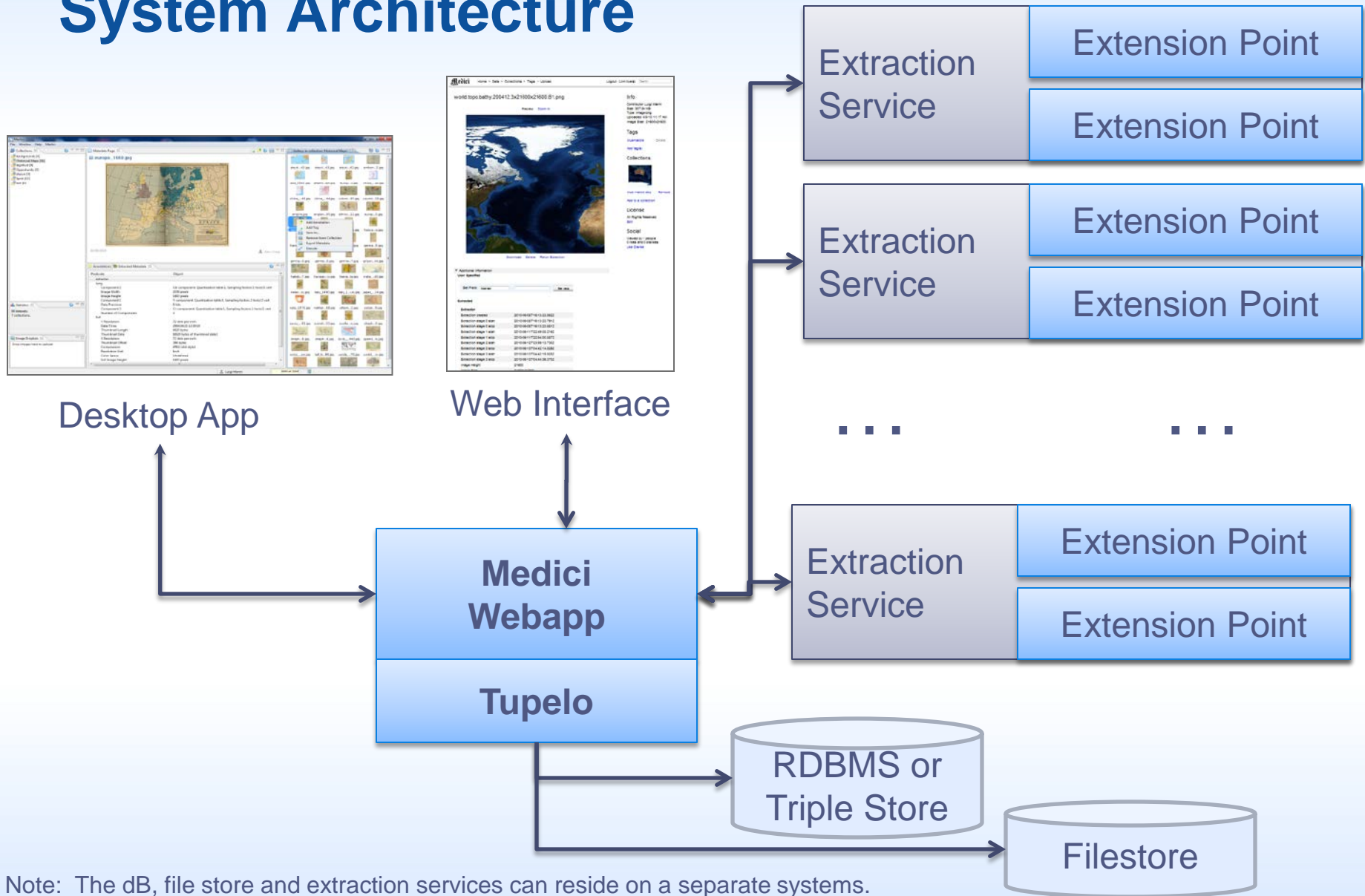


- web 2.0 interfaces
- Semantic web technologies (RDF)
- cloud-based processing and preprocessing

# Client features

- Upload / download

- Search / browse

- Tag / comment

- Create collections

- Geo-locate data (map view)

- Content-type-specific previewing
  - e.g., zoomable images (Seadragon), playable movies (jwplayer), rotatable 3D objects (HTML5)

- Define a specific taxonomy

- Access statistics, provenance

- Citable persistent URLs

- Set copyright and license attributes
  - View only, prevent download

- Define dataset relationships

# System Architecture



Desktop App

Web Interface

Extraction Service — Extension Point / Extension Point

Extraction Service — Extension Point / Extension Point

. . .          . . .

**Medici Webapp**

**Tupelo**

Extraction Service — Extension Point / Extension Point

RDBMS or Triple Store

Filestore

Note: The dB, file store and extraction services can reside on a separate systems.
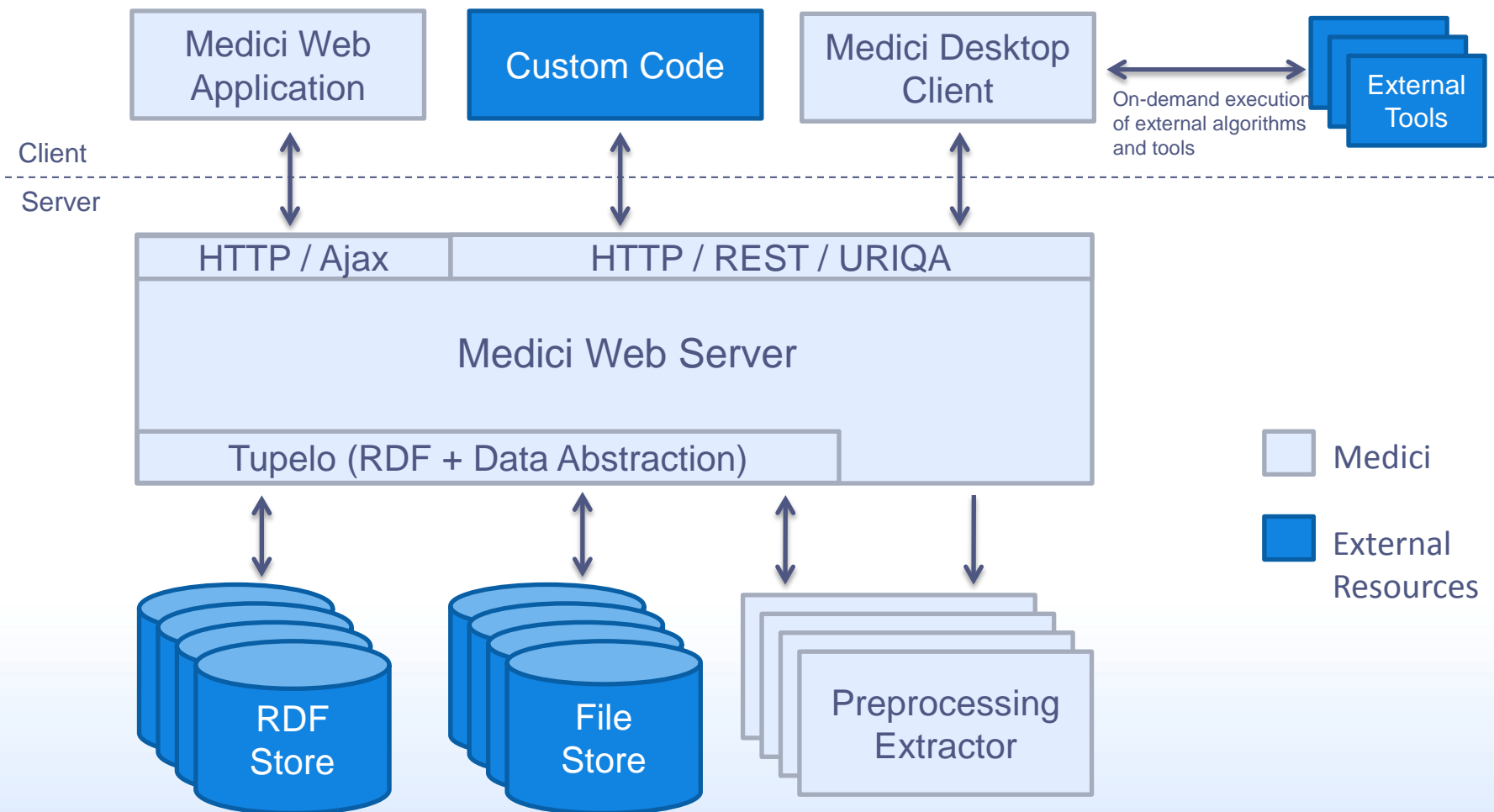
# Extraction services

- Multiple, extensible pre-processing pipelines
- Asynchronous, distributed, triggered by upload
  - Processing selected on basis of file content type (MIME type)
  - Recursive (products can trigger additional extractions)
- Used to produce web-viewable previews
  - Image pyramids, audio/video previews, thumbnails, pdf to plain text
- Used for domain-specific pre-processing, e.g.,
  - Metadata extraction (e.g., FITS headers, geolocation)
  - Feature detection
  - Specialized OCR for non-standard textual types (e.g., 18th-century manuscripts)

NCSA

# Medici Technologies

- Web application
  - Google Web Toolkit
  - Java Servlets
  - Plain Javascript
  - Viewers: Flash, Java Applet, HTML, etc.
  - Apache Lucene
  - Mysql
- Extraction Service
  - Eclipse RCP (Java)
  - Large collection of external applications
- Desktop Client
  - Eclipse RCP (Java)
  - Cyberintegrator Workflow Management System

# Software Architecture

# Medici Communities

- Cyprus Institute (Digital Cultural Heritage)
  - 3D object archive for artifacts
- Datanet: sead.ncsa.illinois.edu
- Medici-demo.ncsa.illinois.edu
  - An open public server, upload requires account
- InvertNet.org
  - Digitization of Biological Collections
- Digging into Data
  - University of Sheffield, MATRIX Center
  - Given a set of images of historical artefacts, discover what salient characteristics make an artist different from others using computational image analysis
  - Enable statistical learning about individual and collective authorship.
- Walker Institute – Rule of Law
  - Repository of reports, video, satellite images for the Rule-of-Law in different locations around the world
- 18Connect – OCR of 18th century Manuscripts
  - Institute for Computing in Humanities, Arts, and Social Science
  - Extraction service to OCR manuscript images using Gamera OCR toolkit

# Acknowledgements

- Institute for Advanced Computing Applications and Technologies (IACAT)
  - UIUC Campus collaboration
- NIH - (Image repository)
- iChass - (Digging into Data, 18Connect)
- NSF - (InVertnet, Datanet:SEAD)
- EPA / Seagrant
- Cyprus Institute Collaboration

# For more information please visit

http://medici.ncsa.illinois.edu

email us at

medici@ncsa.illinois.edu

join the discussion at

medici-users@ncsa.illinois.edu

NCSA