



# Pathomics Based Biomarkers, Tools and Methods

Joel Saltz

Department of Biomedical Informatics

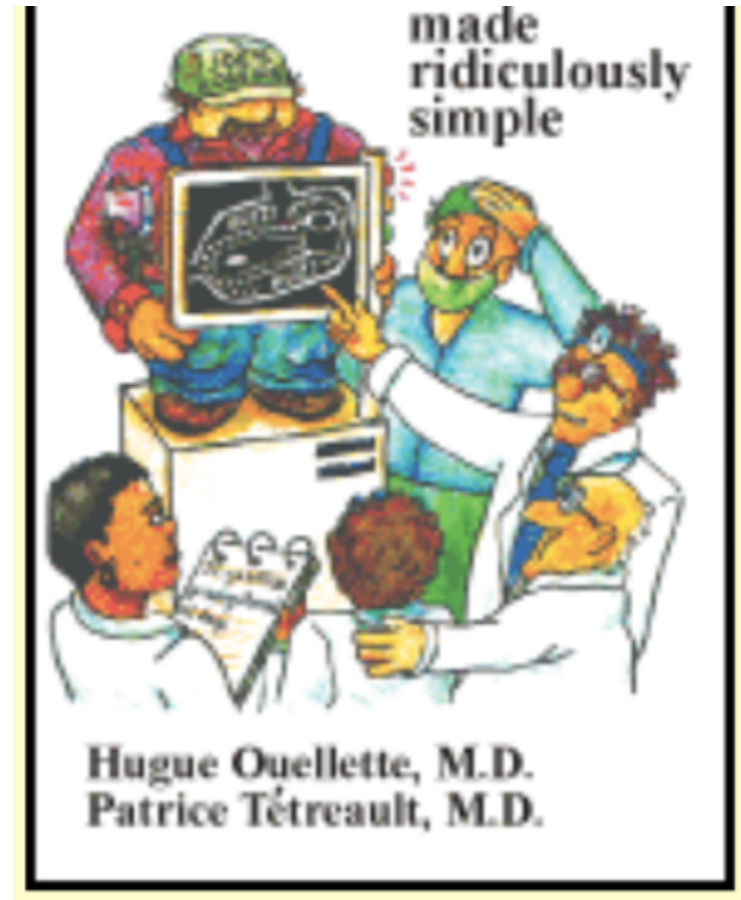
Stony Brook University

Imaging Community Call

October 24, 2016

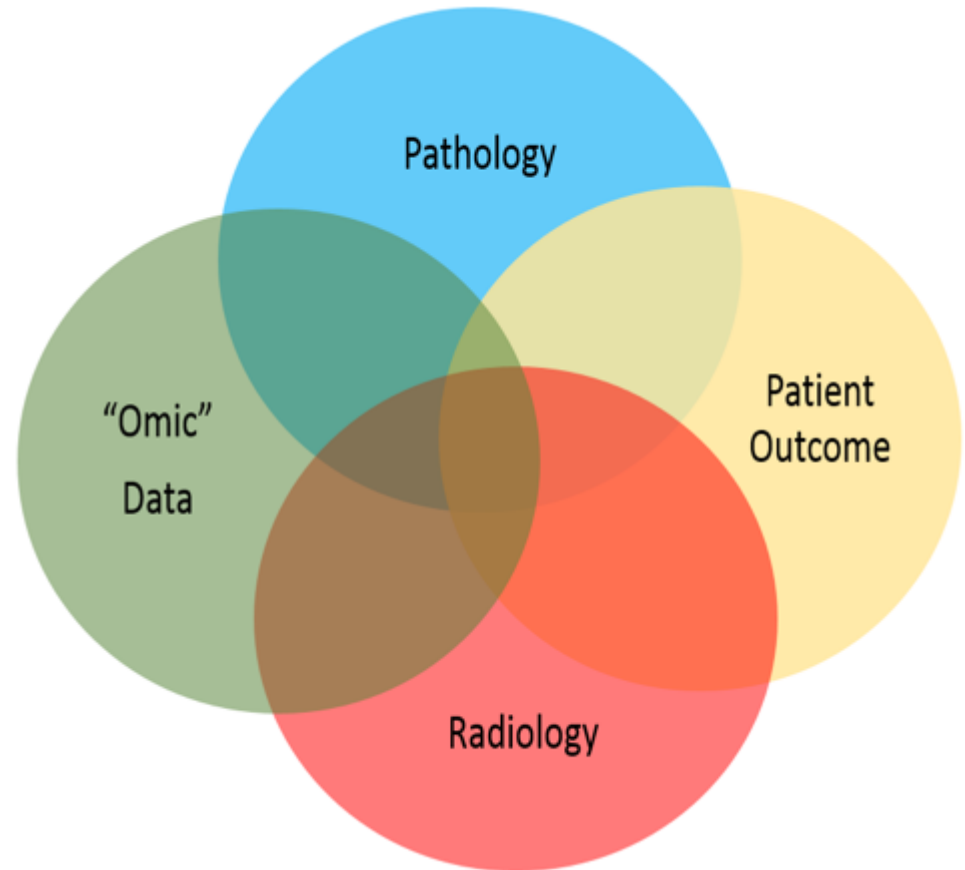


## Multi-Scale Precision Medicine



# Multi-scale Integrative Analysis in Biomedical Informatics

- Predict treatment outcome, select, monitor treatments
- Reduce inter-observer variability in diagnosis
- Computer assisted exploration of new classification schemes
- Multi-scale cancer simulations



# Pathomics, Radiomics

Identify and segment trillions of objects – nuclei, glands, ducts, nodules, tumor niches ... from Pathology, Radiology imaging datasets

Extract features from objects and spatio-temporal regions

Support queries against ensembles of features extracted from multiple datasets

Statistical analyses and machine learning to link Radiology/Pathology features to “omics” and outcome biological phenomena

Principle based analyses to bridge spatio-temporal scales – linked Pathology, Radiology studies

# Radiomics

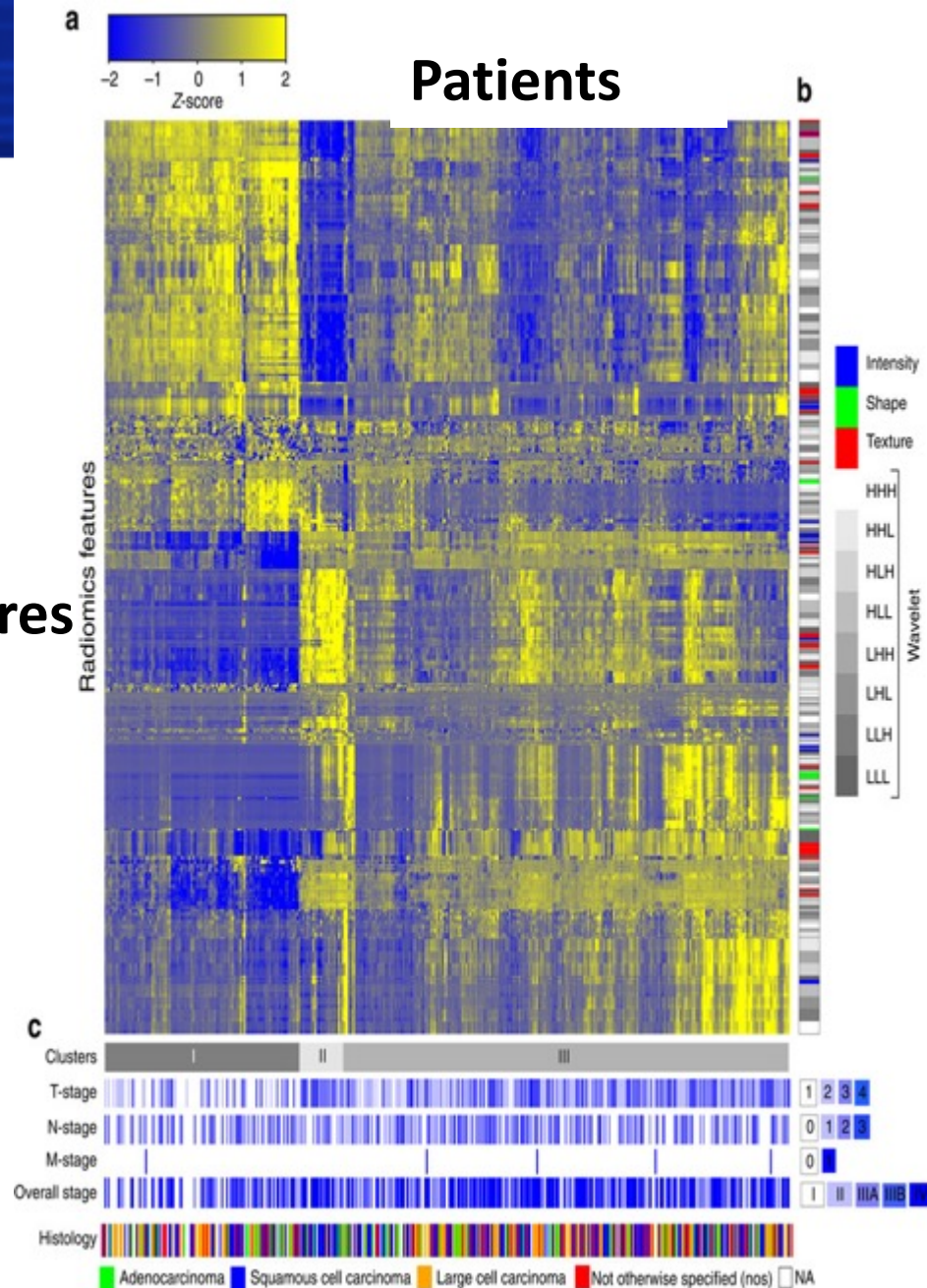
## Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach

Hugo J. W. L. Aerts et. Al.

*Nature Communications* 5, Article number: 4006

doi:10.1038/ncomms5006

Features



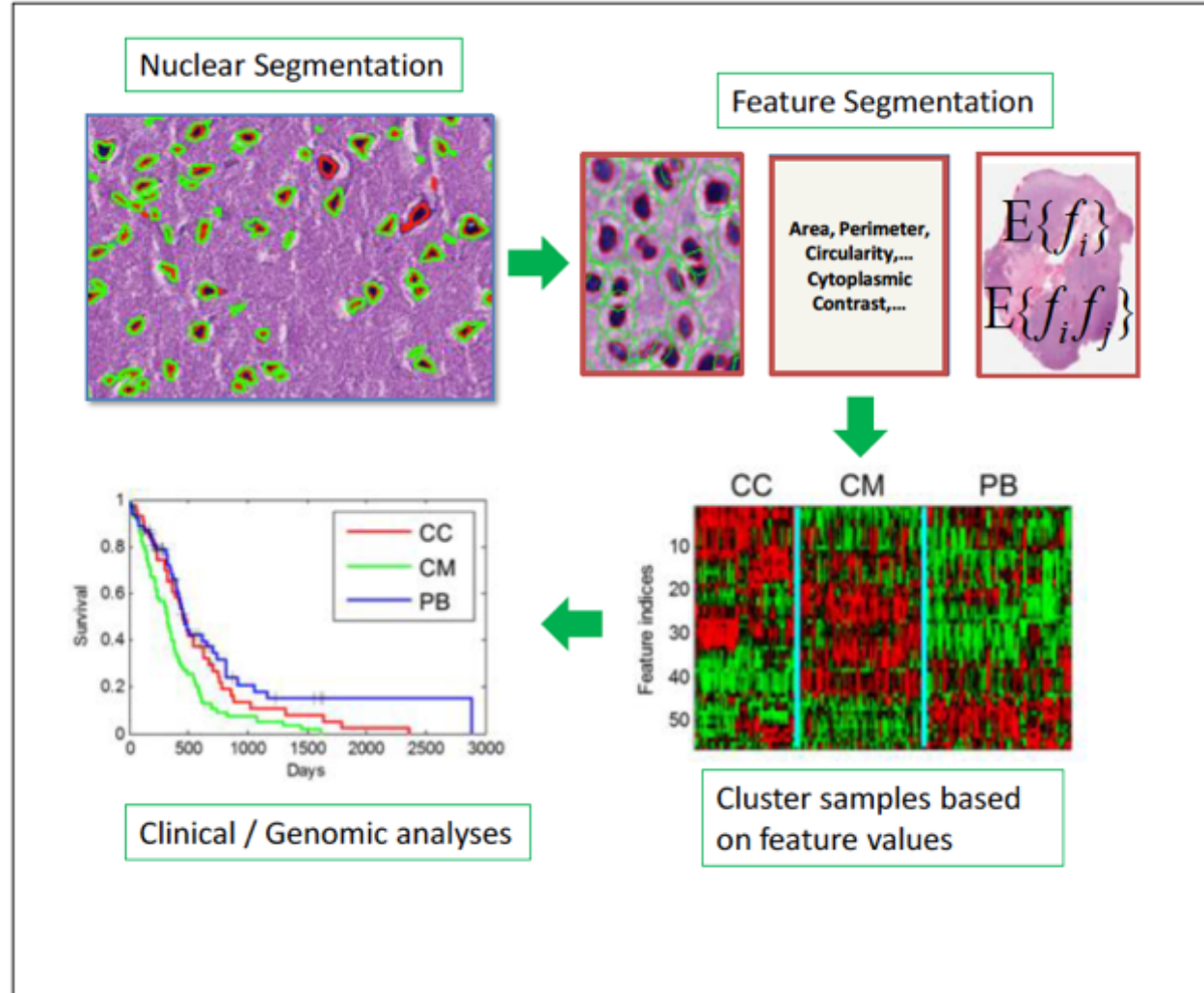
# Pathomics

## Integrative Morphology/"omics"

Quantitative Feature Analysis in Pathology: Emory In Silico Center for Brain Tumor Research (PI = Dan Brat, PD= Joel Saltz)

NLM/NCI: Integrative Analysis/Digital Pathology R01LM011119, R01LM009239 (Dual PIs Joel Saltz, David Foran)

J Am Med Inform Assoc. 2012 **Integrated morphologic analysis for the identification and characterization of disease subtypes.**



Lee Cooper, Jun Kong

- **Specific Aim 1** Analysis **pipelines** for multi- scale, integrative image analysis.
- **Specific Aim 2: Database** infrastructure to manage and query Pathomics features.
- **Specific Aim 3:** HPC software that **targets clusters, cloud computing, and leadership scale systems.**
- **Specific Aim 4:** Develop **visualization** middleware to relate Pathomics feature and image data and to integrate Pathomics image and “omic” data.

# SEER Virtual Tissue Repository

- Lynne Penberthy MD, MPH NCI SEER
- Ed Helton PhD NCI CBIIT Clinical Imaging Program
- Ulrike Wagner CBIIT Clinical Imaging Program
- Radim Moravec NCI PhD, NCI SEER
- Ashish Sharma PhD Biomedical Informatics Emory
- Joel Saltz MD, PhD Biomedical Informatics Stony Brook
- Tahsin Kurc PhD Biomedical Informatics Stony Brook

***Vision – Enable population/epidemiological cancer research that leverages rich cancer phenotype information available from Pathology tissue studies***

***NCIP/Leidos 14X138 and HHSN261200800001E - NCI***



# SEER Virtual Tissue Repository

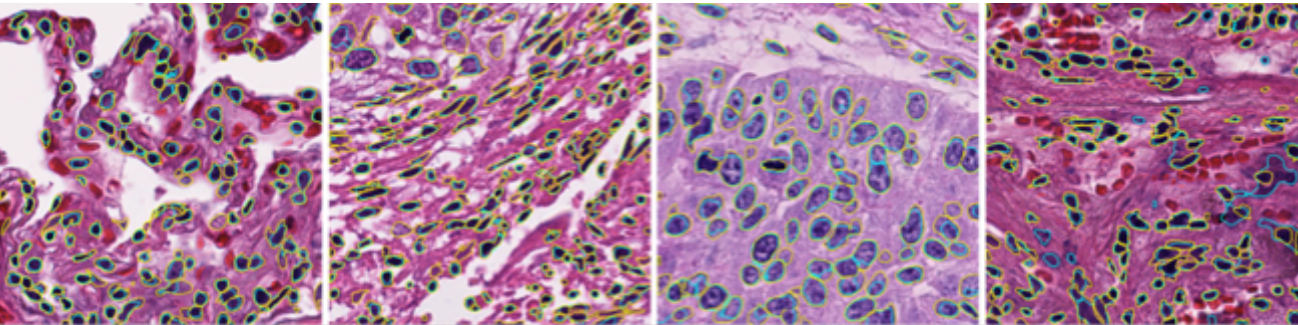
- SEER registries are a potential source of information about unusual outcomes and rare cancers
- Leverage Pathology labs which store FFPE tumors, slides and digital images
- Link to SEER data – track long term outcomes
- Accrue linked clinical data, Pathology slides from SEER sites

# SEER VIRTUAL TISSUE REPOSITORY

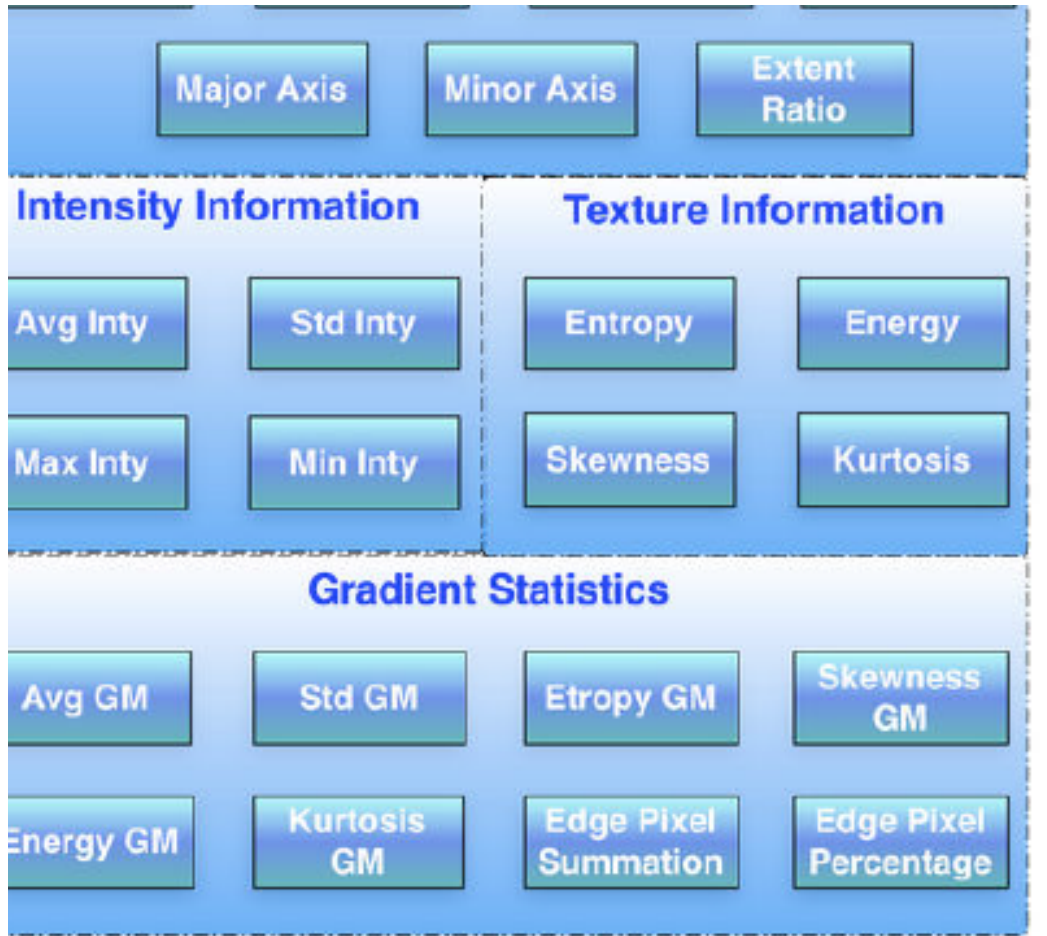
- Create linked collection of de-identified clinical data and whole slide images
- Extract features from a sample set of images (pancreas and breast cancer).
- Enable search, analysis, epidemiological characterization
- Pilot focus on extreme outcome Breast Cancer, Pancreatic Cancer cases
- Display images and analyzed features

# Robust Nuclear Segmentation

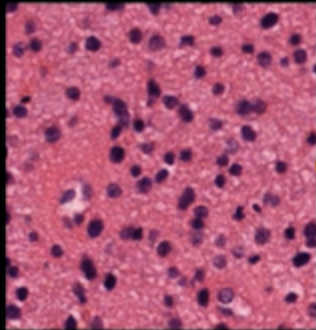
- Robust ensemble algorithm to segment nuclei across tissue types
- Optimized algorithm tuning methods
- Parameter exploration to optimize quality
- Systematic Quality Control pipeline encompassing tissue image quality, human generated ground truth, convolutional neural network critique
- Yi Gao, Allen Tannenbaum, Dimitris Samaras, Le Hou, Tahsin Kurc



# Cell Morphometry Features



## Whole Slide Images (WSI)



Segmentation  
Parameters

## Compute Cluster



Process the images for subjects  
selected

Compute object-level (nucleus-  
level) image features

Compute aggregated patient-  
level image features from  
object-level features

## FeatureDB

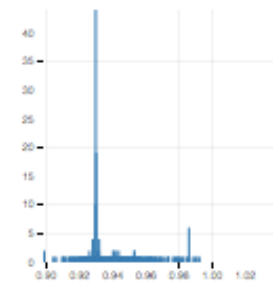
- Load object-level imaging features and segmentation results
- Load patient-level imaging features along with a selected subset of clinical and genomic data (e.g. gene mutations, days to death, vital status)

# Feature Explorer Suite

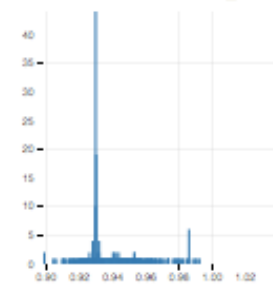
- Explore Relationship Between Imaging Features, Outcome, "omics"
- Explore relationships between features and explore how features relate to images

# Feature Explorer - Integrated Pathomics Features, Outcomes and “omics” – TCGA NSCLC Adeno Carcinoma Patients

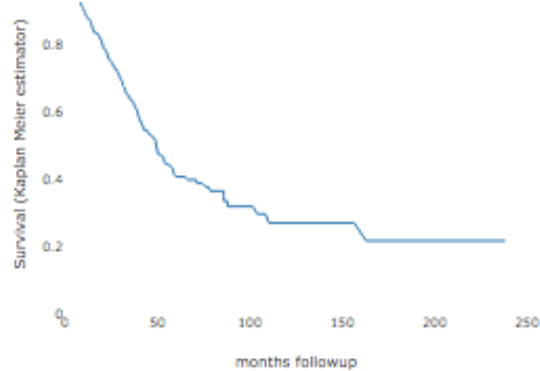
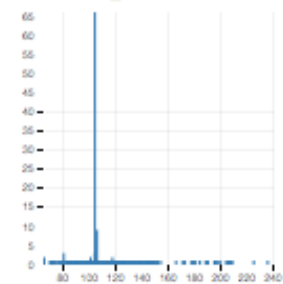
Var 1: Roundness\_median



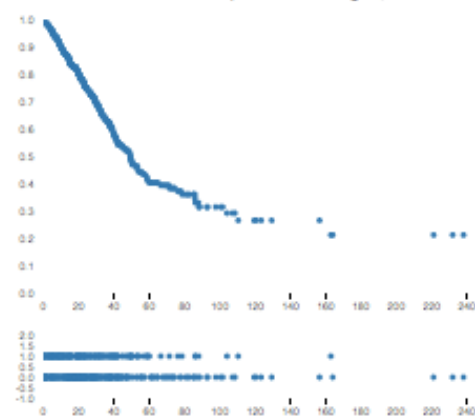
Var 1 Zoom: Roundness\_median



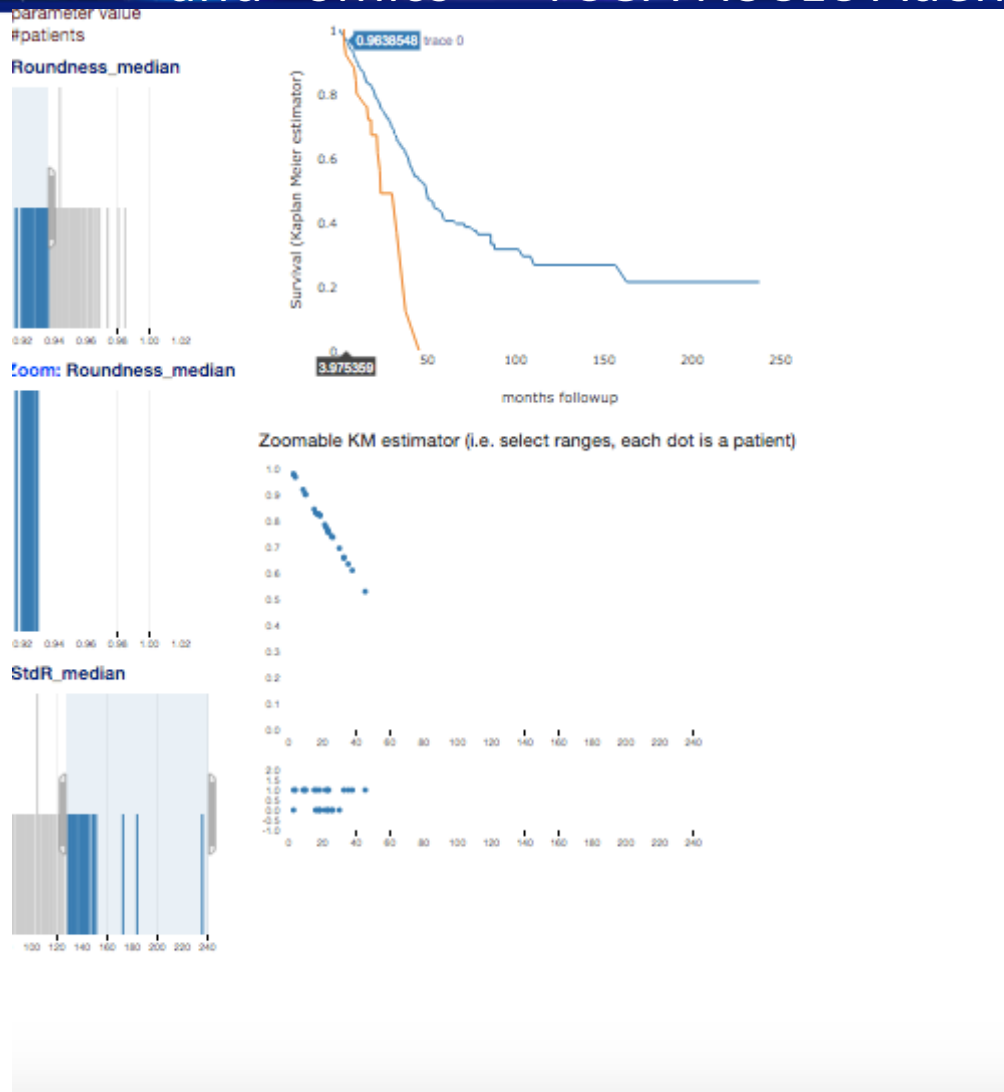
Var 2: StdR\_median



Zoomable KM estimator (i.e. select ranges, each dot is a patient)

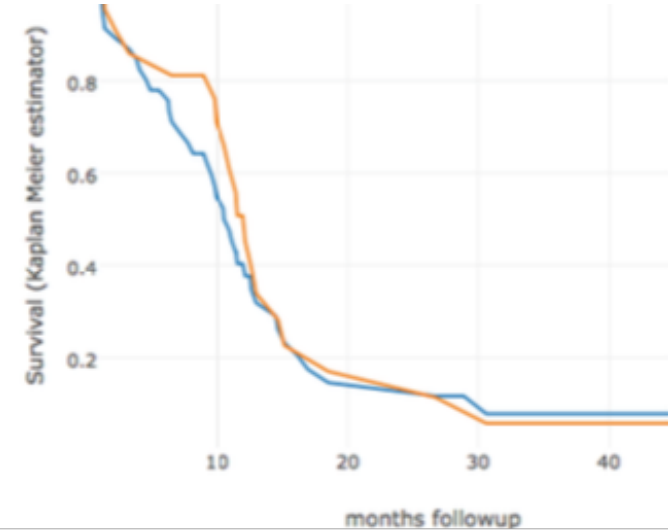
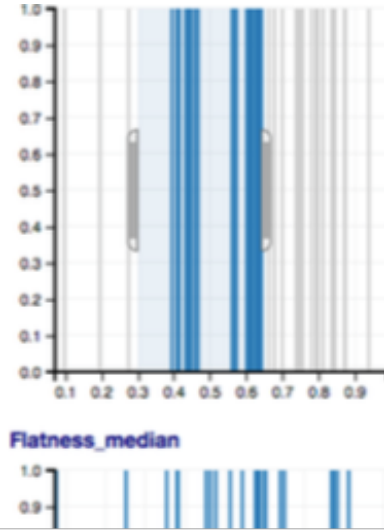
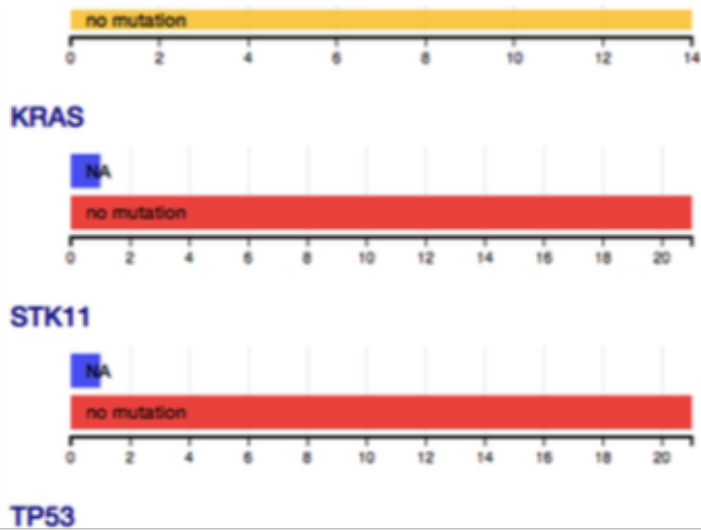


# Feature Explorer - Integrated Pathomics Features, Outcomes and "omics" – TCGA NSCLC Adeno Carcinoma Patients

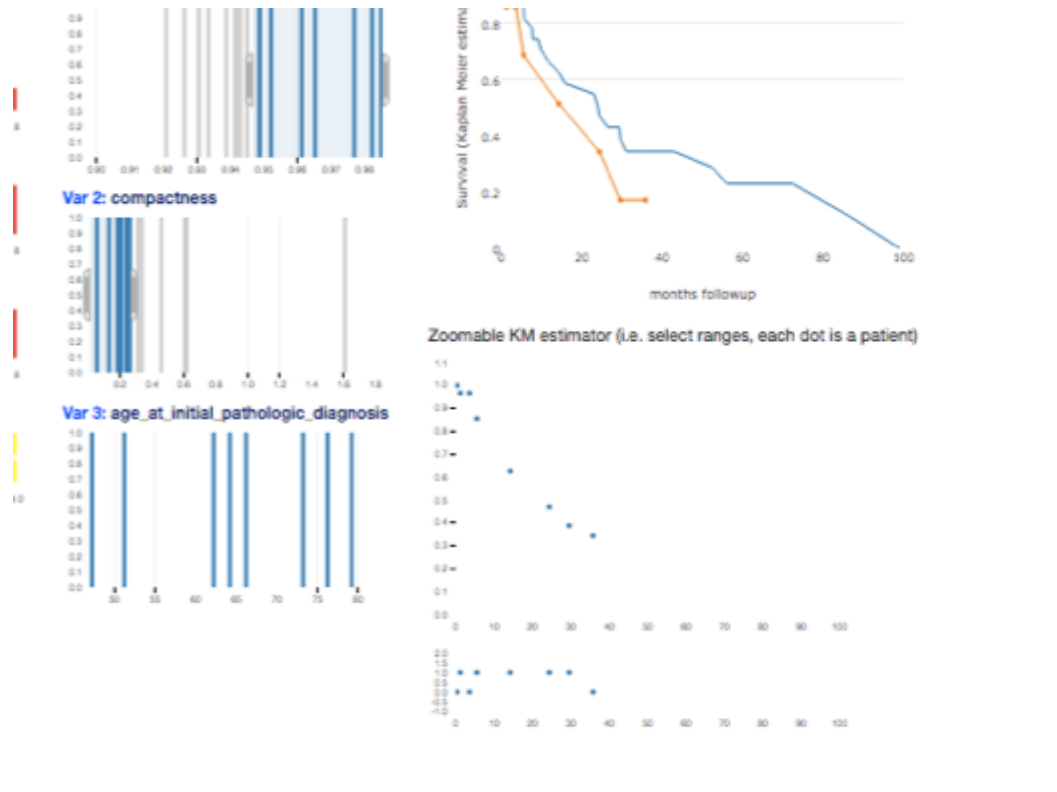




# Collaboration with MGH – Feature Explorer – Radiology Brain MR/Pathology Features



# Collaboration with SBU Radiology – TCGA NSCLC Adeno Carcinoma Integrative Radiology, Pathology, “omics”, outcome




Mary Saltz, Mark Schweitzer SBU Radiology

# Pathomics

## Relationship Between Image and Features

**Step 1:** Choose a case from the TCGA atlas (case #20)

**FeatureScope: u24 case preview prototype** 

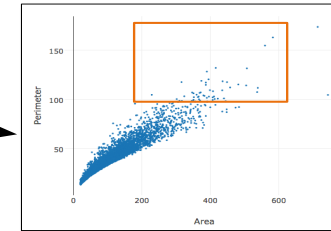
Preview of using `featurescope` to explore different patient case ids.  
For an interactive visualization of Pathology results (including links to the cBio record) see `TCGAScope`.

---

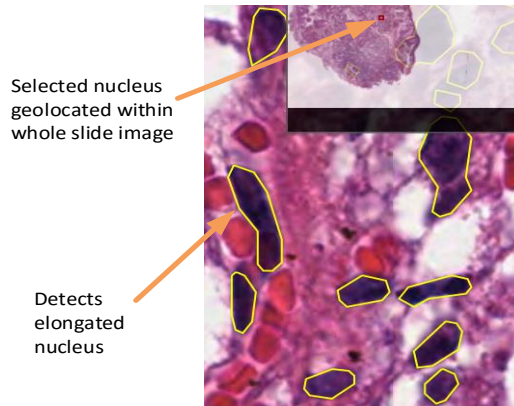
**preview Case IDs**

20\_TCGA-55-7574-01Z-00-DX1 (cbio) random\_seed:0.816, feature\_sample\_size:1000 | [featurescope of sampled features](#)

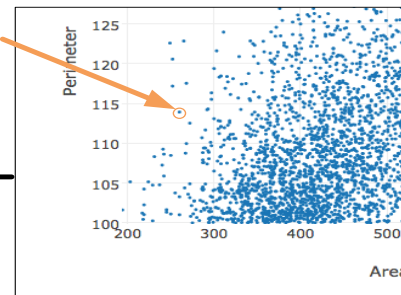
**Step 2:** Select two features of interest; x axis (*area*), Y axis (*perimeter*)



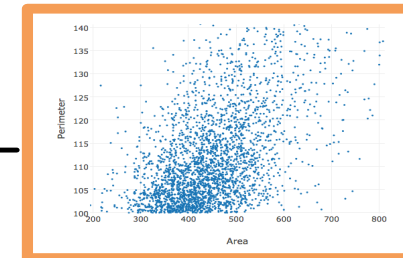
**Step 5:** Evaluate the features selected in the context of the specific nucleus and where this nucleus is located within the whole slide image



**Step 4:** Pick a specific nucleus of interest. Each dot represents a single nucleus



**Step 3:** Zoom in on region of interest



The tool provides visual context for feature evaluation. This technique maps both intuitive features (i.e. size, shape, color) and non-intuitive features (i.e. wavelets, texture) to the ground truth of source images through an interactive web-based user interface.

# FeatureScape: u24 case preview prototype

Preview of using [featurescape](#) to explore different patient case ids.

For an interactive visualization of Pathology results (including links to the cBio record) see [TCGAScope](#).

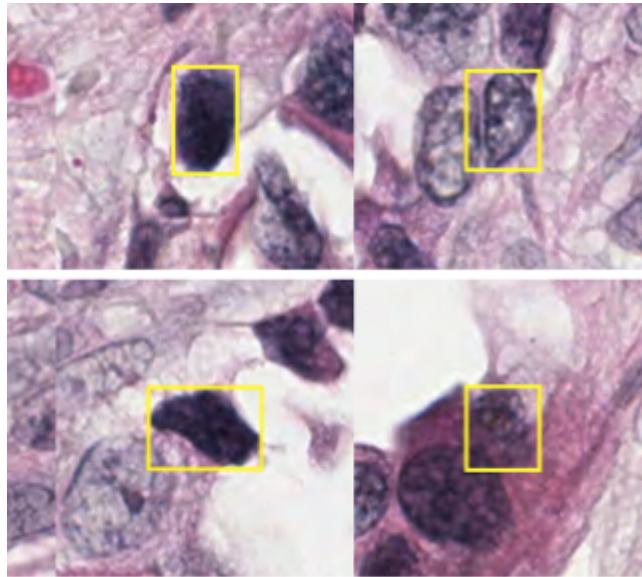
## preview Case IDs

1. TCGA-34-2605-01Z-00-DX1, (cbio) random seed: <b>0.559</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
2. TCGA-38-4625-01Z-00-DX1, (cbio) random seed: <b>0.628</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
3. TCGA-38-4626-01Z-00-DX1, (cbio) random seed: <b>0.700</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
4. TCGA-38-4628-01Z-00-DX1, (cbio) random seed: <b>0.016</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
5. TCGA-38-4629-01Z-00-DX1, (cbio) random seed: <b>0.185</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
6. TCGA-38-6178-01Z-00-DX1, (cbio) random seed: <b>0.317</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
7. TCGA-38-A44F-01Z-00-DX1, (cbio) random seed: <b>0.906</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
8. TCGA-50-5044-01Z-00-DX1, (cbio) random seed: <b>0.055</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
9. TCGA-50-5045-01Z-00-DX1, (cbio) random seed: <b>0.946</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
10. TCGA-50-5045-01Z-00-DX2, (cbio) random seed: <b>0.551</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
11. TCGA-50-5055-01Z-00-DX1, (cbio) random seed: <b>0.127</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
12. TCGA-34-5232-01Z-00-DX1, (cbio) random seed: <b>0.208</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
13. TCGA-50-5055-01Z-00-DX2, (cbio) random seed: <b>0.321</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
14. TCGA-50-5066-01Z-00-DX1, (cbio) random seed: <b>0.711</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
15. TCGA-50-5066-02Z-00-DX1, (cbio) random seed: <b>0.008</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
16. TCGA-50-5942-01Z-00-DX1, (cbio) random seed: <b>0.031</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
17. TCGA-50-5946-01Z-00-DX1, (cbio) random seed: <b>0.768</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
18. TCGA-50-6590-01Z-00-DX1, (cbio) random seed: <b>0.668</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>
19. TCGA-50-6591-01Z-00-DX1, (cbio) random seed: <b>0.498</b> , feature sample size: 1000	<a href="#">featurescape of sampled features</a>

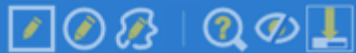




# Sample Nuclei from Gated Region

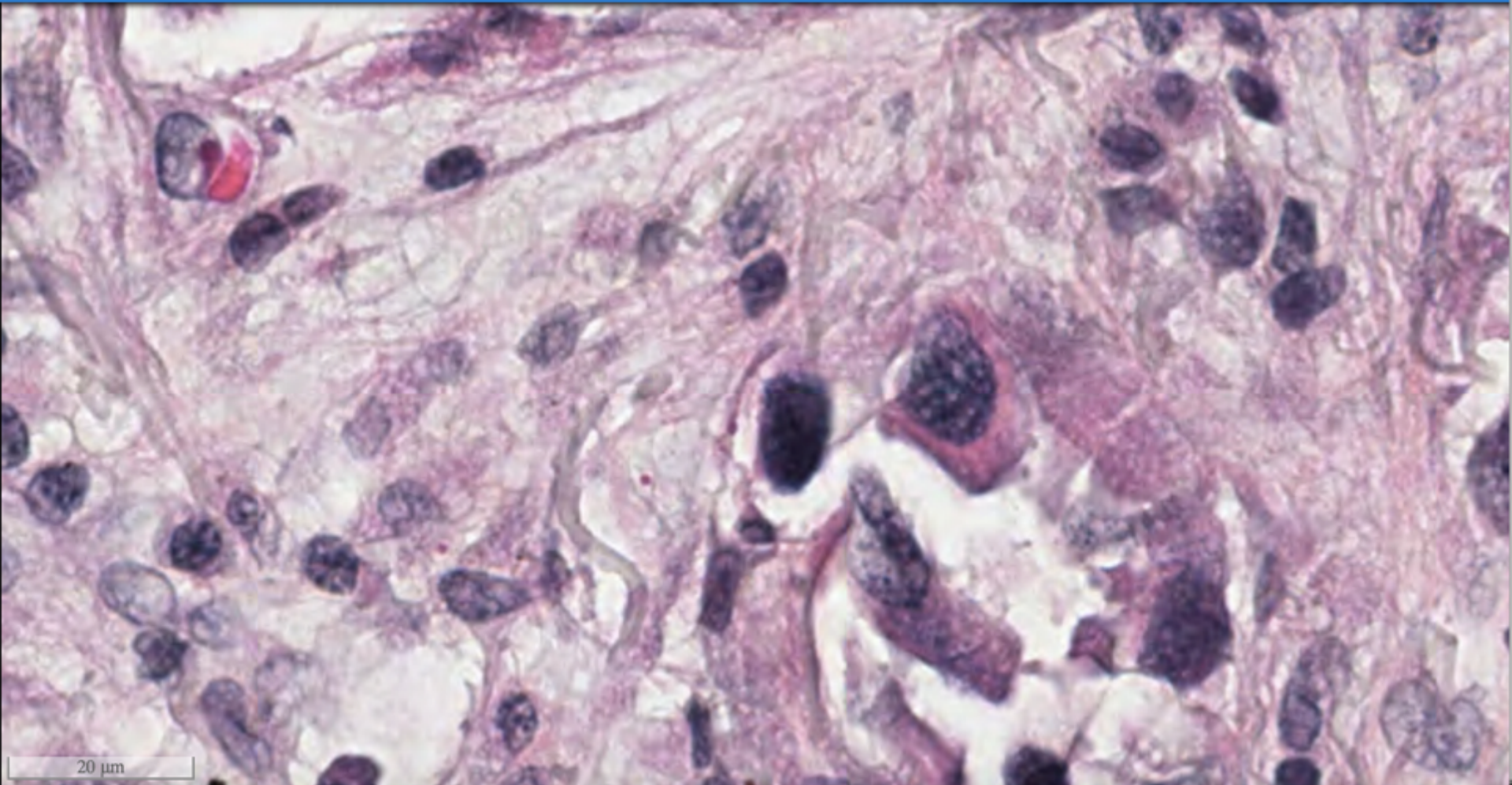


# Gated Nuclei in Context



caMicroscope

SubjectID :TCGA-38-4628-01Z-00-DX1

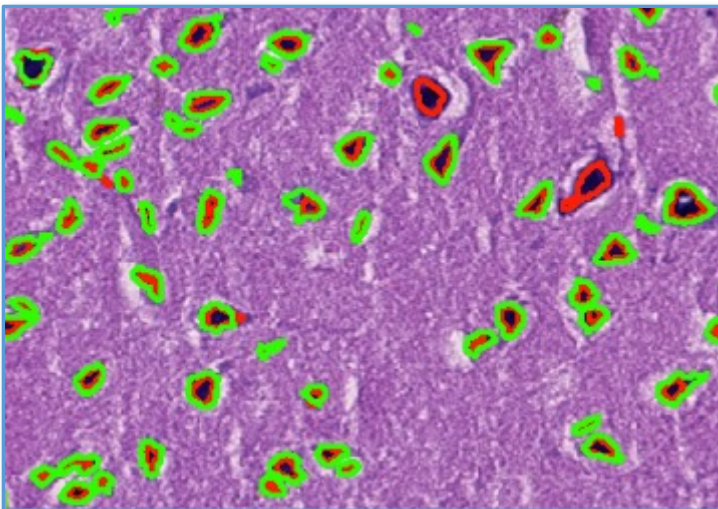


20  $\mu$ m



- High quality image analysis algorithms are essential to support biomedical research and diagnosis
  - Validate algorithms with human annotations
  - Compare and consolidate different algorithm results

e.g.: what are the distances and overlap ratios between markup boundaries from two algorithms?



Cross matching of two spatial data sets

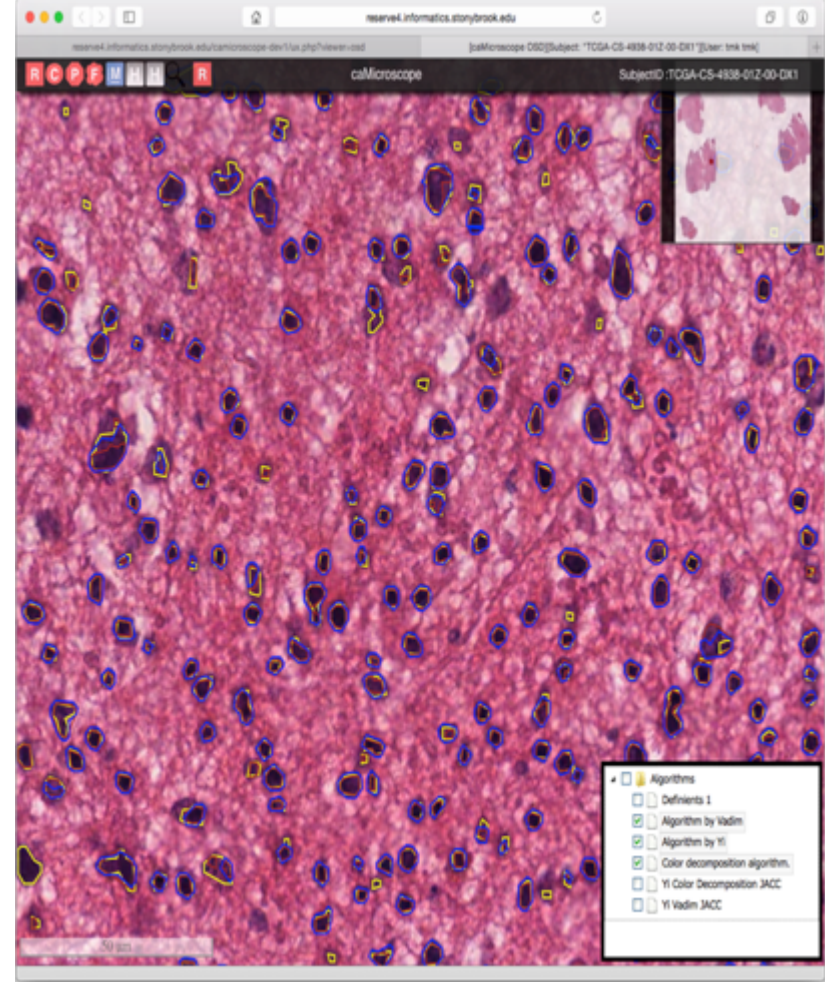
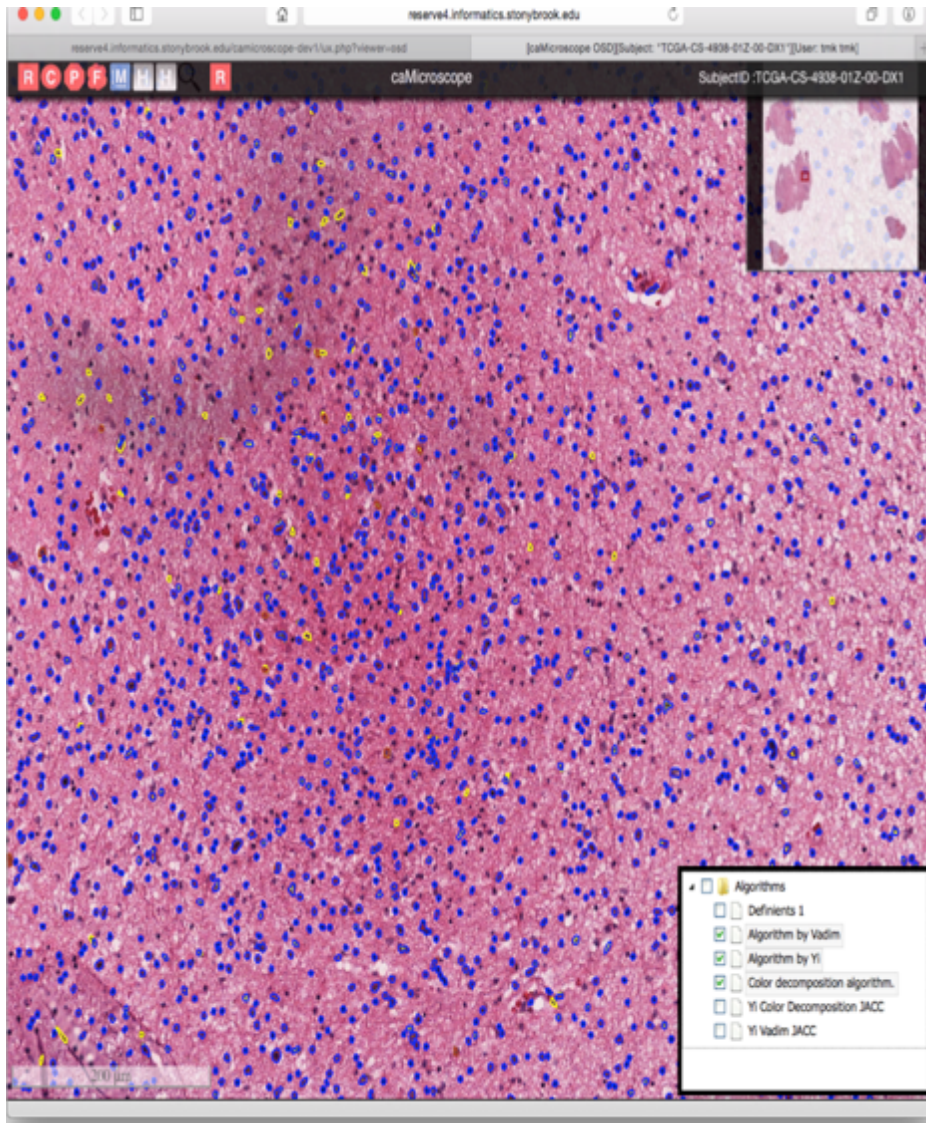
Green: algorithm 1

Red: algorithm 2

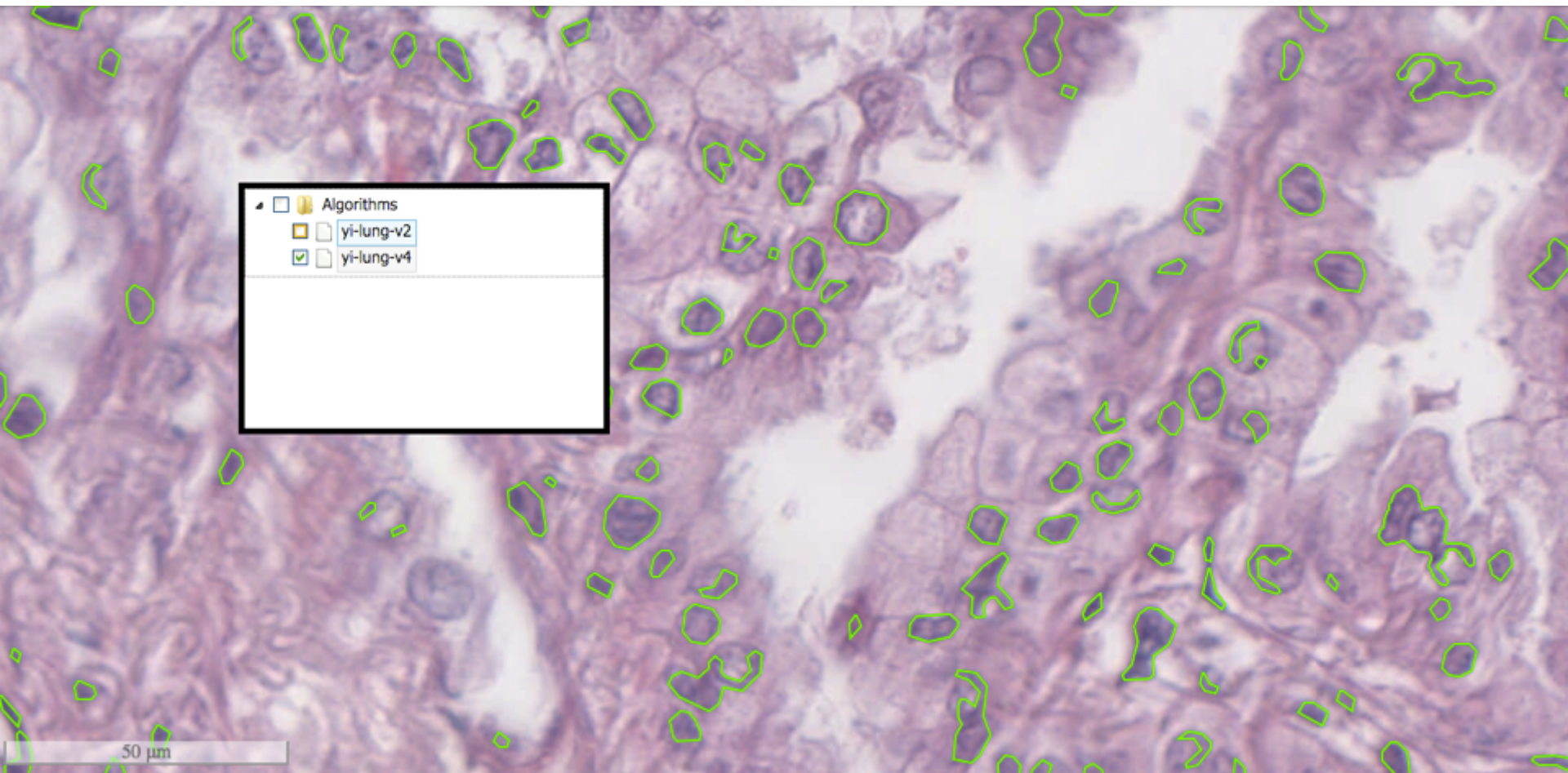


MICCAI 2014  
Brain Tumor

# caMicroscope/MongoDB - Multiple Algorithm Comparison; Generate and Curate Pathomics Feature set



# Compare Algorithm Results



# Heatmap – Depicts Agreement Between Algorithms

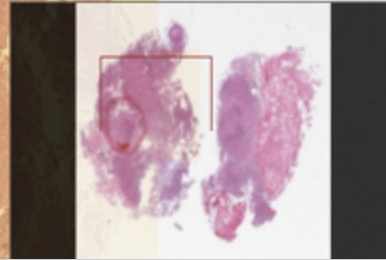
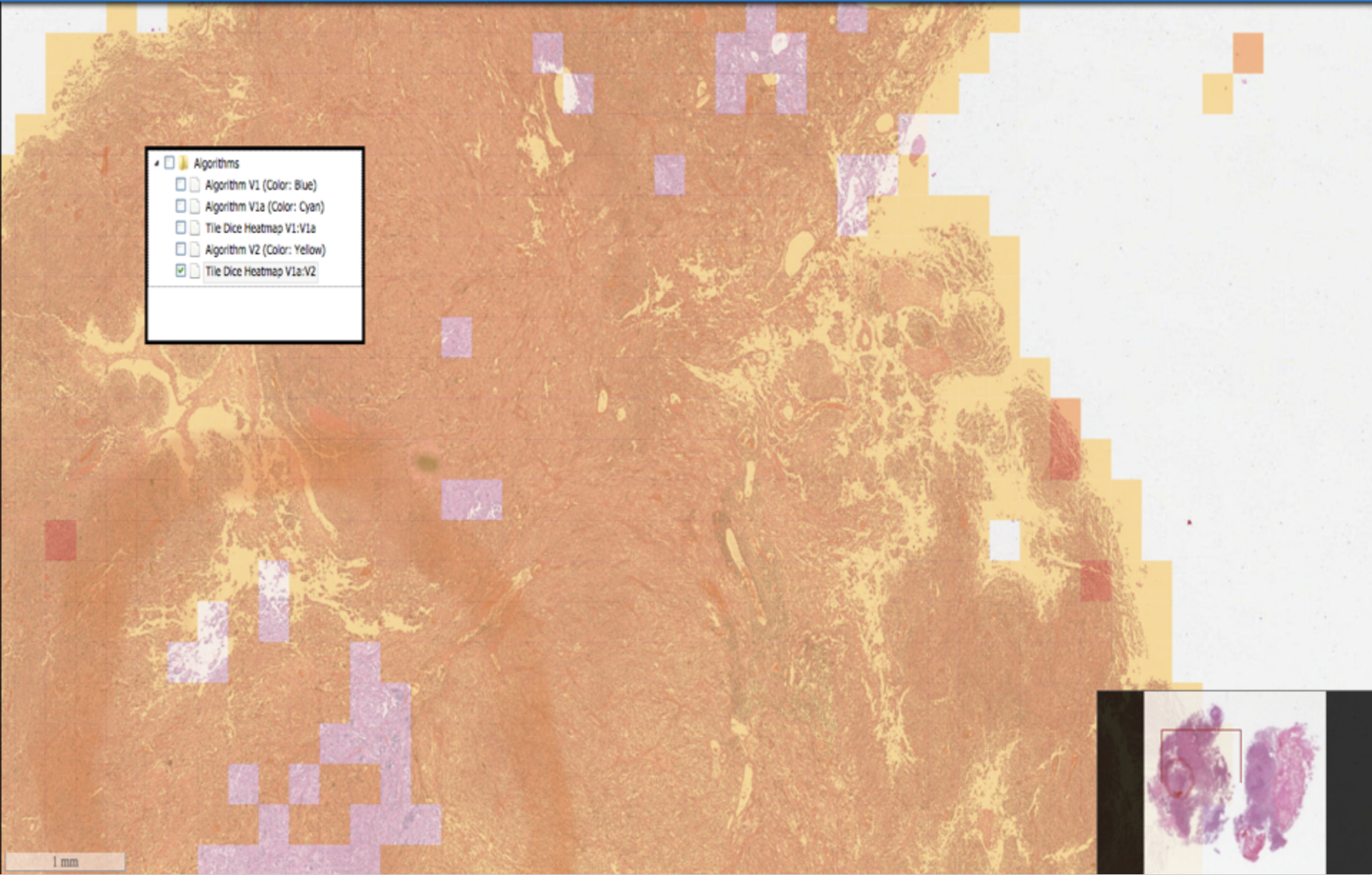


0 1

caMicroscope

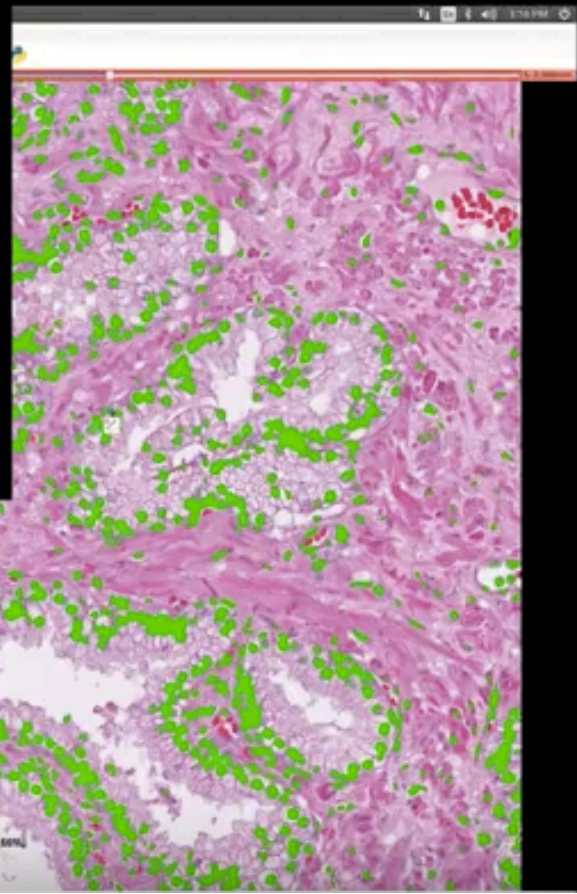
SubjectID :TCGA-12-1090-01Z-00-DX1

- Algorithms
  - Algorithm V1 (Color: Blue)
  - Algorithm V1a (Color: Cyan)
  - Tile Dice Heatmap V1:V1a
  - Algorithm V2 (Color: Yellow)
  - Tile Dice Heatmap V1a:V2



# 3D Slicer Pathology – Generate High Quality Ground Truth

ITCR - Tools to Analyze Morphology and Spatially Mapped Molecular Data



Threshold Plot  
Press Y for non-automatic segmentation on the current image using given parameters.  
Stop Quick TCGA Segmenter

• Nucleus Segmentation Parameters (H. Gao)

Offset Threshold:

Curvature Weight:

Size Threshold:

Size Upper Threshold:

Size Lower Threshold:

• Parameters

Enable Screenshots:

Screenshot scale factor:

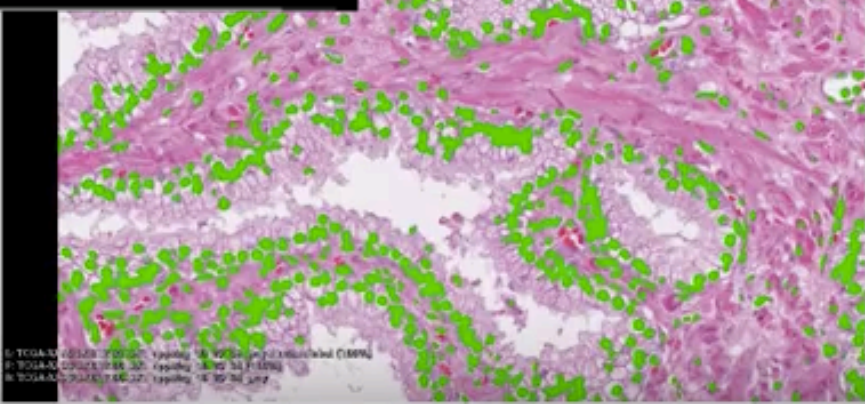
• Data Probe

Red 11, 202.0, P 820.3, S 0.0) Real Sp: 1.0

0 TCGA-K2-890K...\_ave\_label (1811, 820, 0) background (0)

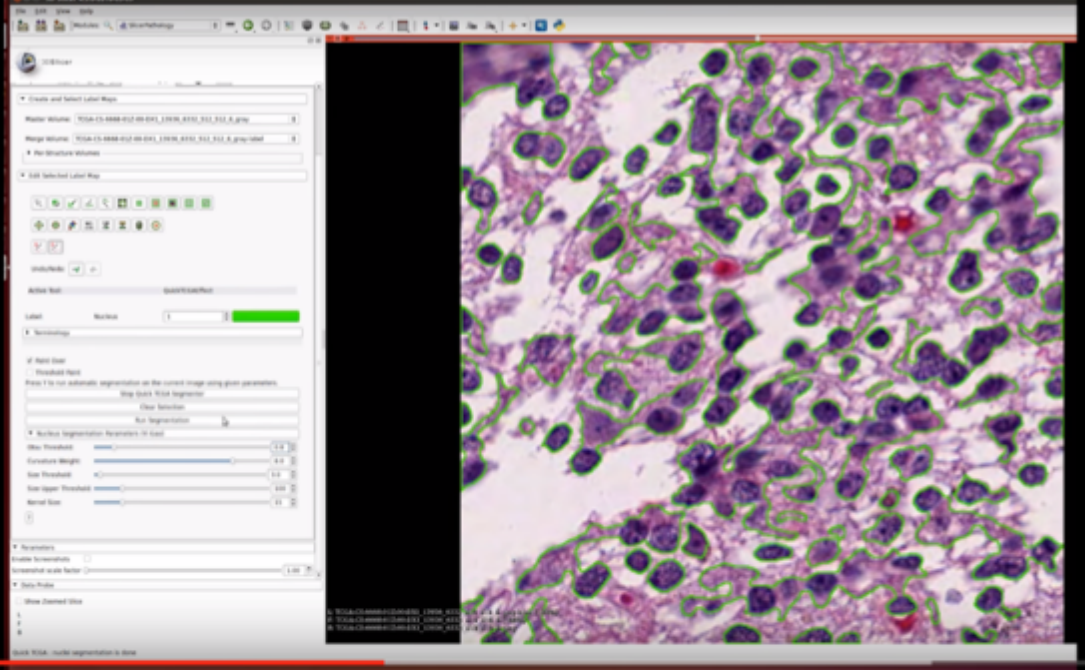
1 TCGA-K2-890K...\_40\_29\_38 (1811, 820, 0) 201, 142, 130

2 TCGA-K2-890K...\_9\_70\_gray (1811, 820, 0) 75

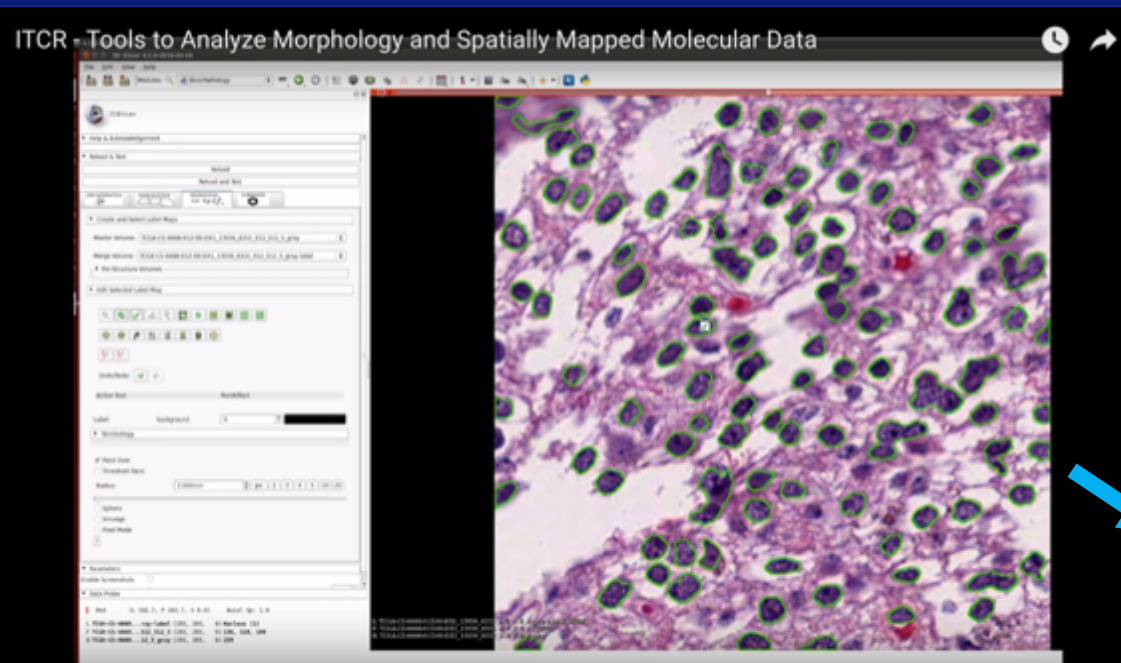


# Apply Segmentation Algorithm

ITCR - Tools to Analyze Morphology and Spatially Mapped Molecular Data



# Adjust algorithm parameters, manual fine tuning



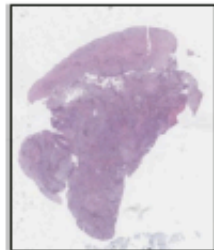
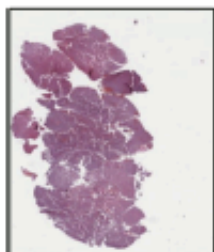
# Auto-tuning and feature extraction

- Goal – correctly segment trillions of objects (nuclei)
- Adjust algorithm parameters
- Autotuning – finds parameters that best match ground truth in an image patch
- Region template runtime support to optimize generation and management of multi-parameter algorithm results
- Eliminates redundant computation, manages locality
- Active Harmony – Jeff Hollingsworth!!
- Collaboration – George Teodoro, Tahsin Kurc

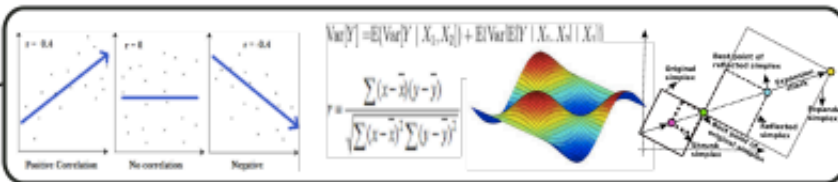


# Sensitivity Analysis (SA) and Auto-tuning methods

WSI input dataset

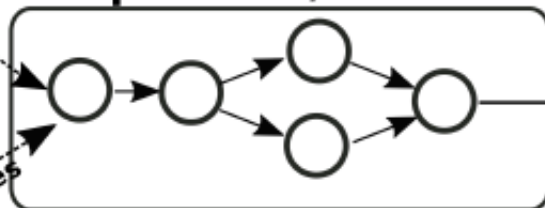


One or More Application Parameter Sets



Target Spatial metric: Dice, Jaccard, Intersect, etc.

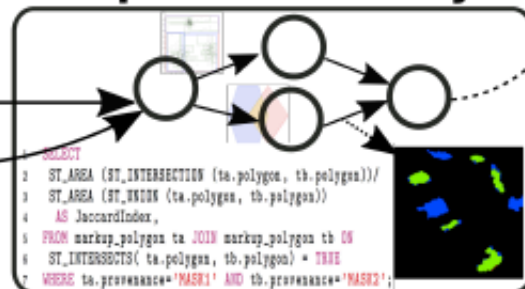
## Application Workflow Composition/Instantiation



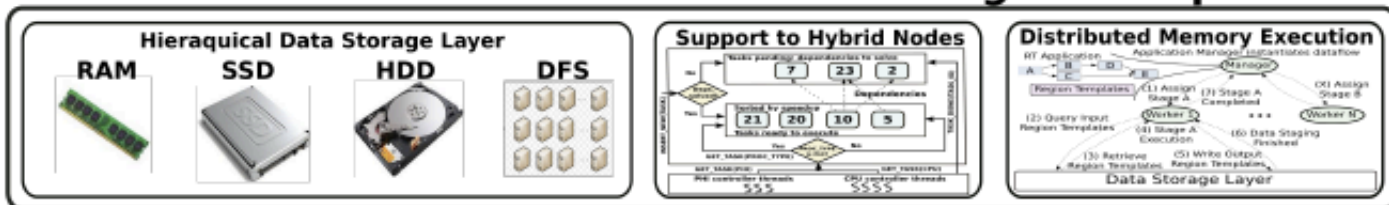
Segmentation Computed

Reference Segmentation

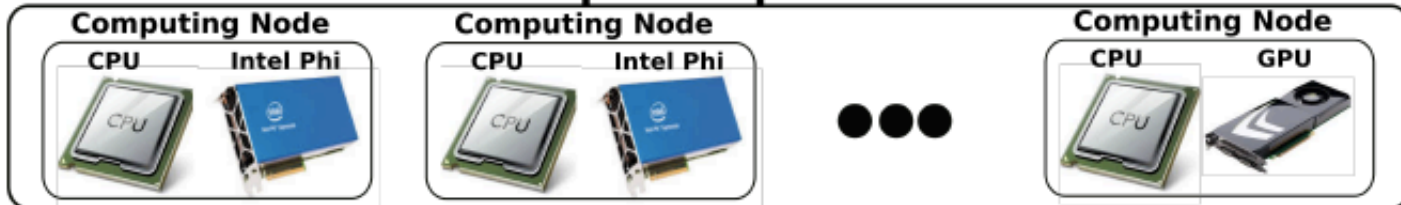
## Spatial Query-based Comparative Analysis



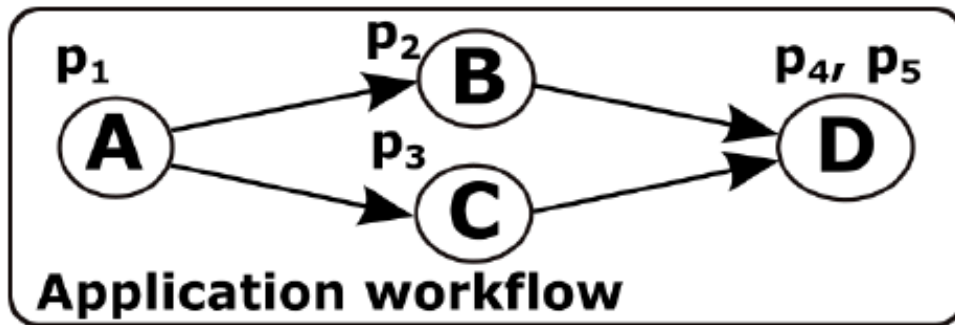
## Scalable and Efficient Execution with Region Templates



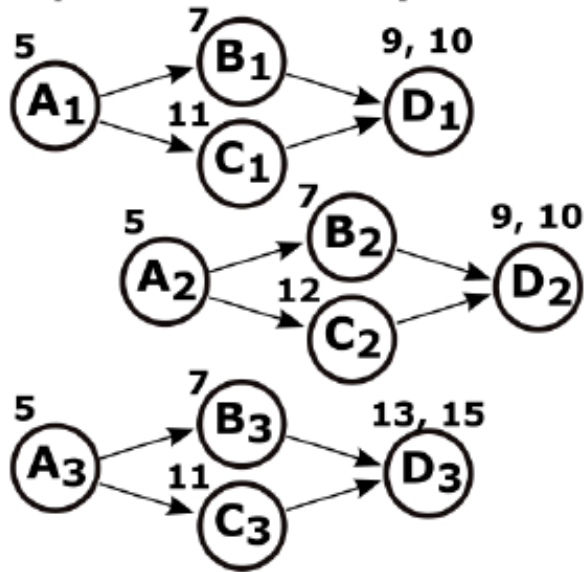
## Supercomputer



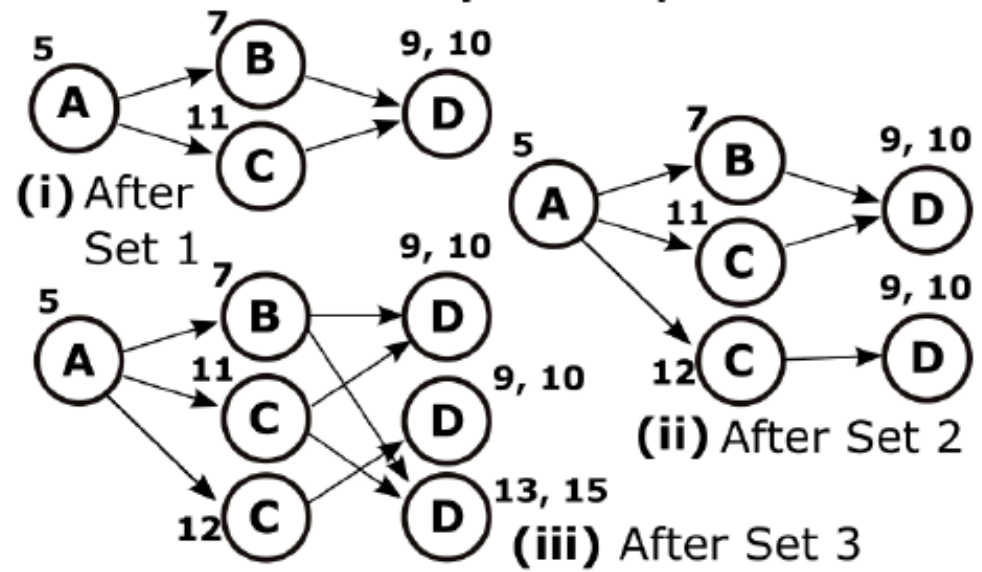
# E=Eliminate Duplicate Computations



Replica based composition

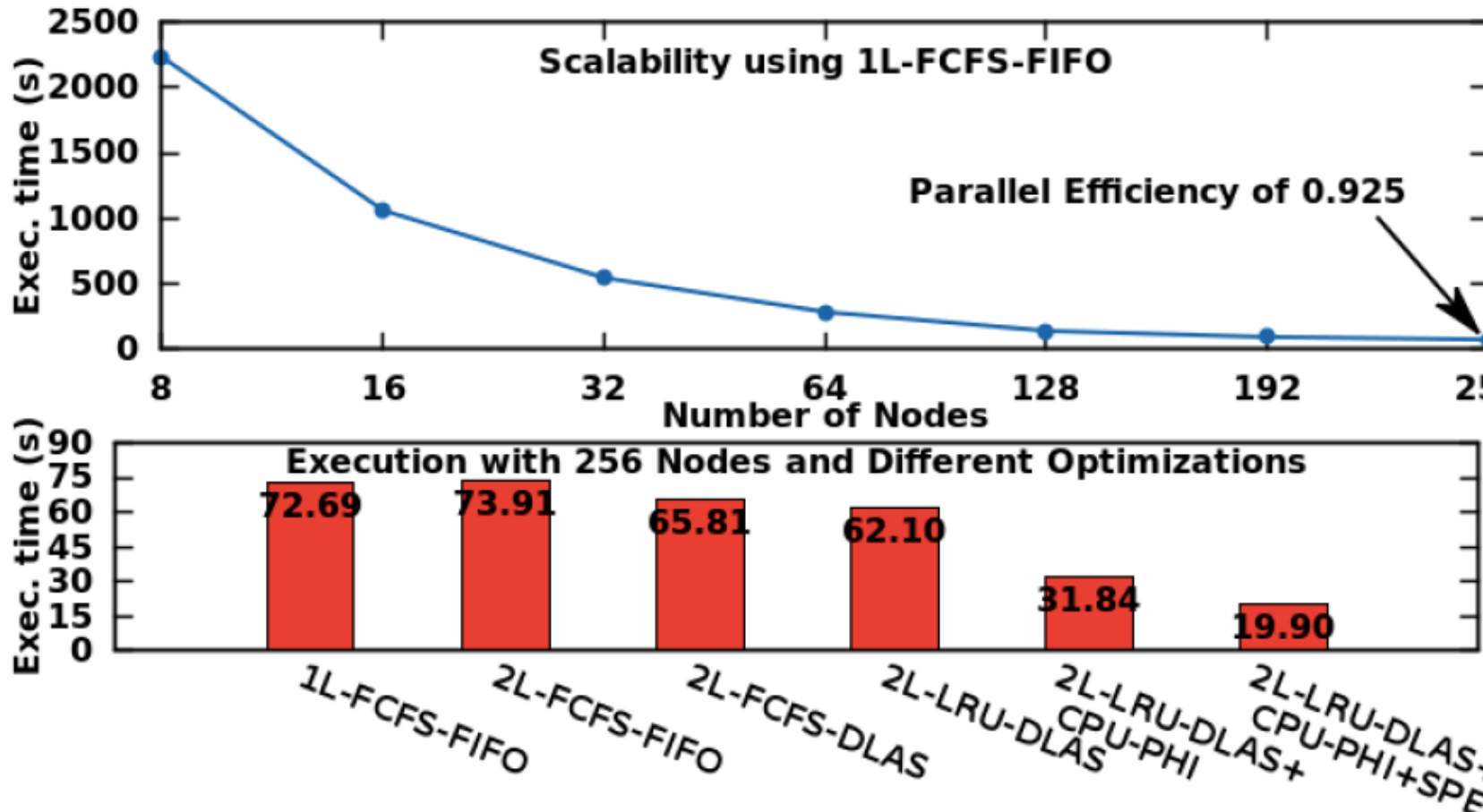


Compact composition



# Performance Optimization

256 nodes of Stampede. Each node of the cluster has a dual socket Intel Xeon E5-2680 processors, an Intel Xeon Phi SE10P co-processor and 32GB RAM. The nodes are inter-connected via Mellanox FDR Infiniband switches.



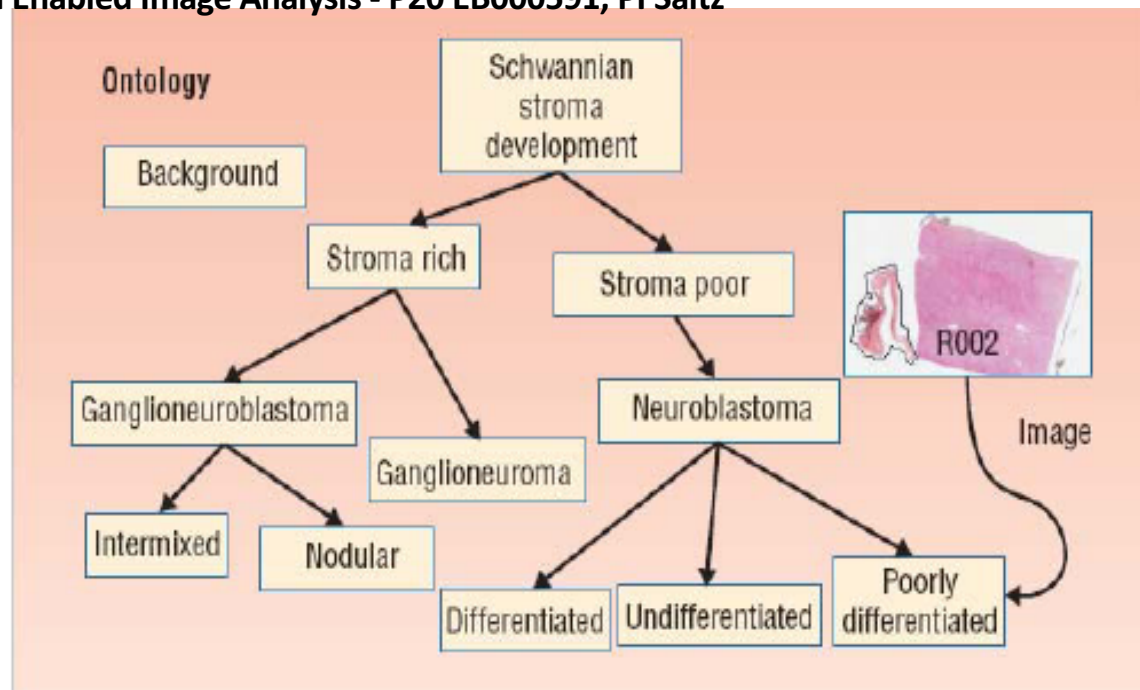
# Classification

- Automated or semi-automated identification of tissue or cell type
- Variety of machine learning and deep learning methods
- Classification of Neuroblastoma
- Classification of Gliomas
- Quantification of lymphocyte infiltration

# Classification and Characterization of Heterogeneity

- BISTI/NIBIB Center for Grid Enabled Image Analysis - P20 EB000591, PI Saltz

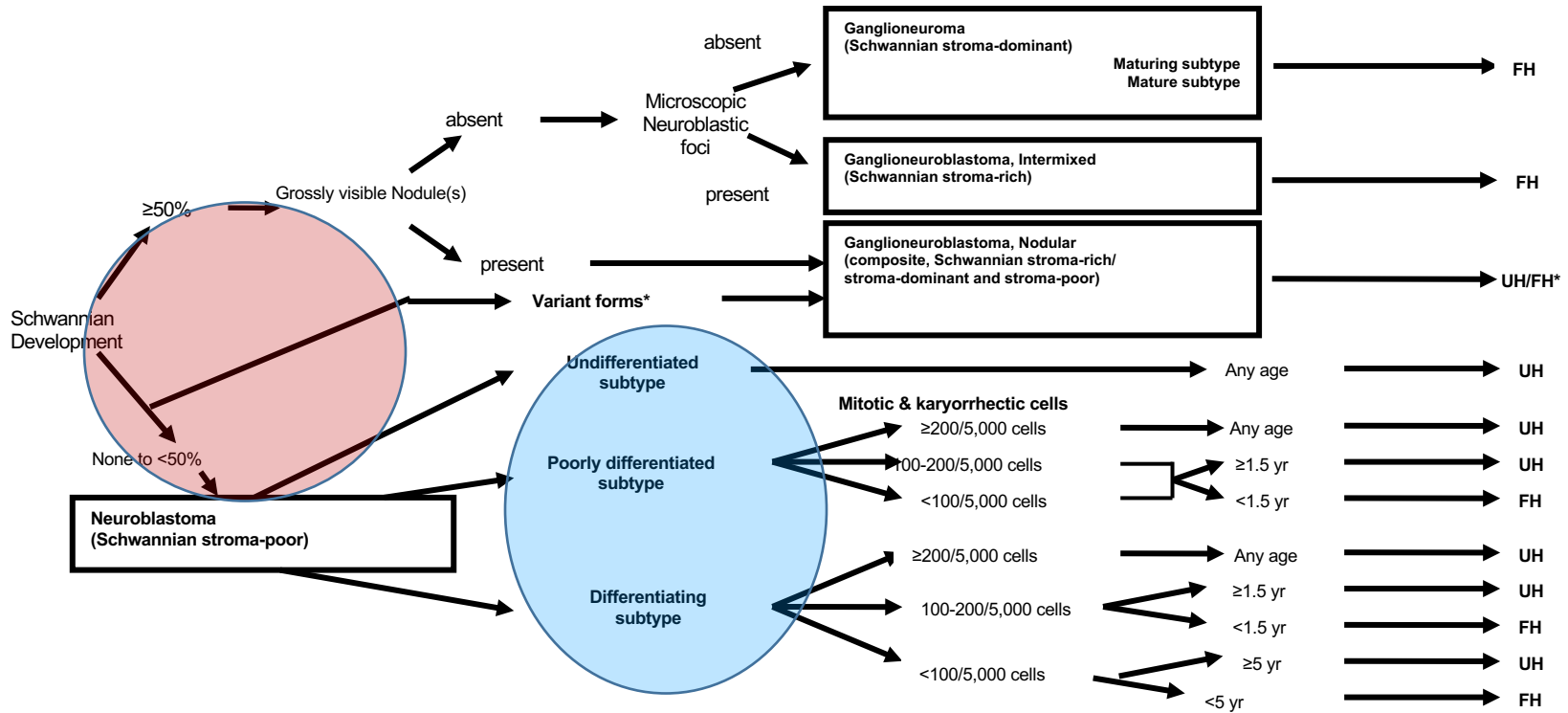
- Analyze images by computer
- Analyze the whole tissue, several slides
- Provide quantitative information to the pathologist
- Reduce inter- and intra-observer



Hiro Shimada, Metin Gurcan, Jun Kong, Lee Cooper Joel Saltz

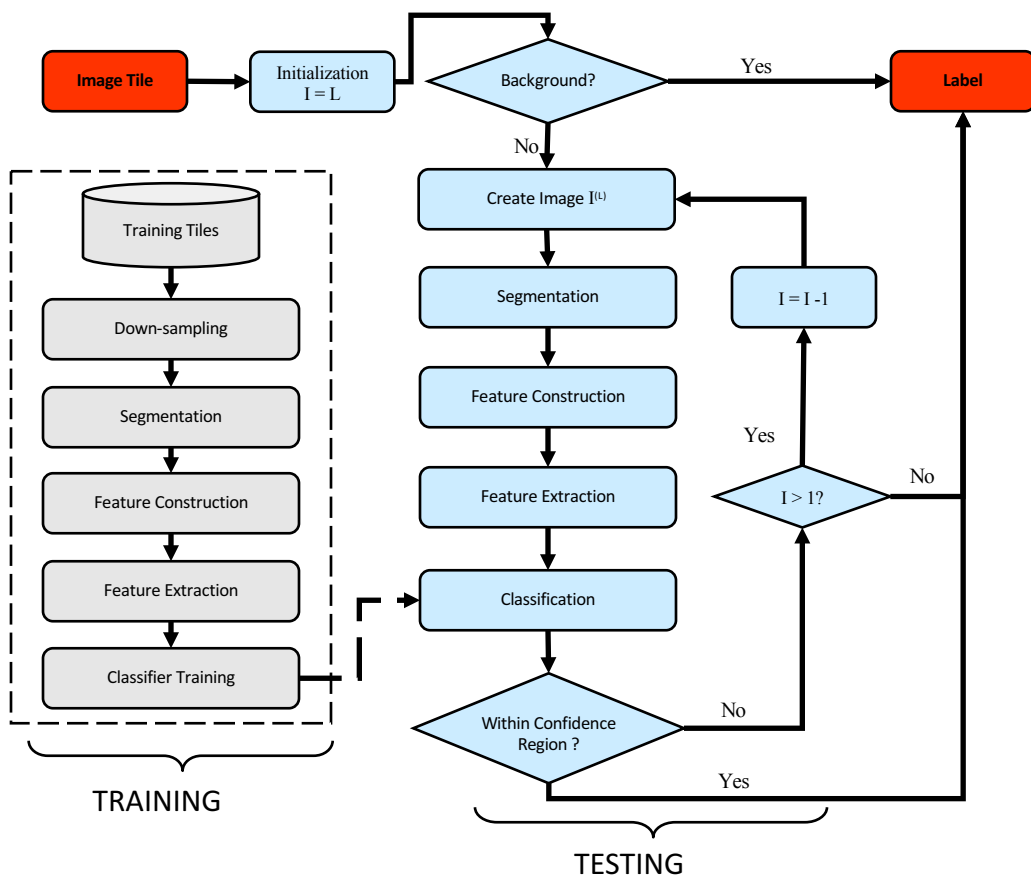
Gurcan, Shamada, Kong, Saltz

# Neuroblastoma Classification

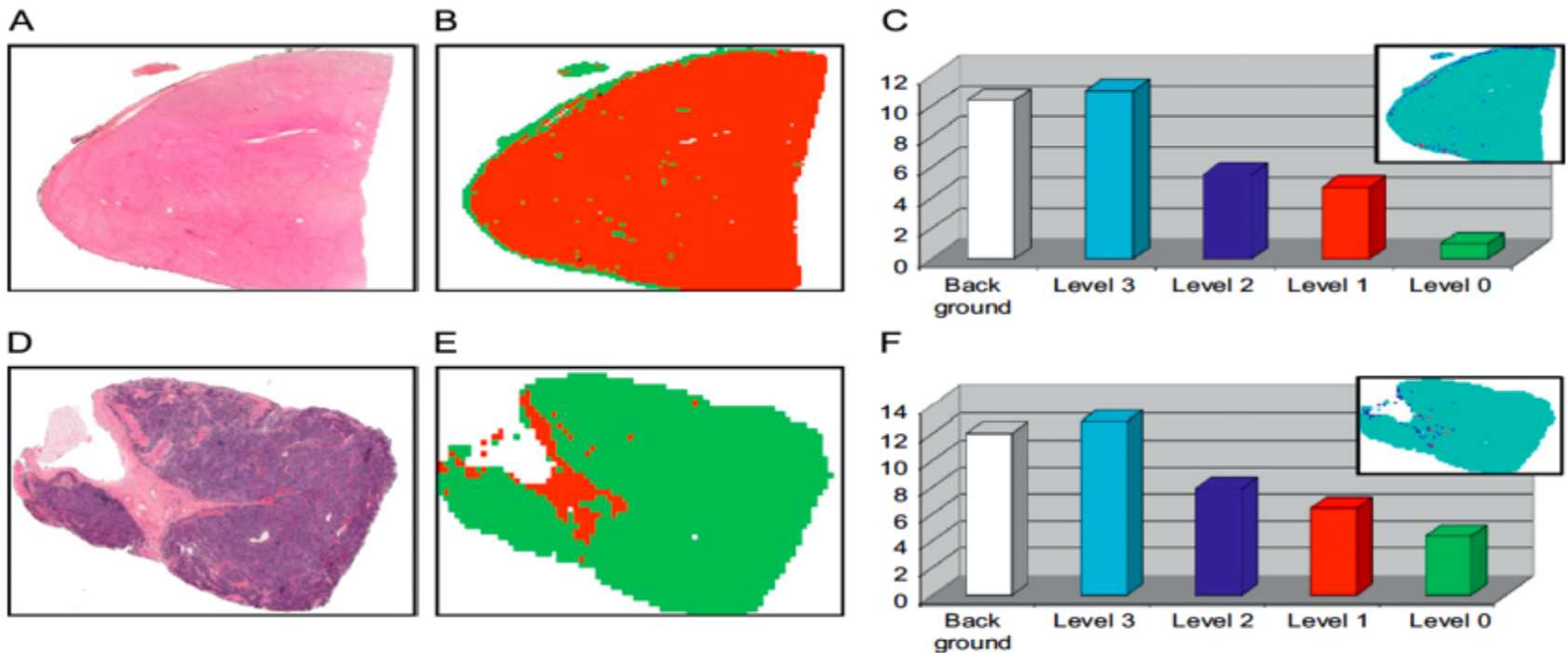


FH: favorable histology UH: unfavorable histology  
 CANCER 2003; 98:2274-81

# Multi-Scale Machine Learning Based Shimada Classification System



- Background Identification
- Image Decomposition (Multi-resolution levels)
- Image Segmentation (EMLDA)
- Feature Construction (2<sup>nd</sup> order statistics, Tonal Features)
- Feature Extraction (LDA) + Classification (Bayesian)
- Multi-resolution Layer Controller (Confidence Region)



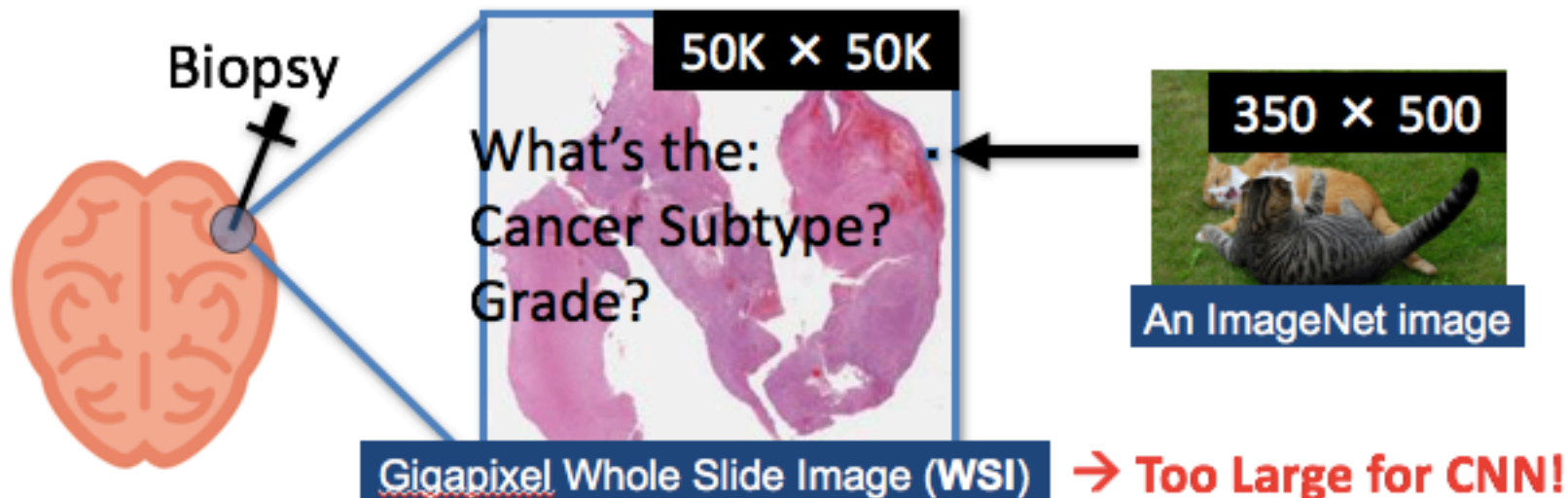
**Fig. 8.** Sample classification results after processing a whole-slide NB image. (a) and (d) are the H&E stained NB slides associated with stroma-rich and stroma-poor by an expert pathologist. (b) and (e) are the classification maps identified by the computerized system where the red color corresponds to stroma-rich regions and the green color corresponds to stroma-poor regions. (c) and (f) are the corresponding decision level statistics that show in log-scale the number of image tiles classified at a certain resolution level. In the resolution level map on the upper right, cyan color represents the lowest resolution and green color represents the highest resolution, respectively.



## Contributions

- Automatic **discriminative patch identification** for patch-CNN training
- A robust, general method to **combine patch-level predictions**

An important application: cancer classification



# Combining Information from Patches



Histogram of Patch-level Classes

## Advantage:

- The prediction of every patch counts  
→ robustness & generalization

SVM or  
Logistic Regression

Cancer  
subtype &  
grade

## Engineering details

- CNN architecture: AlexNet and VGG16
- Patch size: 500x500, Multiple scale
- Dataset: TCGA [[gdc-portal.nci.nih.gov](http://gdc-portal.nci.nih.gov)]
- Number of Patches: 1000 per WSI

# Brain Tumor Classification Results

## Glioma is

- The most common brain cancer
- The leading cause of cancer-related deaths in people under age 20

Methods	Accuracy
VGG16 features + <u>BoW</u> + SVM	0.667
Patch-CNN + Voting	0.710
Patch-CNN + Max-pooling	0.710
Our method	<b>0.771</b>
Pathologists' Agreement [M. Gupta 2015] (on a similar dataset)	0.7-0.8

**Confusion Matrix:** OA is very hard even for pathologists

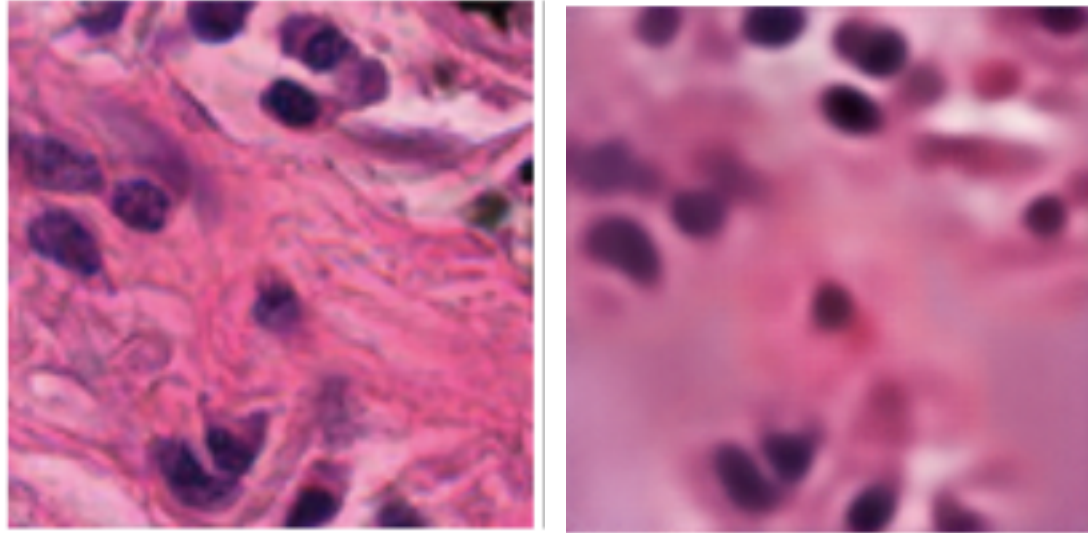
	GBM	OD	OA	DA	AA	AO
Glioblastoma, Grade IV (GBM)	214		2		1	
<u>Oligodendroglioma</u> , Grade II (OD)	1	47	22	2		1
<u>Oligoastrocytoma</u> , Grade II & III (OA)	1	18	40	8	3	1
Diffuse Astrocytoma, Grade II (DA)	3	9	6	20		1
Anaplastic Astrocytoma, Grade III (AA)	3	2	3	3	4	
Anaplastic <u>Oligodendroglioma</u> , Grade III (AO)	2	2	3			1

Le Hou, Dimitris Samaras, Tahsin Kurc, Yi Gao, Liz Vanner, James Davis, Joel Saltz

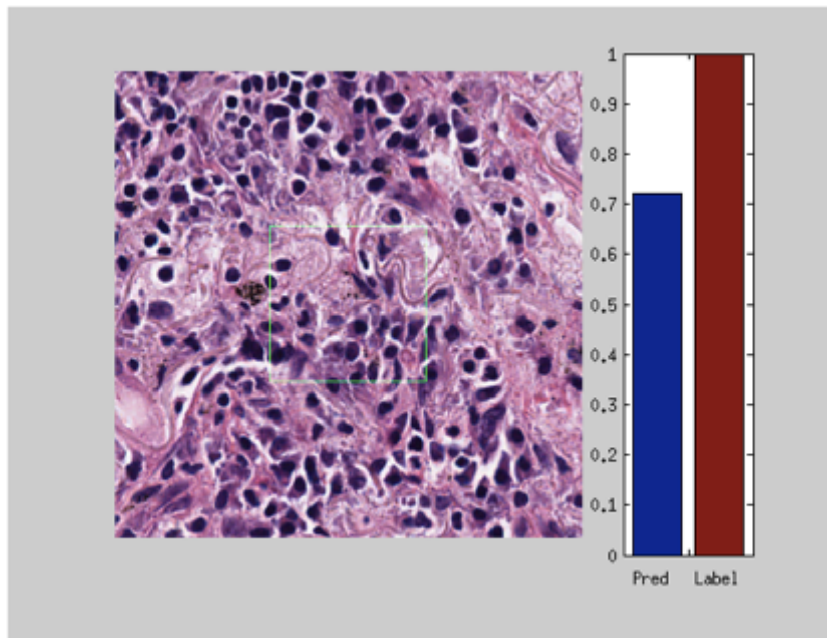
# Tumor Infiltrating Lymphocyte quantification

- Convolutional neural network to classify lymphocyte infiltration in tissue patches
- Convolutional neural network and random forest to classify individual segmented nuclei
- Extensive collection of ground truth
- Joint work with Emory and TCGA PanCanAtlas Immune group

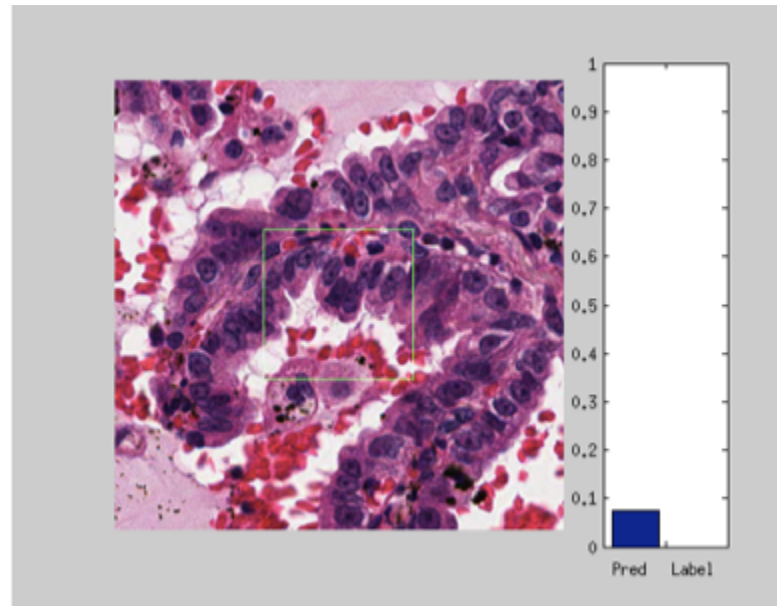
Unsupervised Autoencoder – 100 feature dimensions



# Lymphocyte identification



Lymphocytes Infiltration



No Lymphocyte Infiltration

# Receiver Operating Characteristic – Area Under Curve – 95%

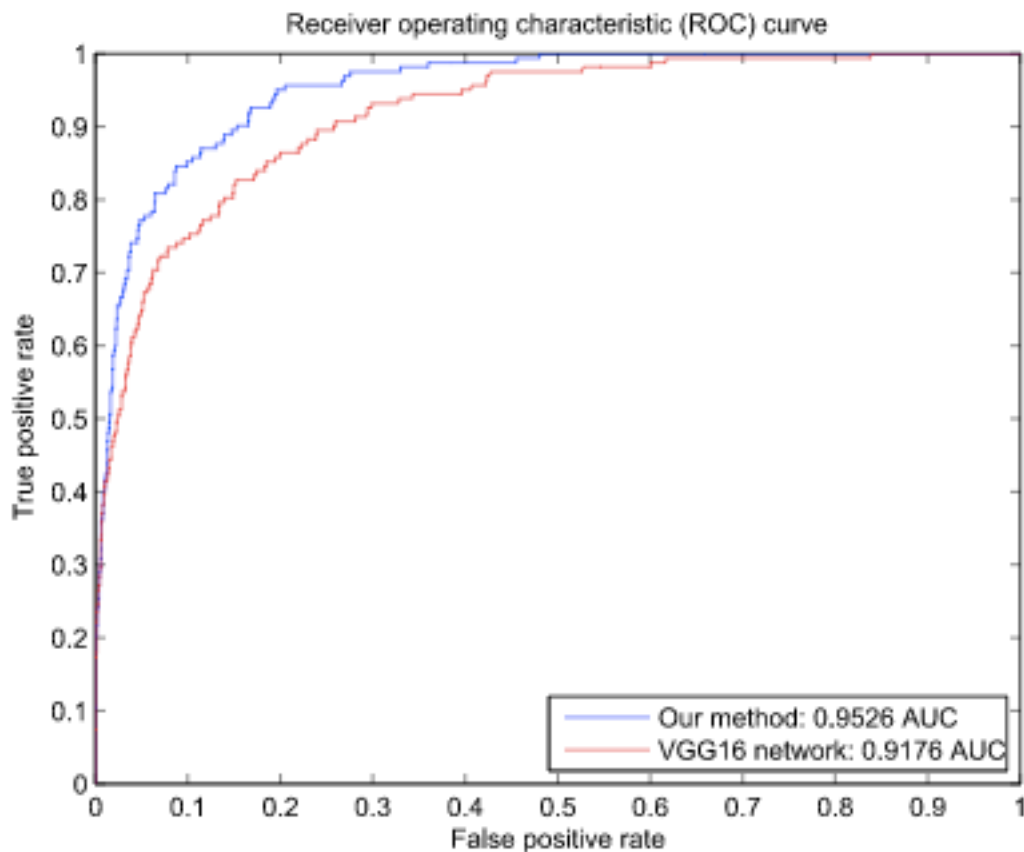
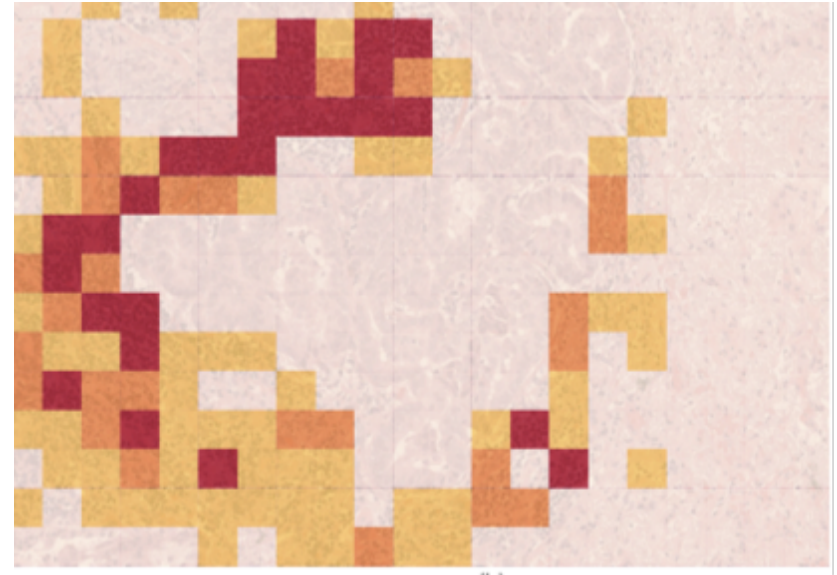
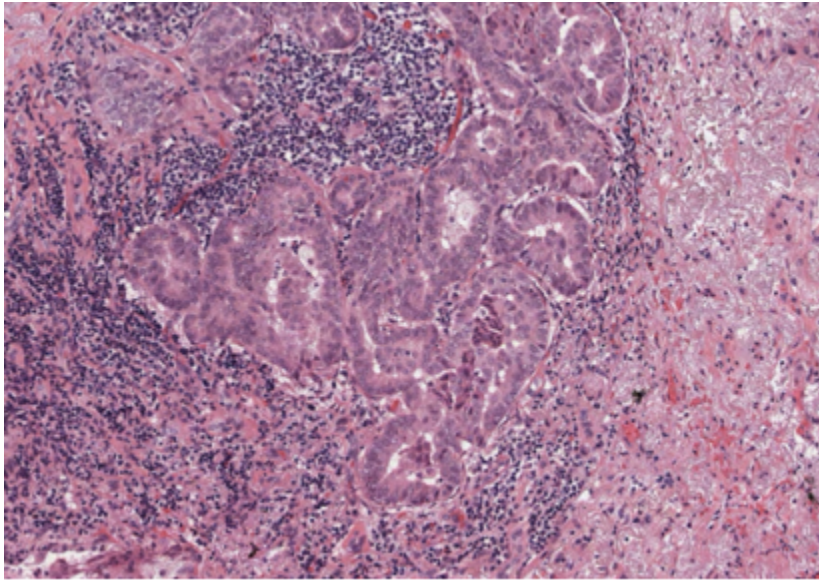


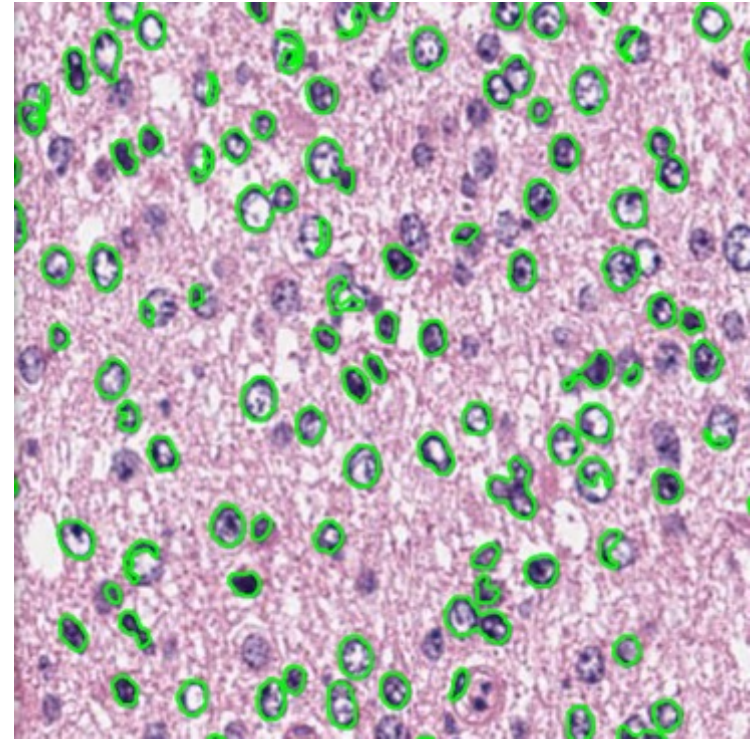
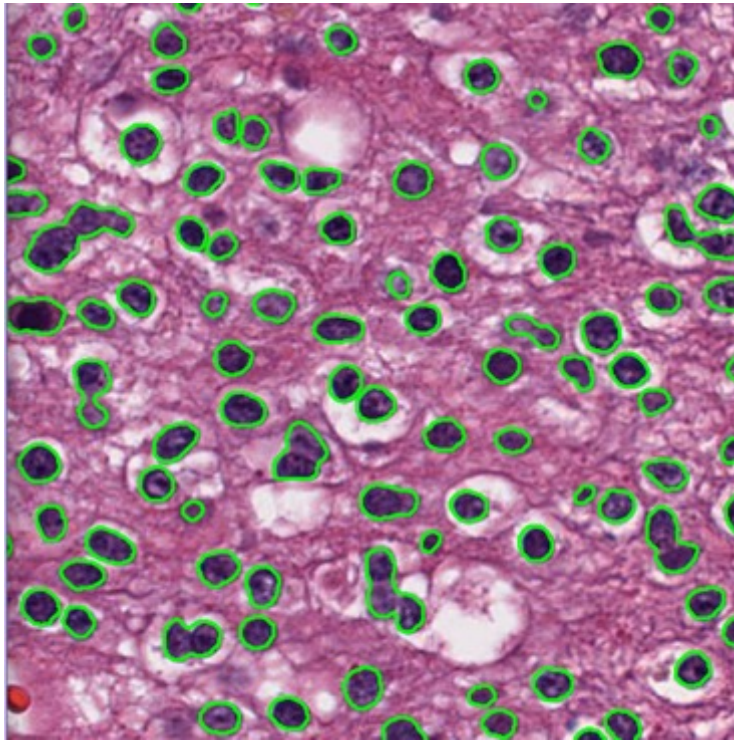
Figure 2 ROC curve

# Lymphocyte Classification Heat Map

Trained with 22.2K image patches  
Pathologist corrects and edits



# Machine Learning and Quality Critiquing



	Good	Bad
Test as Good	2916	33
Test as Bad	28	2094



# Dissemination

- Containers
- Cloud
- TCIA
- HPC via NSF and DOE
- TCGA – PanCanAtlas – Lymphocyte characterization
- Integrated Features/NLP joint with TIES

# ITCR Team

## **Stony Brook University**

Joel Saltz

Tahsin Kurc

Yi Gao

Allen Tannenbaum

Erich Bremer

Jonas Almeida

Alina Jasniewski

Fusheng Wang

Tammy DiPrima

Andrew White

Le Hou

Furqan Baig

Mary Saltz

## **Emory University**

Ashish Sharma

Adam Marcus

## **Oak Ridge National Laboratory**

Scott Klasky

Dave Pugmire

Jeremy Logan

## **Yale University**

Michael Krauthammer

## **Harvard University**

Rick Cummings

# Funding – Thanks!

- This work was supported in part by U24CA180924-01, NCIP/Leidos 14X138 and HHSN261200800001E from the NCI; R01LM011119-01 and R01LM009239 from the NLM
- This research used resources provided by the National Science Foundation XSEDE Science Gateways program under grant TG-ASC130023 and the Keeneland Computing Facility at the Georgia Institute of Technology, which is supported by the NSF under Contract OCI-0910735.

Thanks!