# Big Data in Biomedical Imaging: Pathomics Analysis of Cancer

Tahsin Kurc
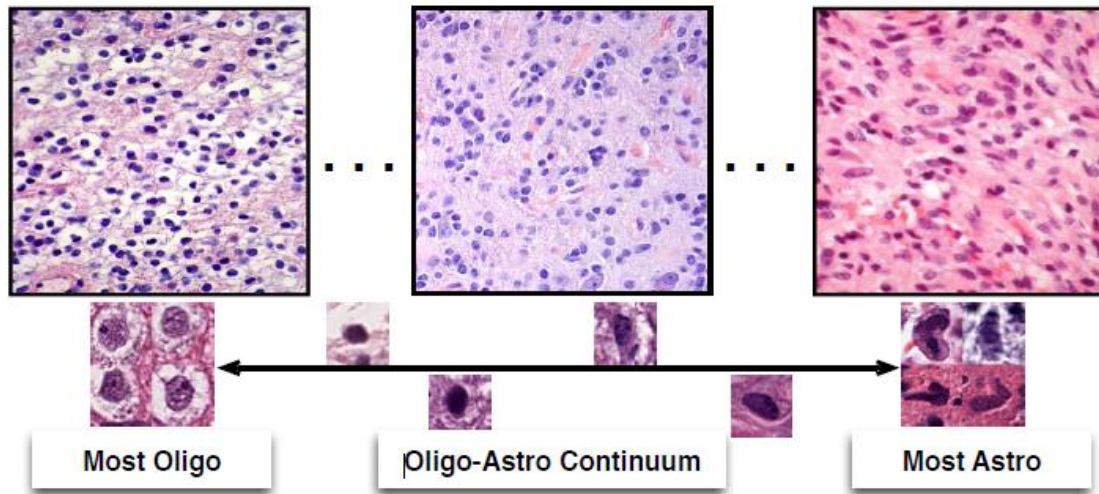
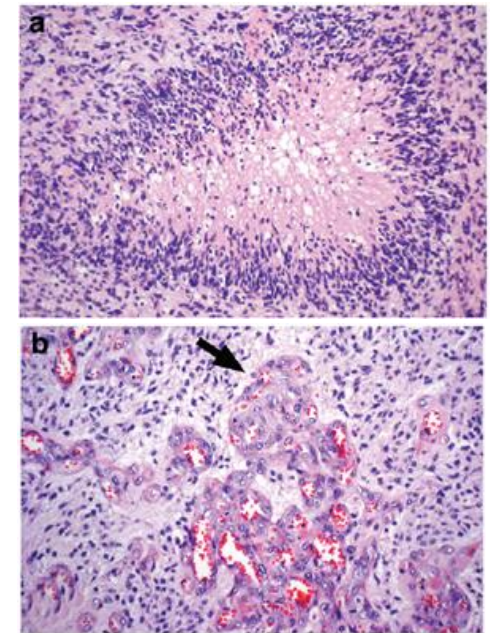Biomedical Informatics Department
Stony Brook University

# Team

- Stony Brook:
  - Joel Saltz, Erich Bremer,  Feiqiao Wang, Tammy Di Prima, Jonas Almeida, Mary Saltz, Le Hou, Vu Nyugen,  Si Wen, Maozheng Zhao, Tianhao Zhao, Raj Gupta, Fusheng Wang, Furqan Baig, Yi Gao, Alina Jasniewski

- Emory University:
  - Ashish Sharma, Ganesh Iyer, Adam Marcus

- Oak Ridge National Laboratory:
  - Scott Klasky, Jeremy Logan, David Pugmire

- University of Brasilia, Brazil
  - George Teodoro, Luís F. R. Taveira, Alba C. M. A. Melo

# Tissue Image Analysis

Different types of nuclei and cells

Regions of necrosis and angiogenesis



Most Oligo

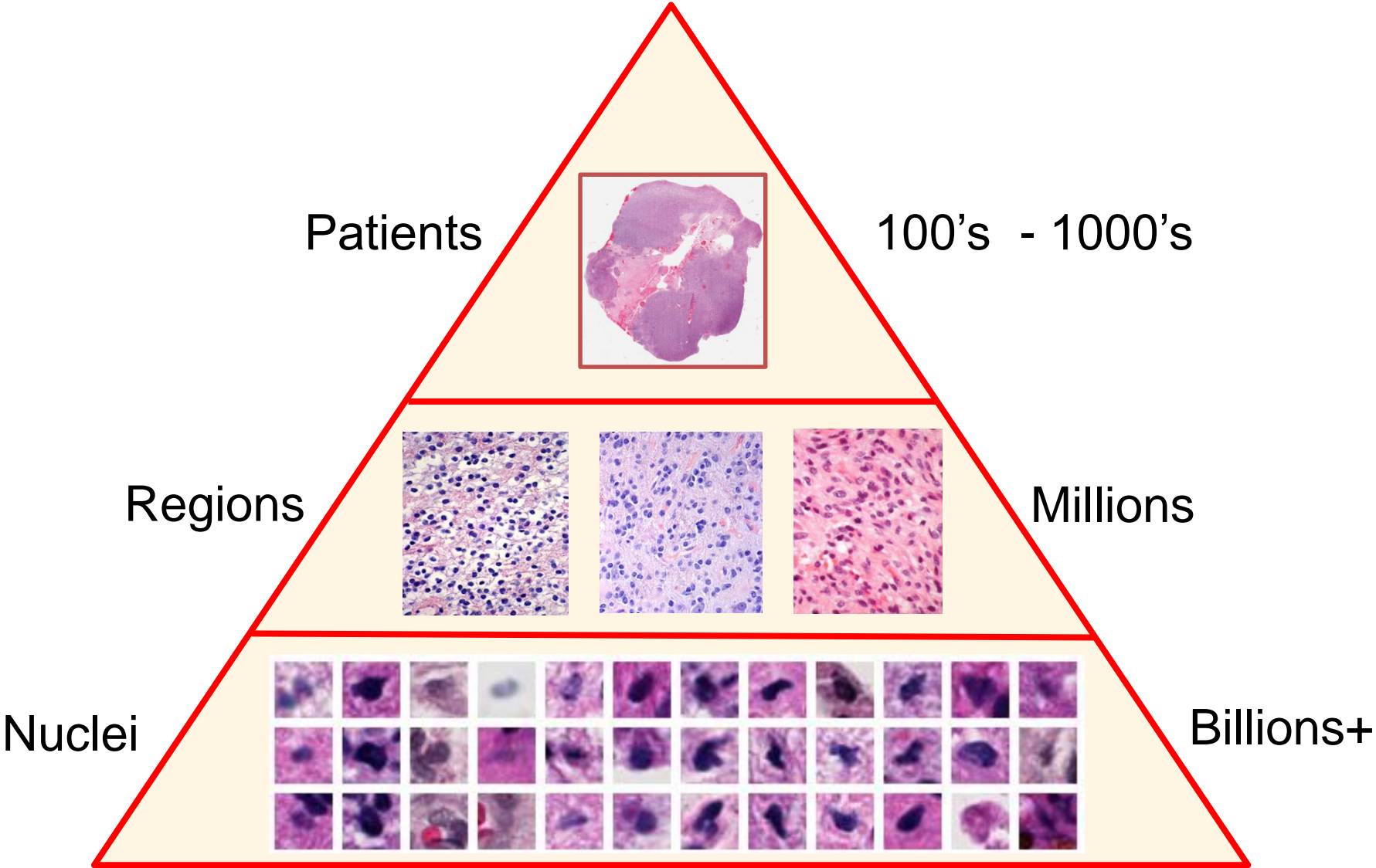Oligo-Astro Continuum

Most Astro

- Tumors are complex and heterogeneous
- Datasets of rich nuclear morphological information
  - Nuclear morphology
  - Maps of tumor infiltrating lymphocytes

# Tissue Imaging

- Confluence of several technological advances is making slide scanning more practical

- Image scanning technology has progressed significantly in recent years
  - Sophisticated auto-focus mechanisms
  - Slide trays for batch scanning of slides

- Time required to scan a slide at high-resolution has reduced from multiple hours to several minutes
  - 200-300 slides / day

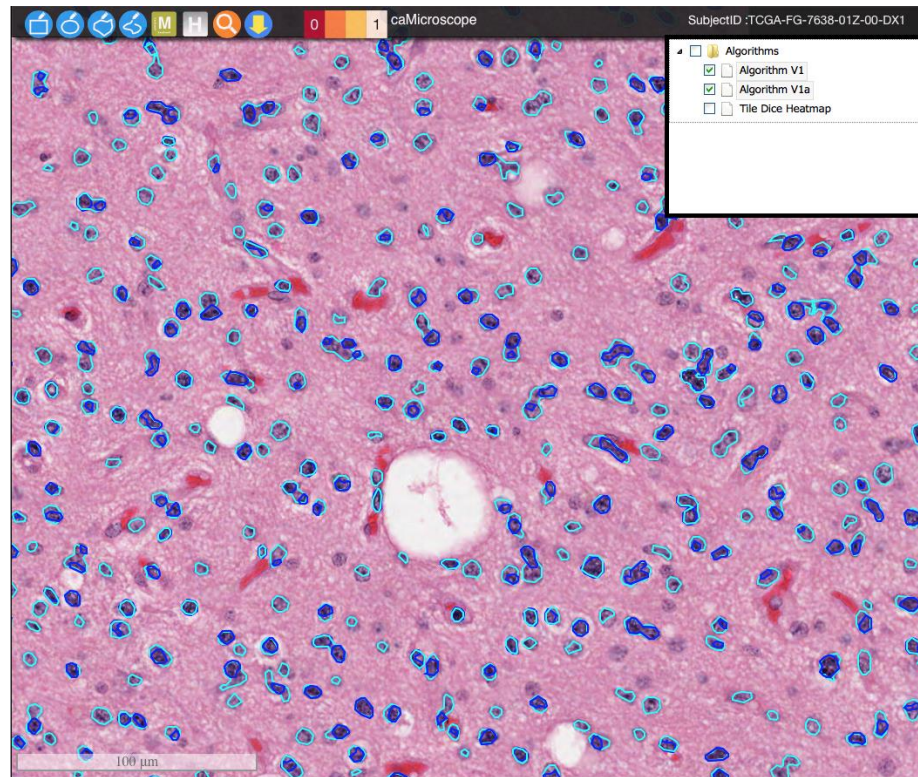- Disk storage is getting cheaper

# Big Data Challenges in Pathomics



Patients — 100's - 1000's

Regions — Millions

Nuclei — Billions+

# TCGA Microscopy Image Data

- Data from 11,000 subjects

- 30,000+ tissue slide images

- Image resolutions ranging from 4000x6000 pixels to 130,000x250,000 pixels
  - On average about 3.6 Billion pixels per image

# Big Data: Analyzing Whole Slide Tissue Images

- Needs large main memory
  - 130,000x250,000x3 (RGB) = 90GB
- Takes from 30 minutes to 10-12 hours to process an image
  - A dataset with 100 images would require 8 days assuming 2 hours per image on average
- Hundreds of thousands to millions of objects in a whole slide tissue image

# Analysis Sensitivity: Generating Robust Feature Sets

- Image analysis pipelines are sensitive to input parameters

# Generating Robust Feature Sets

- Run multiple analyses

- Store, index, interact with results

- Computational comparison of results
  - Sensitivity analysis for algorithm evaluation and development
  - Parameter tuning

- Visual comparison of results
  - Curation

# Feature Sets

- 40-70 features per object (nucleus)

- Analysis of 4000 images

  – About 2 Billion segmented objects

- Multi-analysis of 300+ images

  – 6-10 analyses per image

  – About 2 Billion segmented objects

# Deep Learning

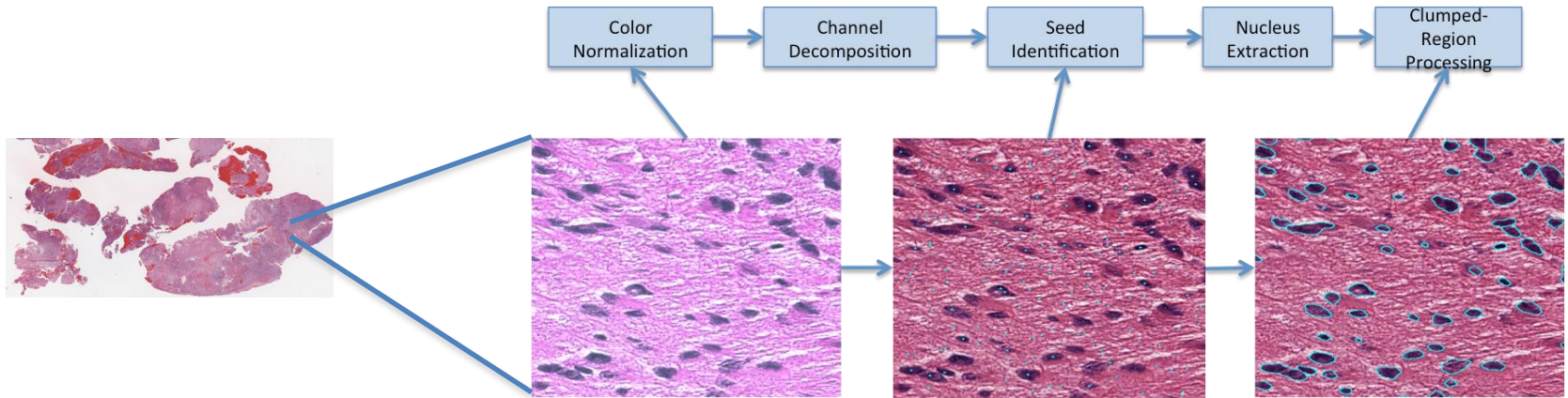- Tumor infiltrating lymphocytes

# TIL Maps

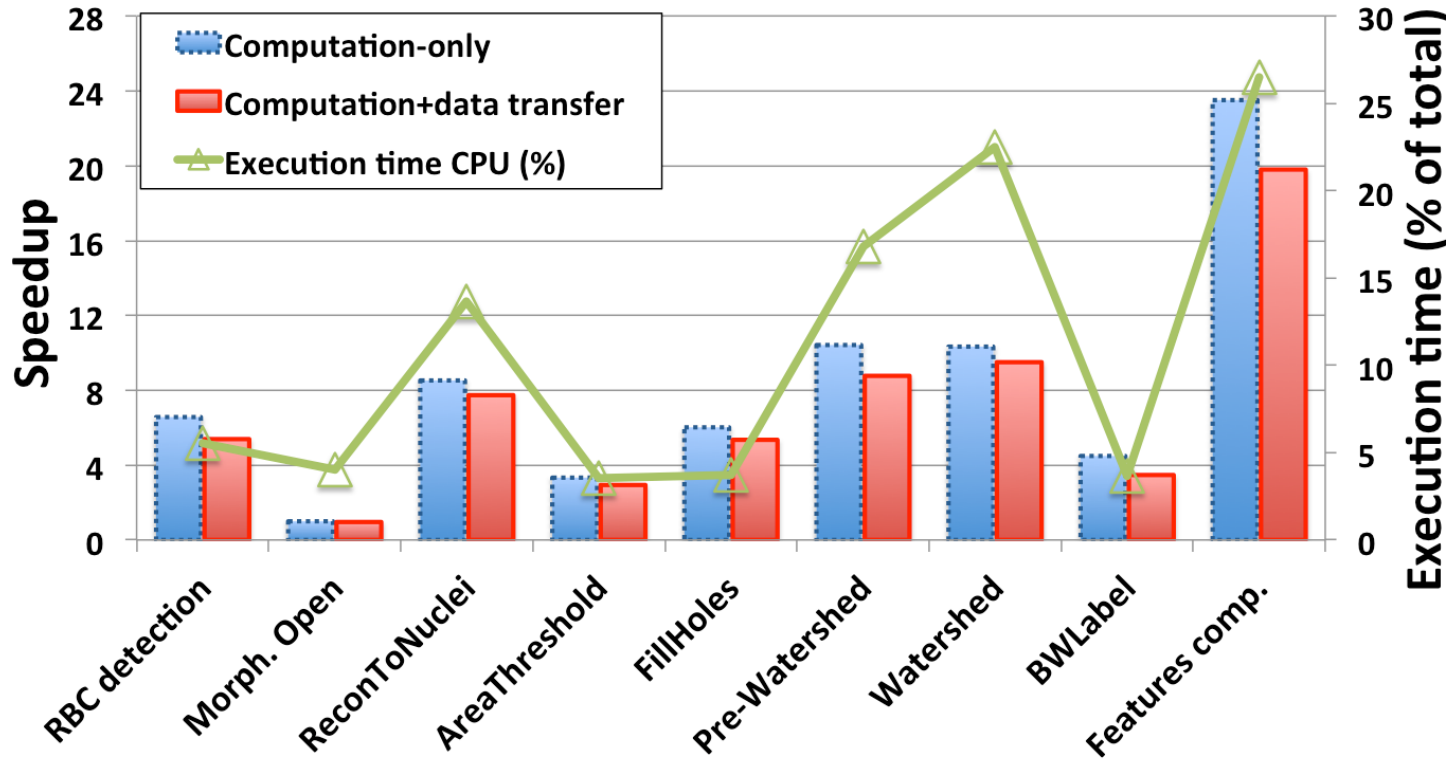# Pathomics Framework

# Leveraging Hybrid High Performance Computing Systems

- Lots of nodes

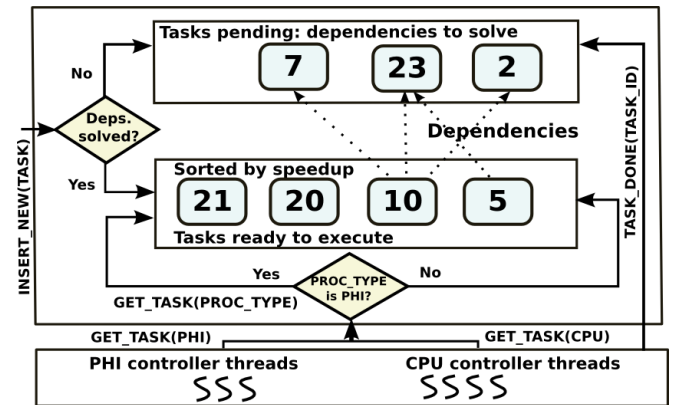- Multi-core CPUs, GPUs, and other hardware accelerators per node

- XSEDE, Supercomputing Centers
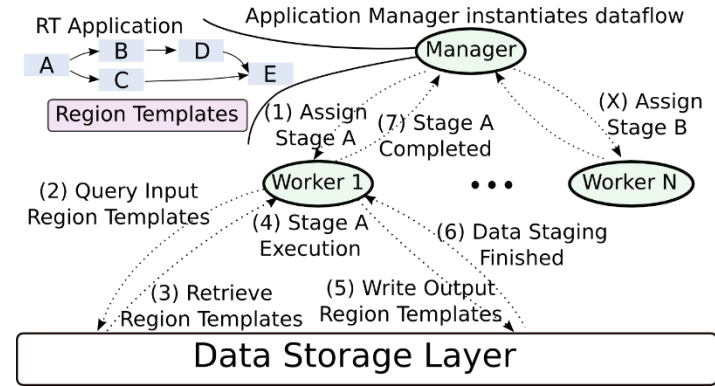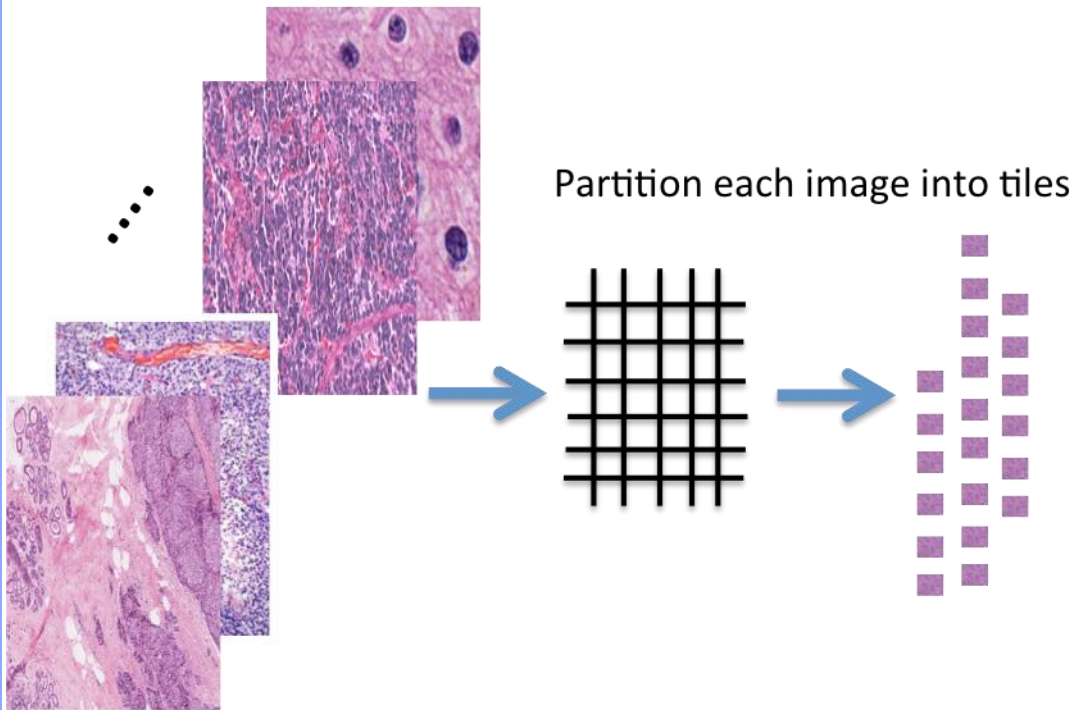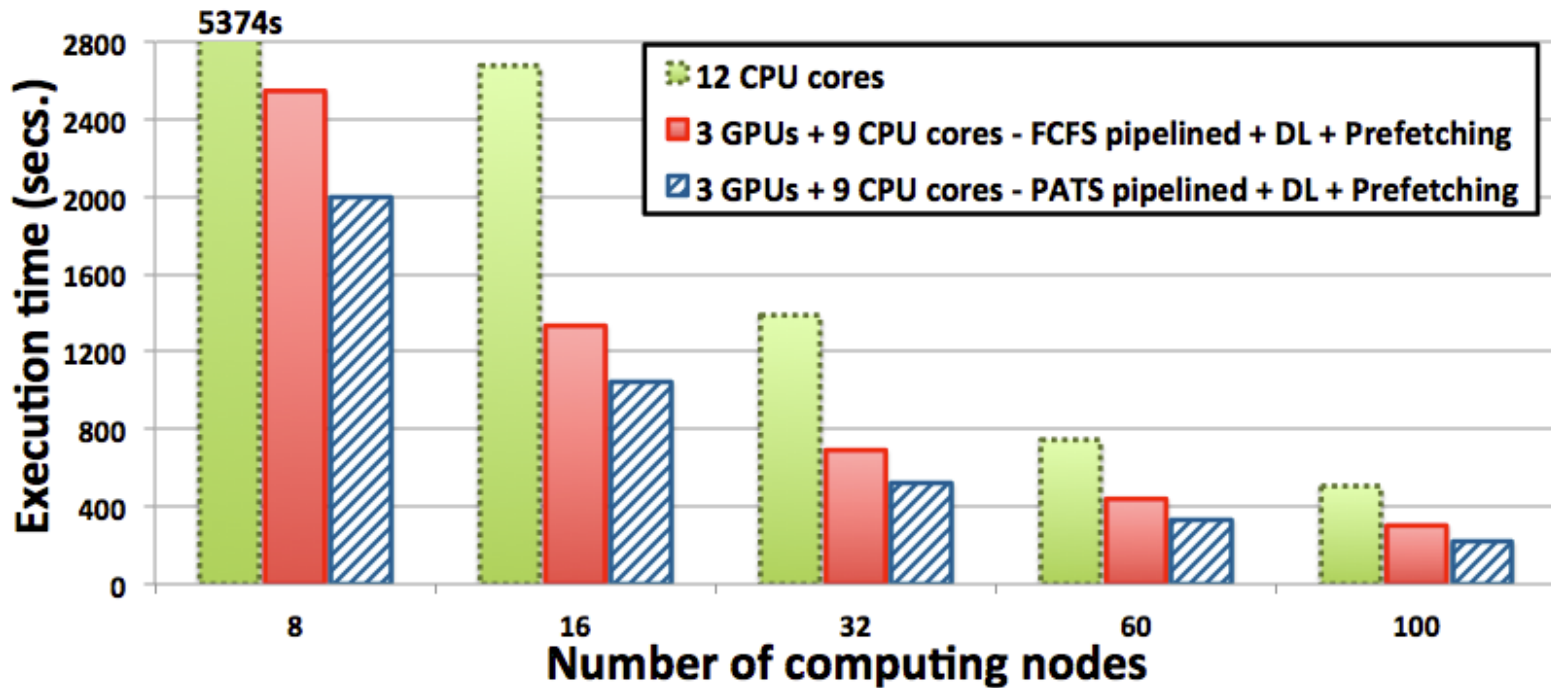
# Nuclear Imaging Features

# Inter-operation Performance Variations



Experiments on Keeneland: equipped w/ M2090 GPUs

# Analysis on HPC Platforms



Partition each image into tiles



- Thousands of CPU cores
- Hundreds of GPUs
- Large distributed memory

# Coordinated Use of CPUs and GPUs



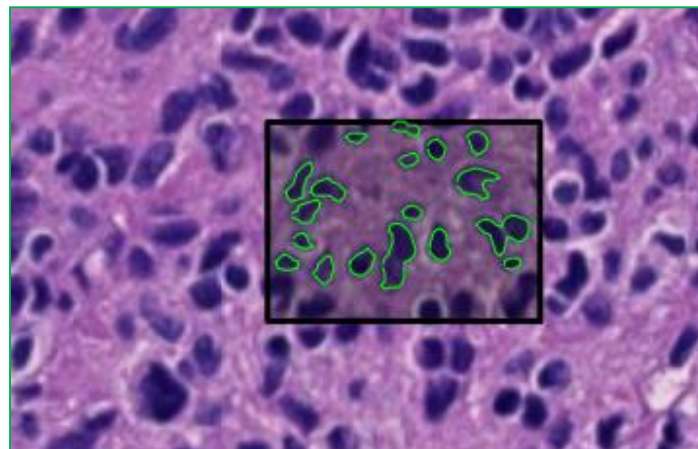36,844 4Kx4K-pixel tiles from 340 whole slide tissue images.

On 100 nodes, less than 4 minutes to process 36K tiles.
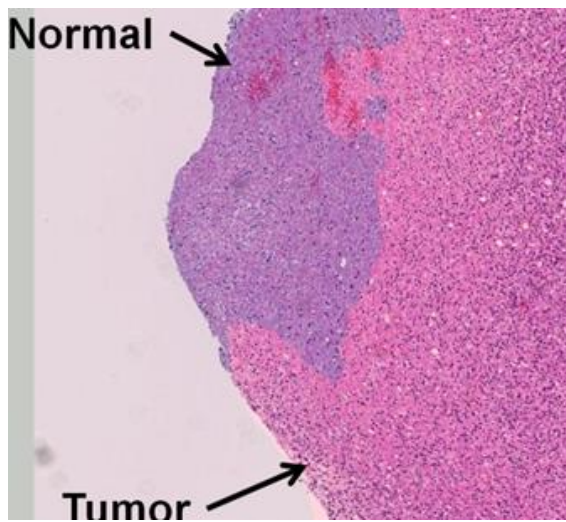
# Support for Spatial Queries
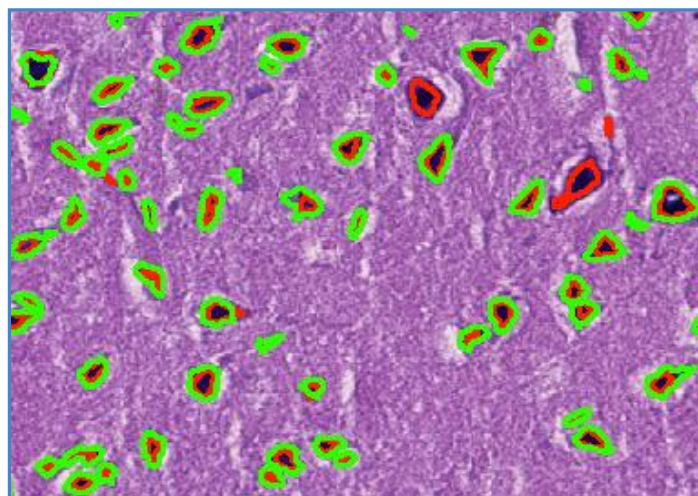
**POINT** query: human marked point inside a nucleus



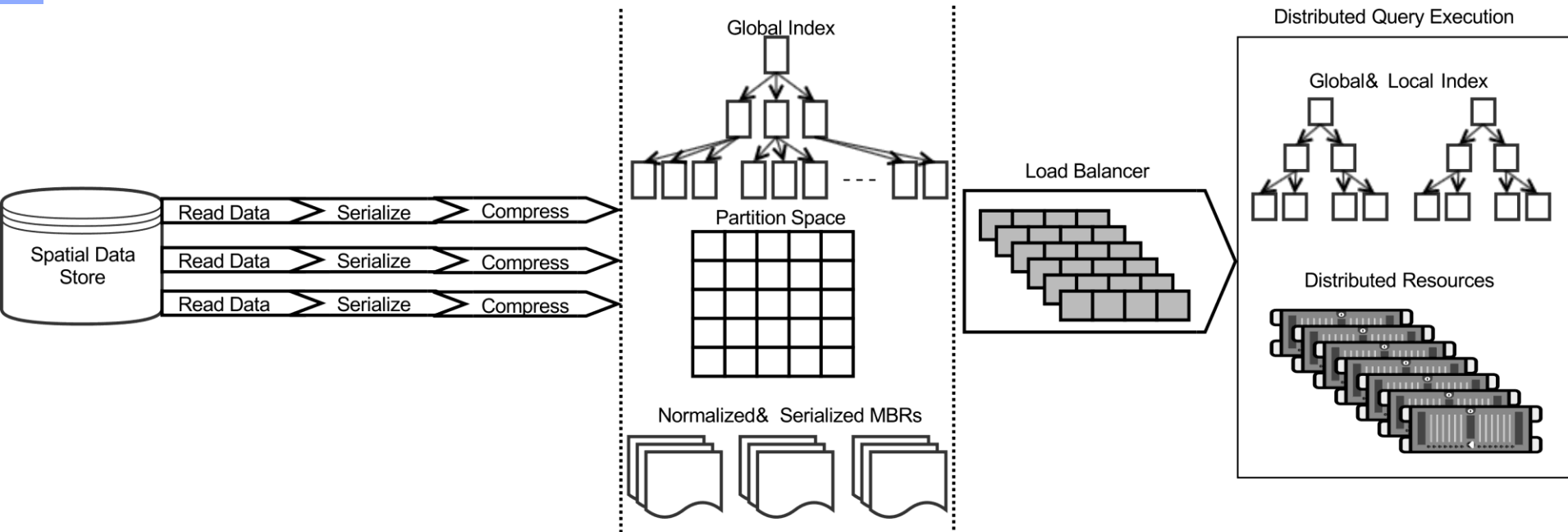**WINDOW** query: return markups contained in a rectangle



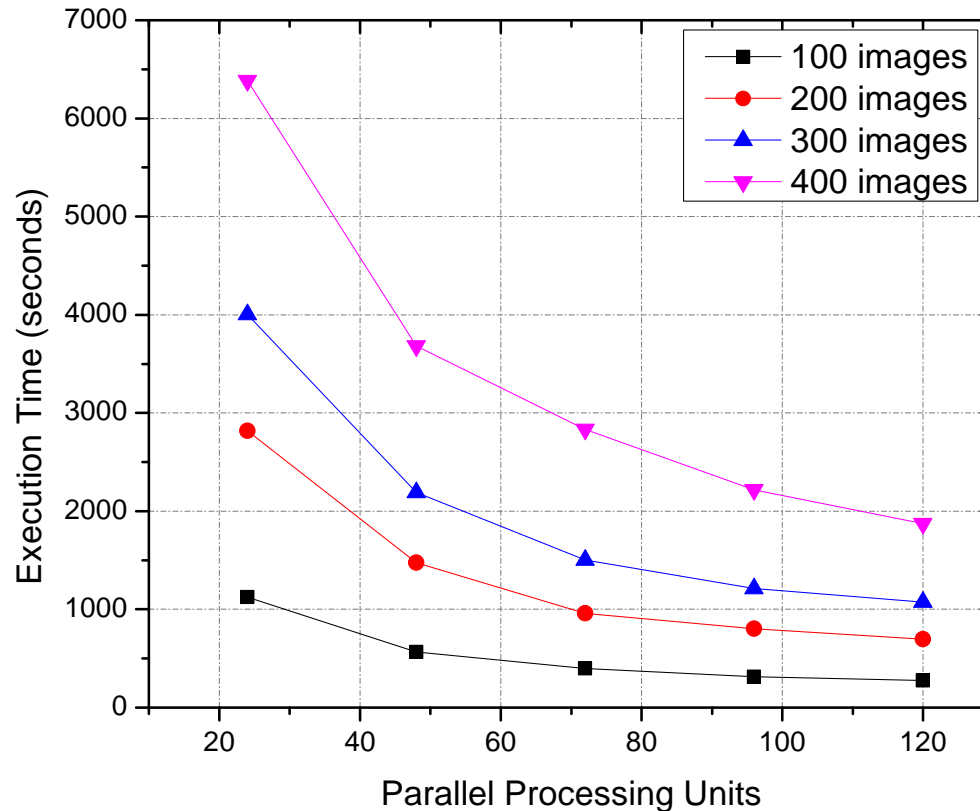**CONTAINMENT** query: nuclear feature aggregation in tumor regions



**SPATIAL JOIN** query: algorithm validation/comparison

# SparkGIS
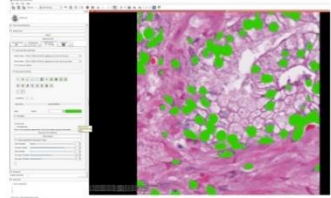
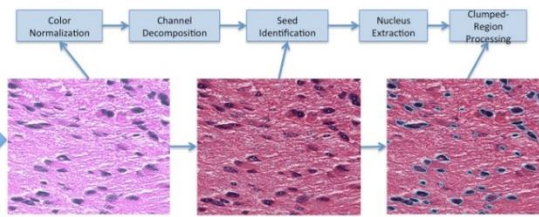# Heatmap Computations: Spatial Joins + Dice/Jaccard Metric

# Going Forward

- Containers
  - Facilitates a modular design.
  - Self-contained, Isolated.
    - Software to support tissue image analysis needs to leverage a variety of existing libraries and tools.
  - Flexibility
    - Move computation to data
    - Move back and forth between Cloud platforms and local resources
- Scripting
  - JavaScript: Take advantage of web browsers
  - Python: Take advantage of large set of libraries

# Pathomics Framework

# Containerized Software for Tissue Image Analysis

# Containerized Software for Tissue Image Analysis

- ***The application service group*** is a single container that hosts a suite of Web applications to view images and interact with analysis results.

- ***The image analysis group*** is made up of three containers, which collectively execute image analysis requests.
  - Analysis service – hosts analysis pipeline
  - Job Manager service – tracks jobs
  - Image tile service – services image tiles for analysis

- ***The data management service group*** is implemented as a set of three containers for data loading, data management, and query processing.
  - Data loader service  -- load image metadata and analysis results
  - Data manager service – manage and index image data, analysis results, features
  - Feature query service – query feature data for visualization and exploration

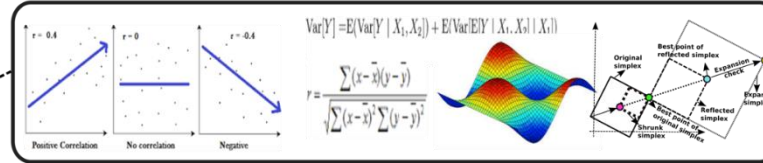Software: https://github.com/SBU-BMI/quip_distro.git

# Thank you.

# Support for Sensitivity Analysis