

VIDIA: a HUBzero Gateway for Data Analytics Education

Jeanette Sperhac and Steven M. Gallo
Center for Computational Research, University at Buffalo

11th Gateway Computing Environments Conference
San Diego Supercomputer Center
San Diego, CA
2 November 2016

CENTER FOR **COMPUTATIONAL RESEARCH**



Outline

1. The collaboration
2. The VIDIA instance
3. Supporting teaching
 - Tools
 - Coursework
4. Enabling research
5. Demonstration

The VIDIA Collaboration



- Academic supercomputing facility
- Buffalo, New York
- 10,000+ cores, 170+ Tflops compute capacity
- Four-year liberal arts college
- Oneonta, New York
- Enrollment 5,900 students

CENTER FOR **COMPUTATIONAL RESEARCH**



Collaboration Goals

- Create a social sciences big data discovery environment
- Support social science teaching and research
- Leverage High Performance Computing (HPC) resources

CENTER FOR **COMPUTATIONAL RESEARCH**



**SUNY
ONEONTA**

How is this platform different?



- Developed for scientific collaboration
- Access via web browser
- Computational horsepower
 - Deployed tools run in virtual containers
 - HPC access via *submit*

CENTER FOR **COMPUTATIONAL RESEARCH**

vidia.ccr.buffalo.edu

VIDIA

Discover

Community

About

Support

Login

Register

VIDIA - Virtual Infrastructure for Data Intensive Analysis

Search



Virtual Infrastructure for Data Intensive Analysis (VIDIA)

This scalable, community-driven infrastructure exposes SUNY students and faculty at Primarily Undergraduate Institutions (PUIs) to data-intensive computing and analysis techniques.

VIDIA is a partnership between the Center for Computational Research at the University at Buffalo, and SUNY Oneonta and is made possible with funding from the 2013 SUNY Innovative Instruction Technology Grants (IITG) program.

Resources

Browse tools, publications, and utilities.

[Center for Computational Research \(CCR\)](#)
[Cluster Status](#)

Current collaborators include

[SUNY Oneonta Big Data Project](#)

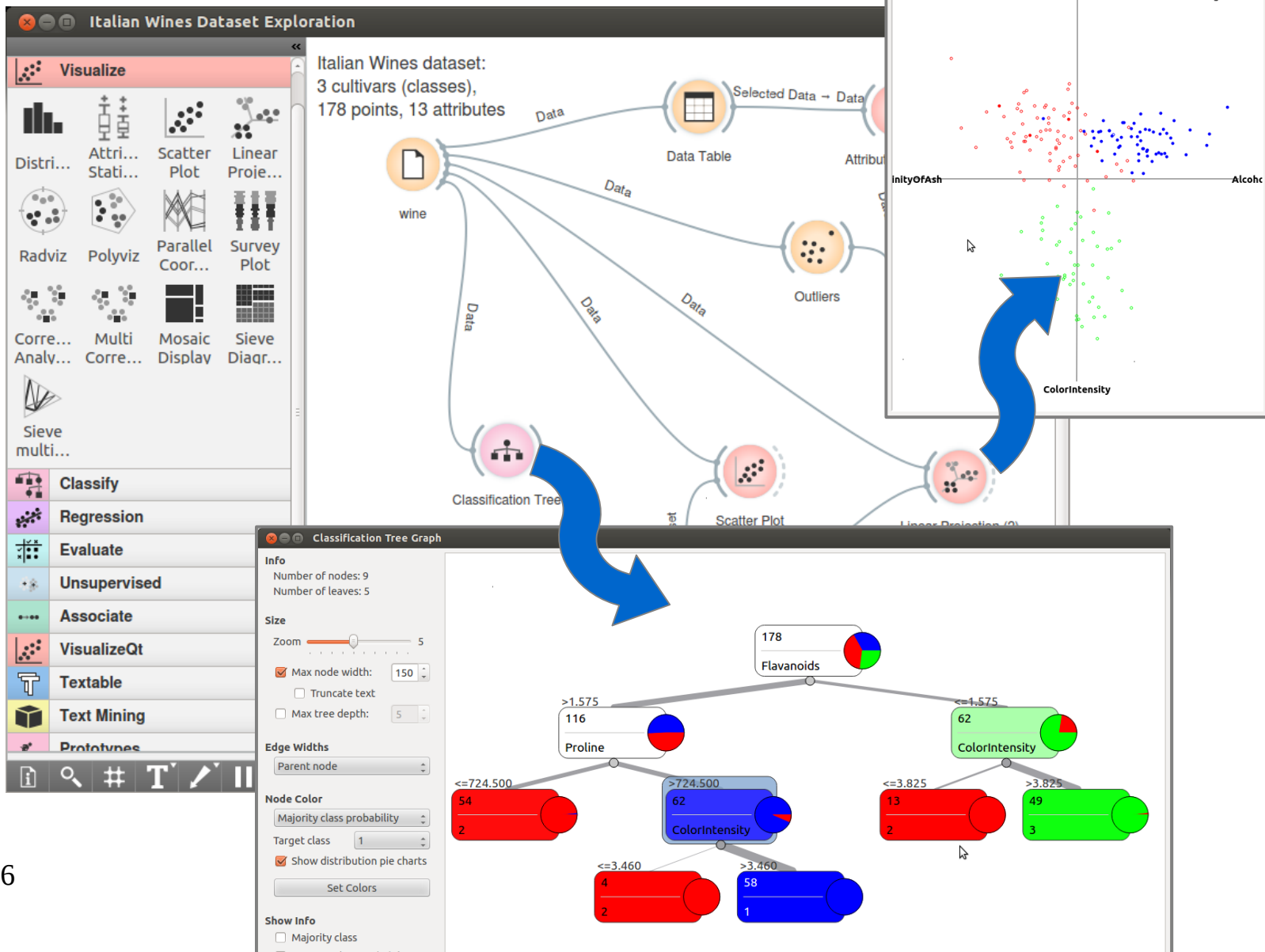
[Oneonta Biochemistry](#)

[SUNY Genesee](#)

[SUNY Upstate](#)

- CCR
- cluster resources
- data analysis
- documentation
- download
- file transfer
- HUBbub
- HUBzero
- Philosophy
- PSPP
- Python
- R
- RapidMiner
- rappture
- RStudio
- statistics
- submit
- tool development
- tutorial
- upload
- using VIDIA
- video
- VIDIA
- visualization
- webDAV

Provide workflow tools for data analysis



Curate datasets of social scientific interest

<input checked="" type="radio"/> Data View <input type="radio"/> Meta Data View <input type="radio"/> Plot View <input type="radio"/> Advanced Charts <input type="radio"/> Annotations				
ExampleSet (3000 examples, 0 special attributes, 3 regular attributes)			View Filter (3000 / 3000): all	
Row No. ▲	Description	Media Type	Sentiment	
134	#fracking is #frackedtags: #fracked #fracking	Images/Videos	Neutral	
135	Fracking: The Bay Delta Conservation Plan Would Provide Water for Mining http://t.co/0qCPIFmLkb #Fracking	Twitter	Neutral	
136	#Fracking #FrackNo	Facebook	Neutral	
137	I don't understand how @BarackObama allows FRACKING to be done when he claims to care so much about HEALTH & HI	Twitter	Negative	
138	I would really like to compare the Utah-fracking maps to the dead-bald-eagles-found-in-Utah map. #fracking #Poison	Twitter	Neutral	
139	RT @robedwards53: Fracking Hormone-Disrupting Chemicals Linked To #fracking Found In Colorado River http://t.co/BSkl	Twitter	Neutral	
140	#hydrofracking #fracking #banfrackingnow #cce #environmentalist #realist #truth #liveright #stopthenonsense #stop	Images/Videos	Neutral	
141	Fiery #Oil #Train Crash in Raging #Shale #Oil Boom State of North Dakota http://t.co/rGc5y8WYqd #fracking @nytimes ht	Twitter	Positive	
142	Conservation and fracking? #no MT @iamgreenbean: .@JerryBrownGov signed law that facilitates #fracking boom in CA I	Twitter	Positive	
143	Fracking the great state of Western Australia#wapol #Auspol #Csg #FrackOff #Fracked #Fracking #Perth http://t.co/y49	Twitter	Positive	
144	RT @onahunttoday: RT @robedwards53: Fracking Chemicals With Hormone-Disrupting Chemicals Linked To #fracking Fou	Twitter	Neutral	
145	Romanian authorities or Chevron? - Sarah in Romania http://t.co/UHnj2okd37 #Pungesti #gazedesist #Chevron #police #	Twitter	Neutral	
146	A #villarrobledo tamb lluiten contra el #fracking #osona #ripolls #igerscatalunya #igerscastillalamancha #albaceteTag:	Images/Videos	Neutral	
147	#naturalgas #fracking #propane #art #landscape #galleryTags: #naturalgas #art #propane #fracking #landscape #g	Images/Videos	Neutral	
148	North Dakota Oil Train Fire Spotlights Risks of Transporting Crude http://t.co/rB3g69clh2 #fracking #shale	Twitter	Neutral	
149	Breaking: Yet another Bakken Shale oil train explodes into flames http://t.co/hGZwXJWQgE #fracking #shale	Twitter	Negative	
150	On but one finite in nature planet #fracking is fricking #madness, we better	Google	Positive	
151	RT @YourAnonNews: .@Pandora_Radio: stop playing misleading advertisements about #fracking in Colorado! http://t.co/	Twitter	Negative	
152	How many people have to light their tap on fire before folks realize fracking is bad? http://t.co/zPXlQp848w #fracking #er	Twitter	Negative	
153	RT @TheFrackingTrap: Fracking: North Dakota Video Shows Flammable Water http://t.co/THJhFkMaE #fracking #shale	Twitter	Neutral	
154	Fracking: North Dakota Video Shows Flammable Water http://t.co/THJhFkMaE #fracking #shale	Twitter	Neutral	
155	Not only does #fracking provide us with ten cent cheaper per gallon gasoline, it also kills us, pollutes the environment, e	Images/Videos	Negative	
156	Fiery Oil Train Crash in Raging Shale Oil Boom State of North Dakota http://t.co/BjZStWvBEt #fracking #shale	Twitter	Positive	
157	Barbing shot for 2012. #piken #fracking #shalegas #natgas #ohio #harrisoncounty #serialTags: #piken #shalegas #a	Images/Videos	Neutral	

Enable access to HPC resources



VIDIA Hardware

HUBzero instance: Dell PowerEdge R720xd

- 2x 6-core Intel Xeon E5-2630

2.30 GHz, 15M cache

- 48 TB raw SATA disk space

~36 TB usable

- 128 GB memory

16x8GB - 1333MHz DIMMS

CENTER FOR **COMPUTATIONAL RESEARCH**

Teaching on HUBzero



- Platform for tool, data, and HPC access
- Easy on campus IT staff
- Access anytime, anywhere
- Resources can be selectively secured
- Students may access resources after course conclusion

CENTER FOR **COMPUTATIONAL RESEARCH**

Initial Teaching Toolset

Features wanted:

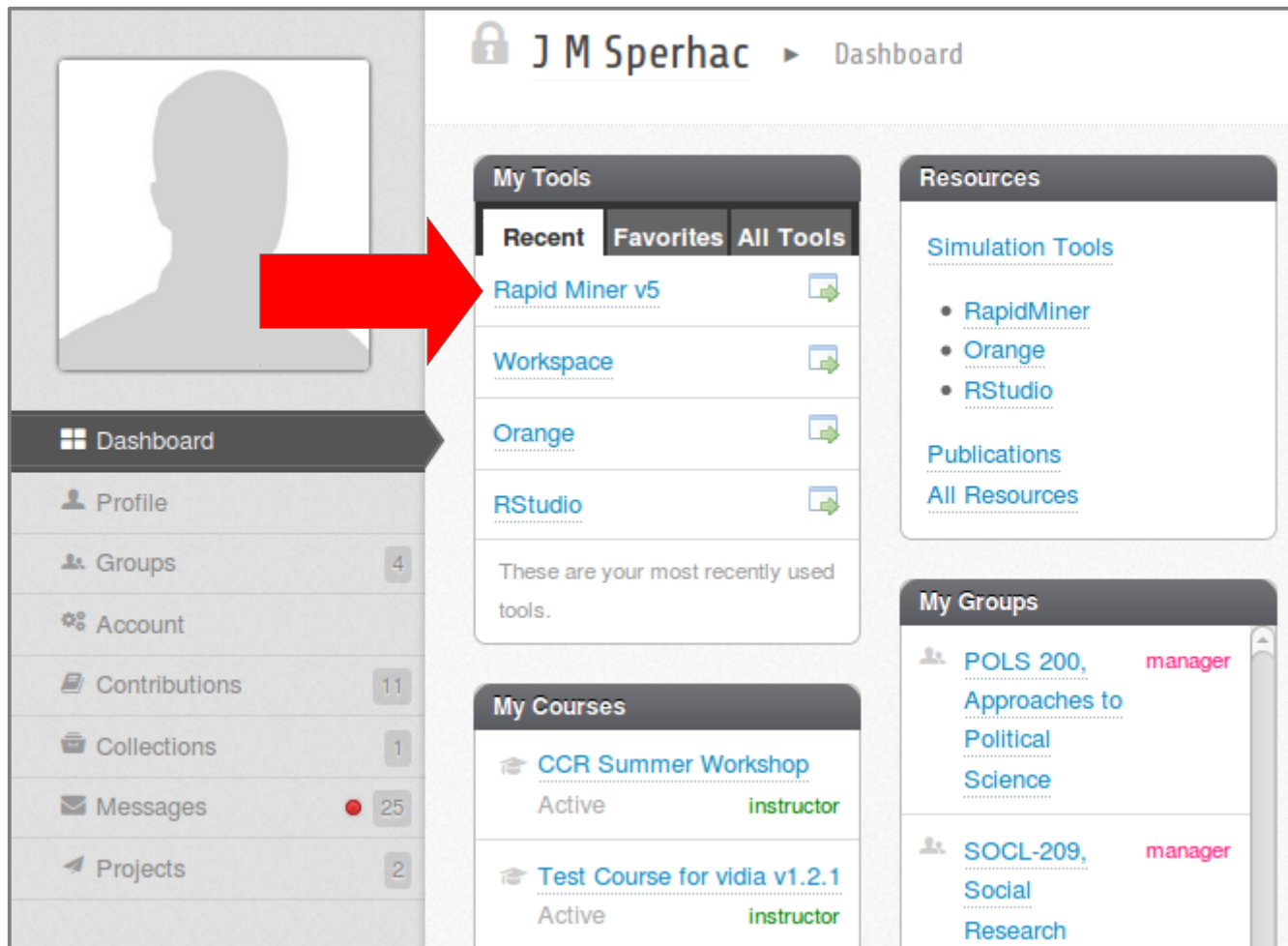
- Graphical User Interface
- Powerful, easy to use
- Open source, extensible

Tools deployed:

- RapidMiner
- Orange
- RStudio
- Python/Octave IDEs
- Gnu PSPP

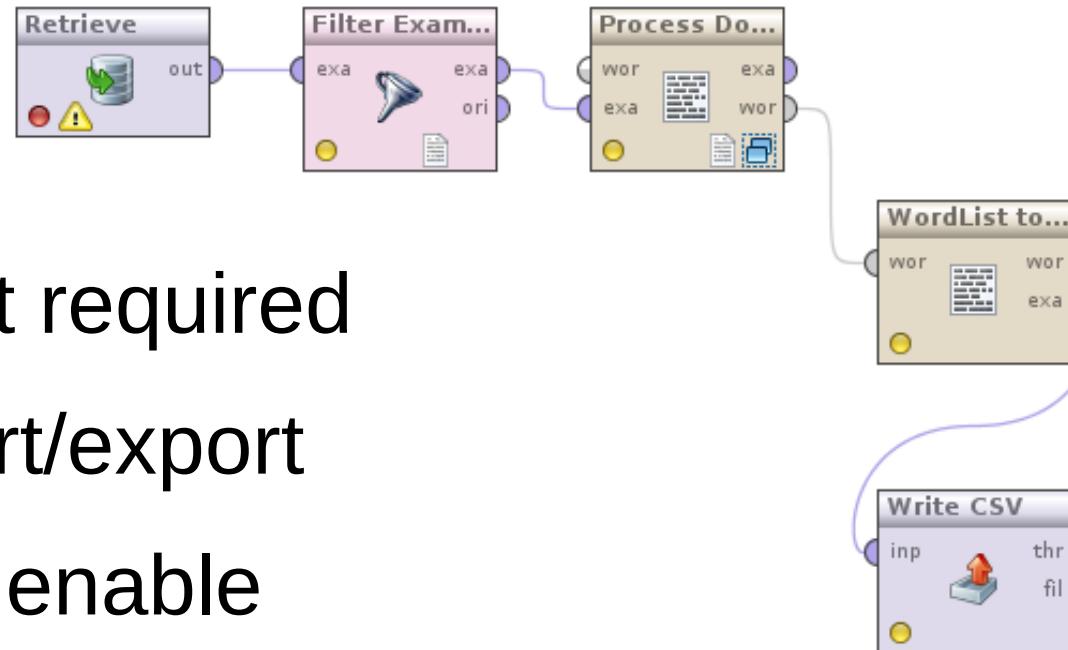


Deployed tool on VIDIA



RapidMiner Data Mining Tool

RapidMiner Tool



- Coding not required
- Data import/export
- Operators enable
 - Construction of analysis workflows
 - Text processing and analysis
 - Plotting and visualization

hub

Local Repository/processes/oneonta-degrees-v2.xl

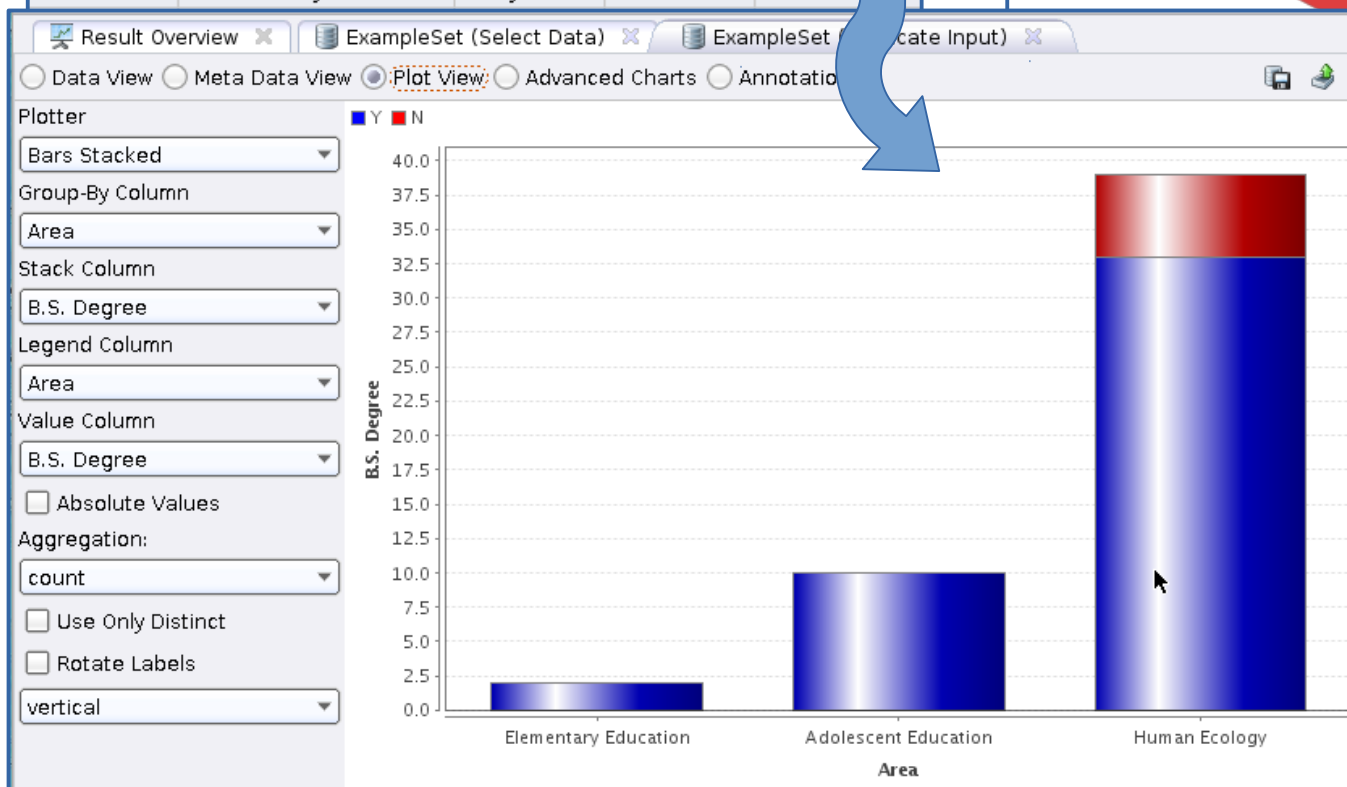
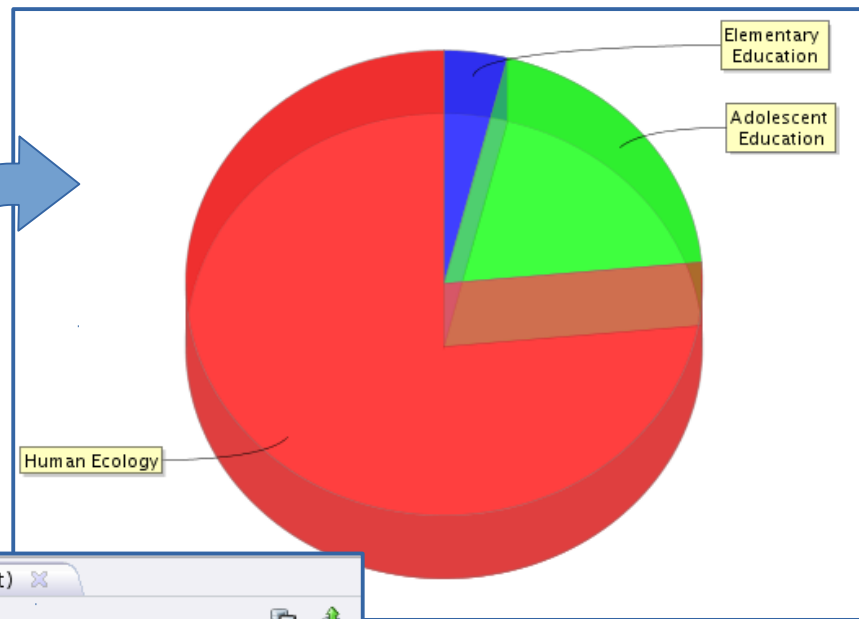
File Edit Process Tools View Help

Result Overview ExampleSet (Select Data)

Data View Meta Data View Plot View Advanced Charts

ExampleSet (51 examples, 0 special attributes, 5 regular attributes)

Row No.	Area	Major	B.A. Degree	B.S. Degree
16	Human Ecology	Africana &	Y	N
17	Human Ecology	Anthropolo	Y	N
30	Human Ecology	Internation	Y	N
31	Human Ecology	Internation	Y	N
40	Human Ecology	Music	Y	N
41	Human Ecology	Music Indu:	Y	N
1	Elementary Education	Childhood I	N	Y
2	Elementary Education	Early Child	N	Y



Development Environment

(write your own tools)

Features:

- Linux container
- Command line
- HPC access
- Source code control
- GUI builder
- Compilers

Supported languages:

- Fortran, C, C++
- Python, Perl, Ruby
- Octave
- Java
- R
- ...



Coursework on VIDIA

- POLS-284: U.S. Foreign Policy
- PHIL-231: Media Ethics
- SOCL-209: Social Research Methods
- SOCL-294: Animals and Society
- SOCL-244: Environmental Sociology

CENTER FOR COMPUTATIONAL RESEARCH



**SUNY
ONEONTA**

Case Study: Animals and Society

- Sociology majors, 200 level
- No programming experience
- Comparative, social scientific analyses

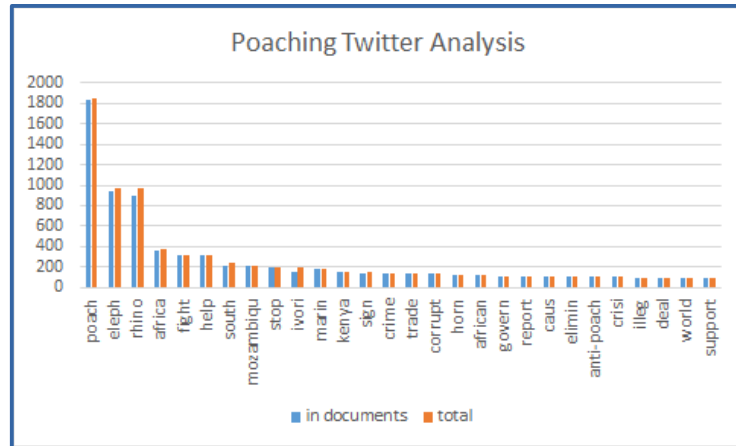
Gather, organize, and interpret social media data

CENTER FOR COMPUTATIONAL RESEARCH




**SUNY
ONEONTA**

Case Study: Animals and Society



Chose search terms from journal articles on animal poaching.
Collected data set of *1831 tweets* on these search terms.
Generated list of commonly occurring terms using RapidMiner.
The keyword *poach* occurred 1857 times. Top *co-occurring terms*:

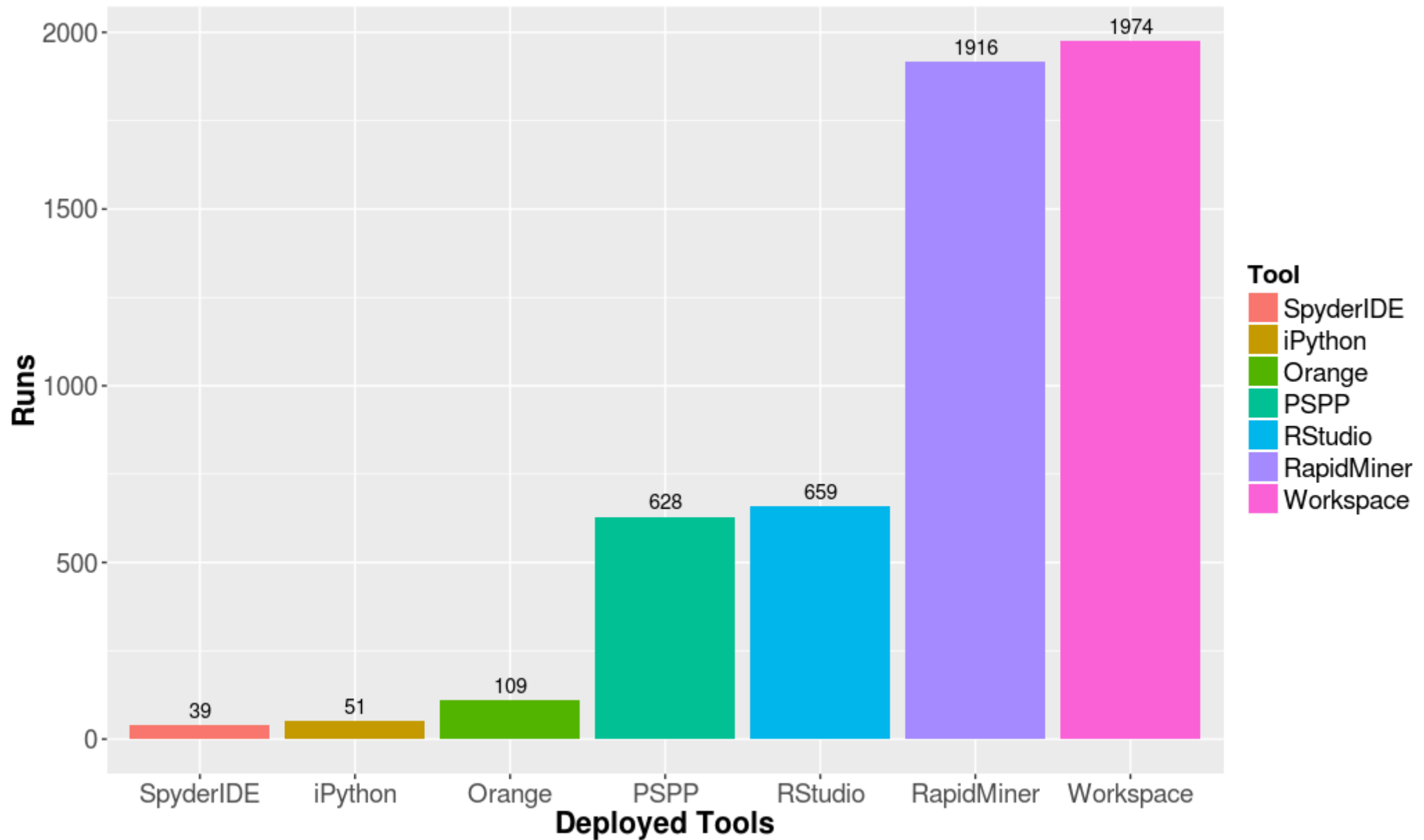
- Elephant (946 tweets/976 total) and Rhino (896/976)
- Africa (363/376) and Fight (313/317)
- Help (311/312) and Stop (192/198)



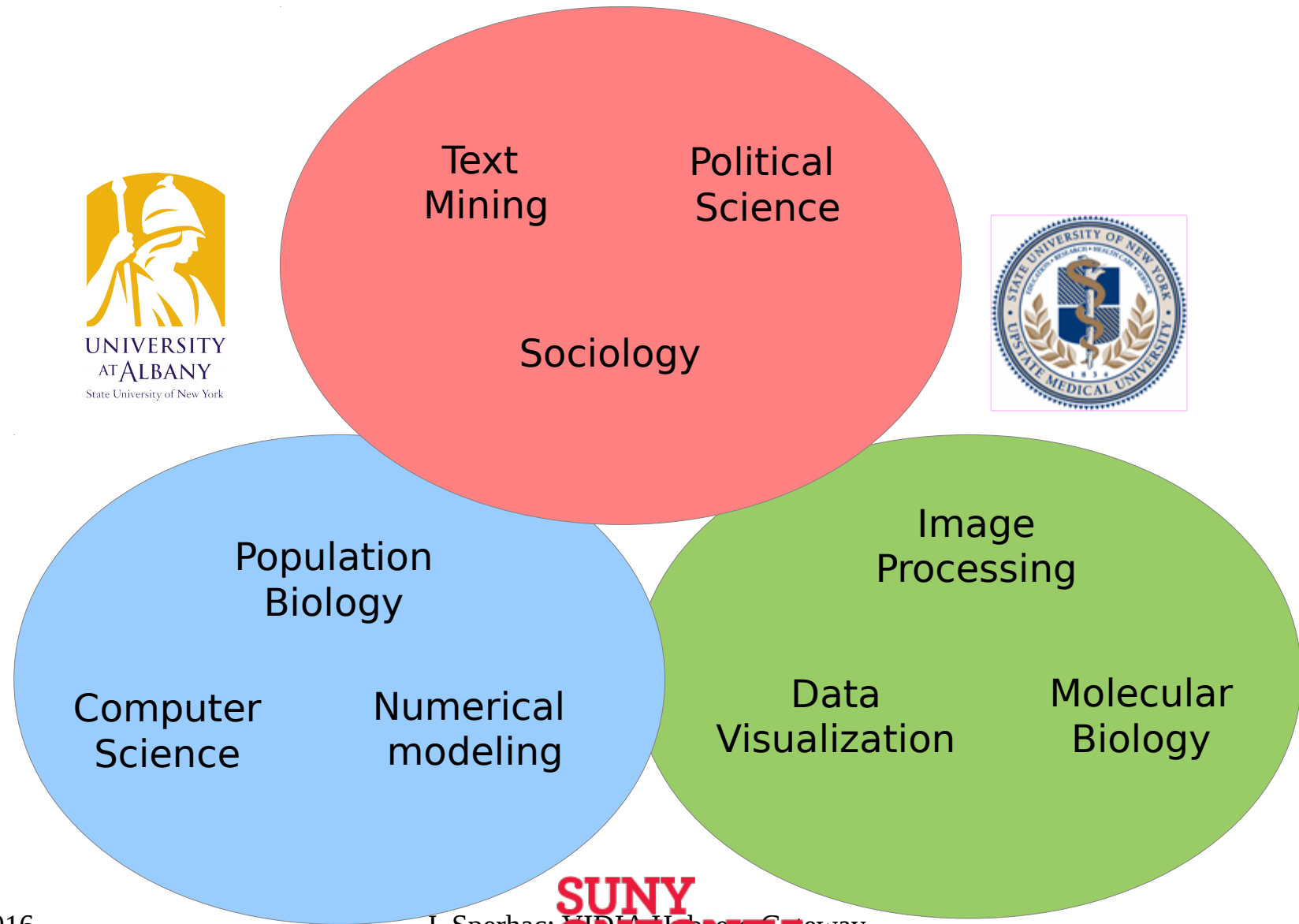
Big Data
[VIDIA] was
the best part
of the class

A SUNY Oneonta student, quoted in W. Wilkerson,
“Using ‘Big Data’ in a Political Science Research
Methods Course”

VIDIA Tool Sessions, to 2016-09



VIDIA supports research...



Enabling Research

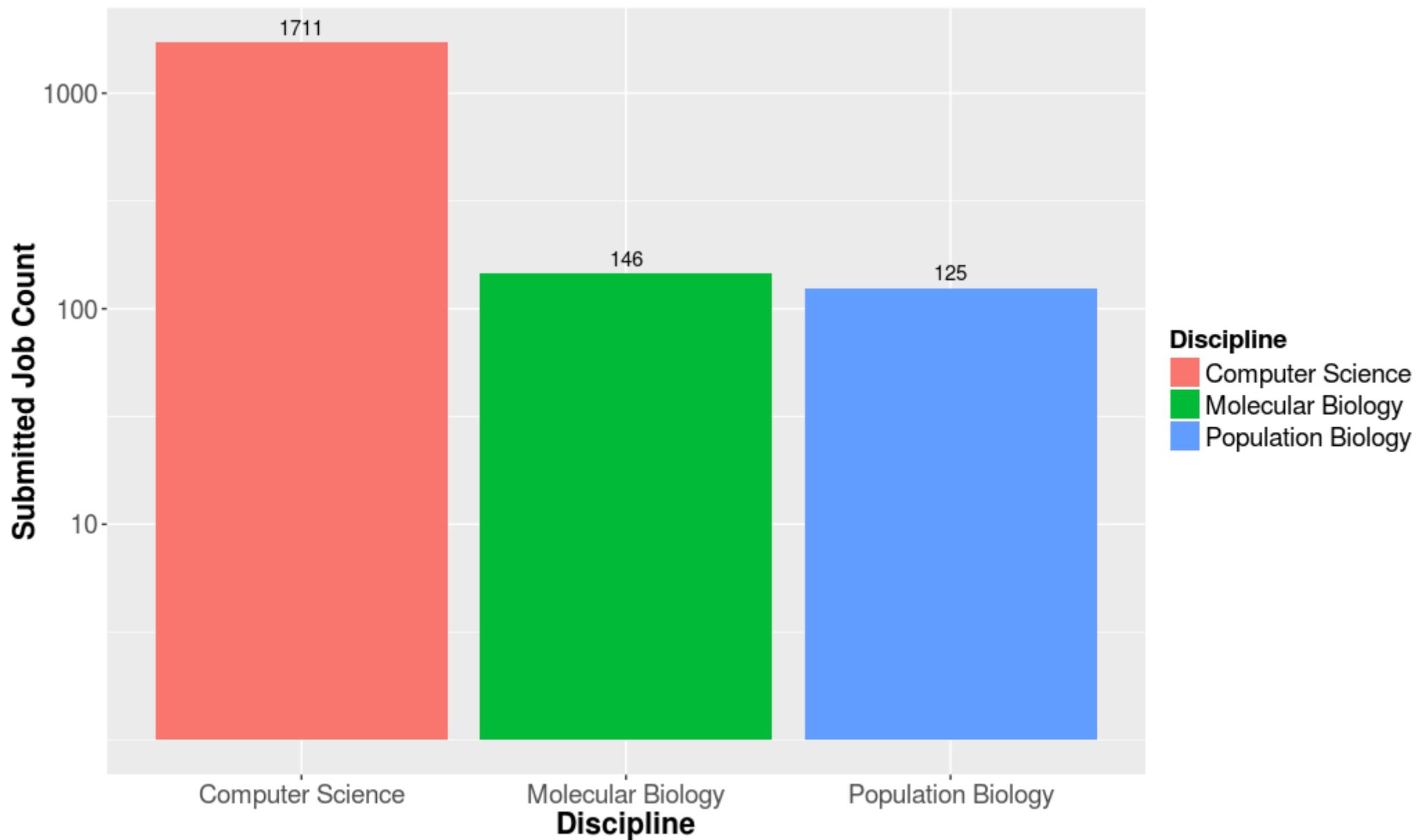
Research	Tool	Discipline
Electron Microscopy Image Analysis	HPC	Molecular Biology
Supreme Court in Public Opinion	R	Political Science
Controversies in Animal Rights	RapidMiner	Sociology
Software Anomaly Detection	HPC	Computer Science
American Shad Migration Modeling	R, HPC	Population Biology

CENTER FOR **COMPUTATIONAL RESEARCH**

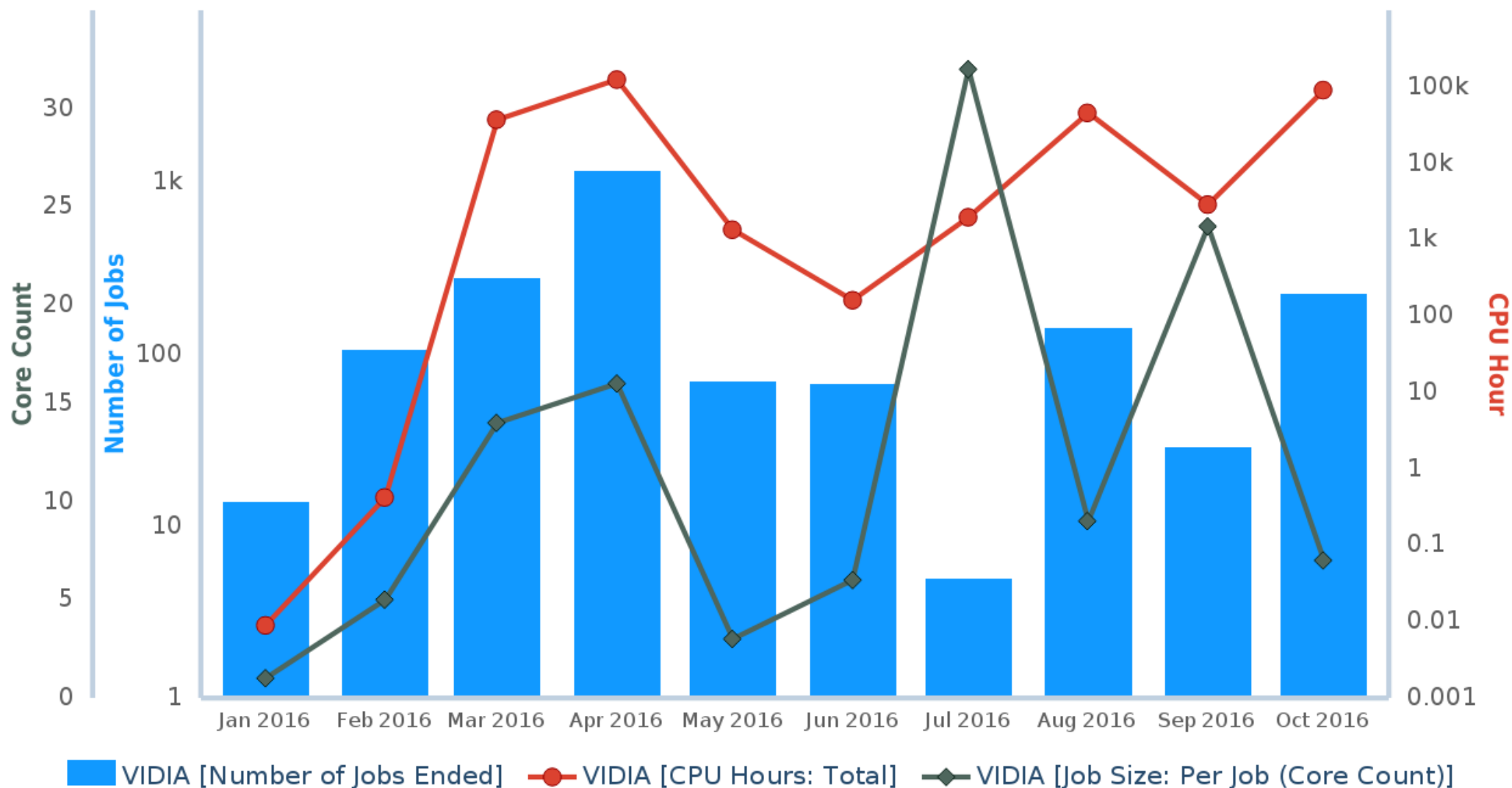


VIDIA HPC Cluster Usage

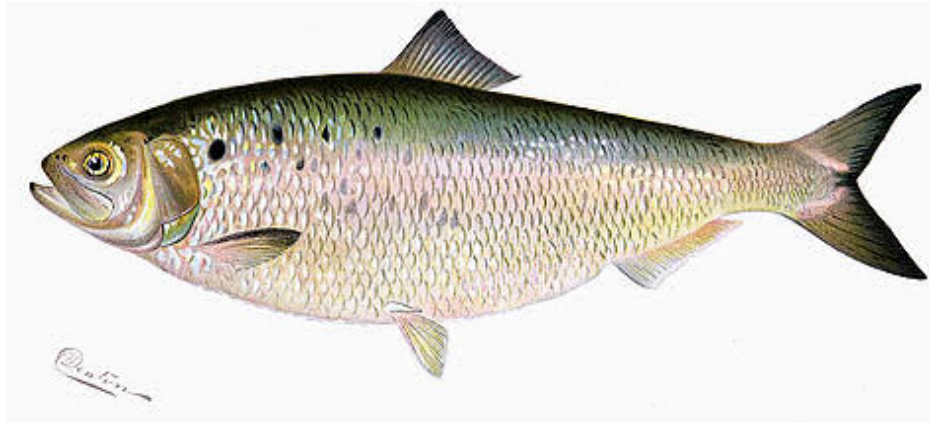
2016-01 to 2016-10



VIDIA HPC Cluster Usage 2016



American Shad Migration Model



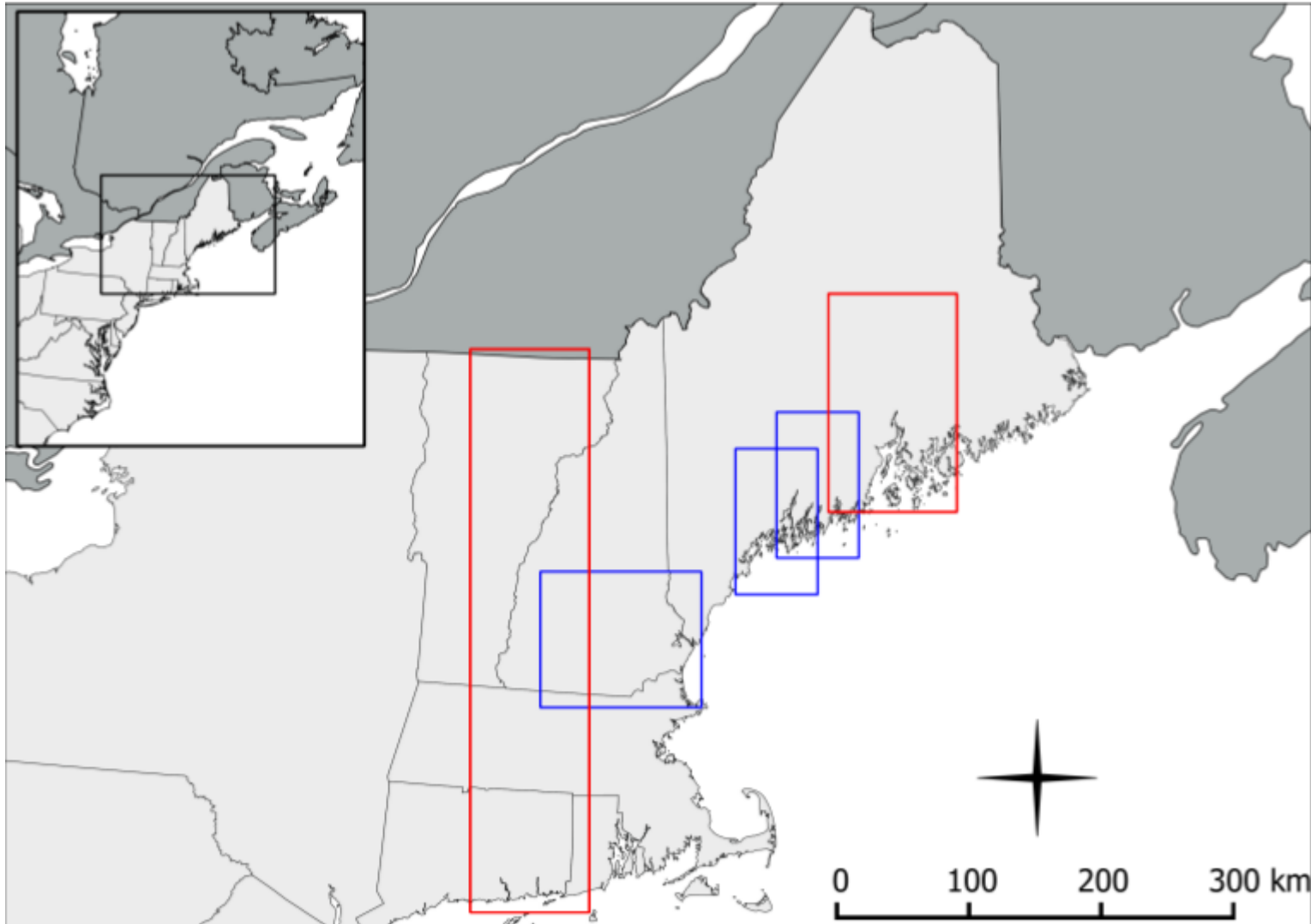
- Daniel Stich, Oneonta
- Data from literature, surveys
- Statistical models
- Advice for dam managers, fisheries

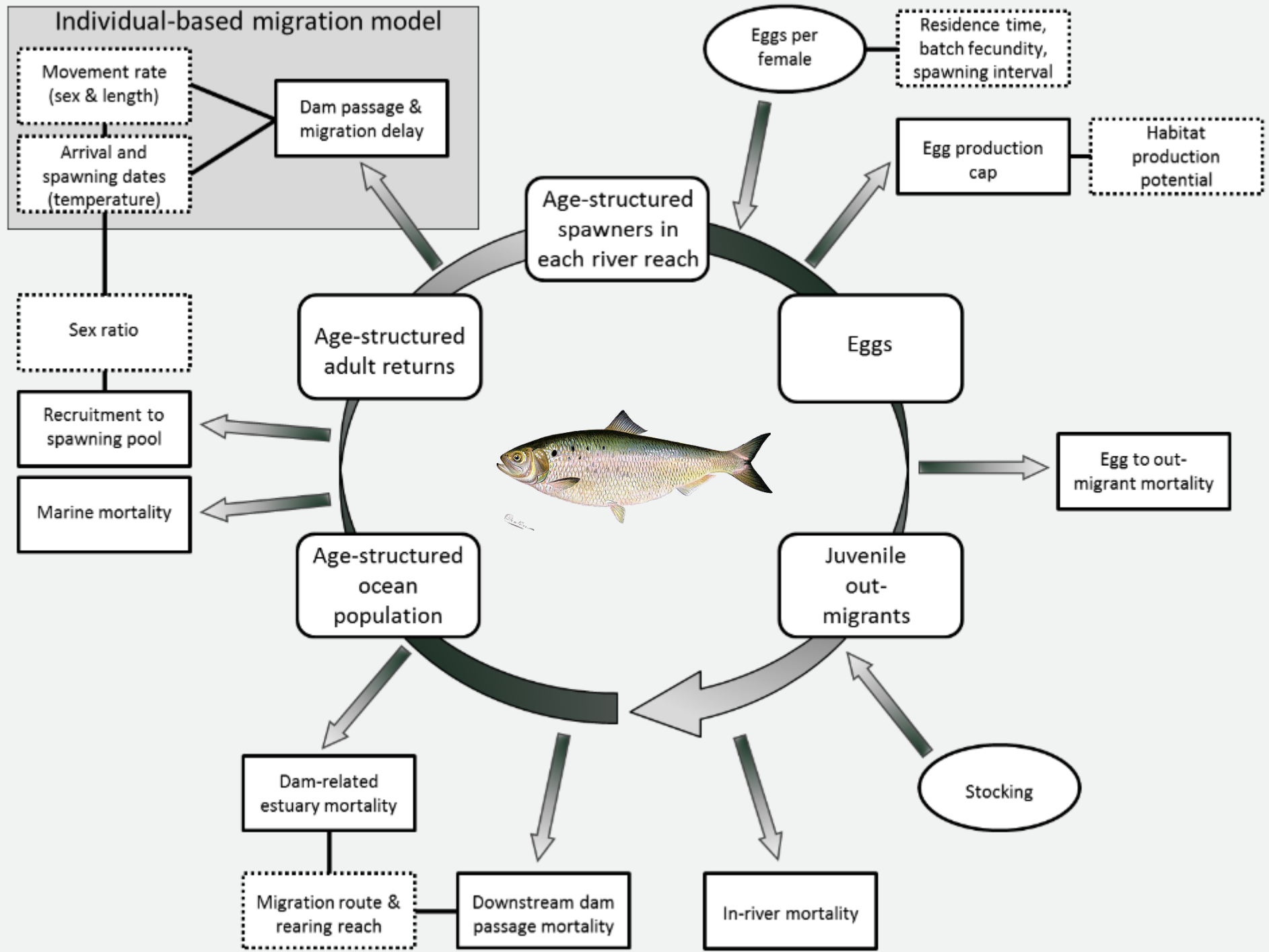
CENTER FOR **COMPUTATIONAL RESEARCH**



**SUNY
ONEONTA**

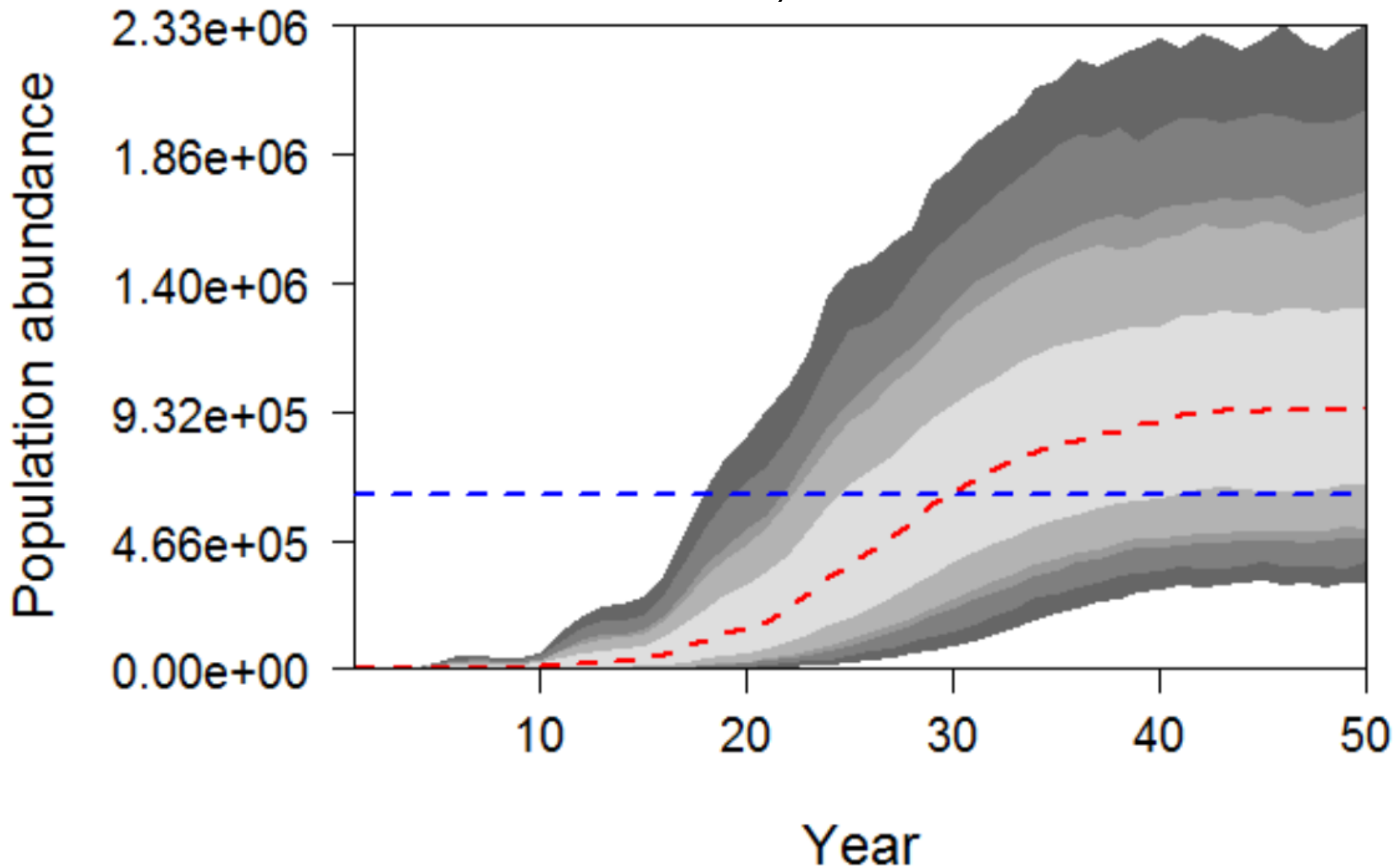
Shad migration: river systems



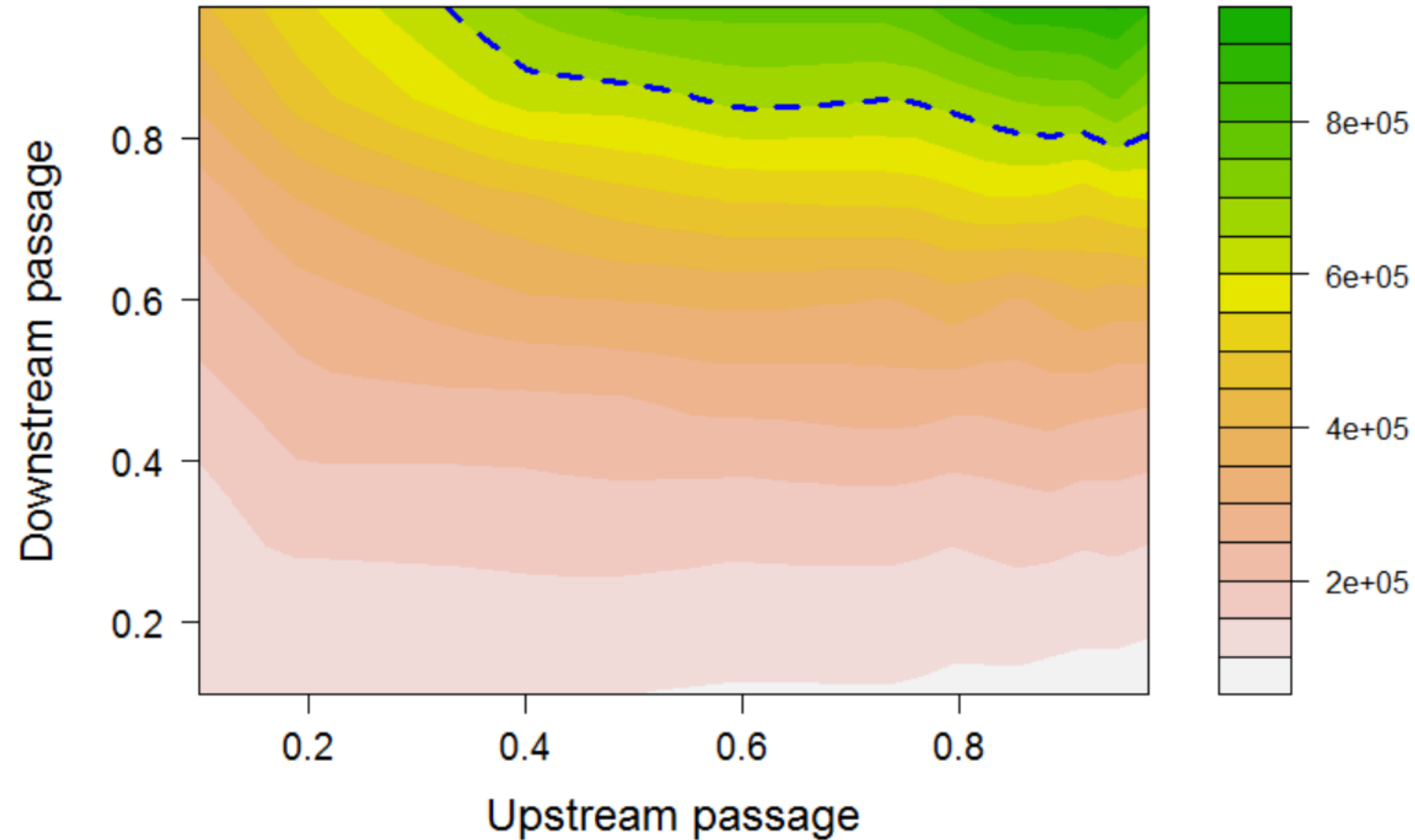


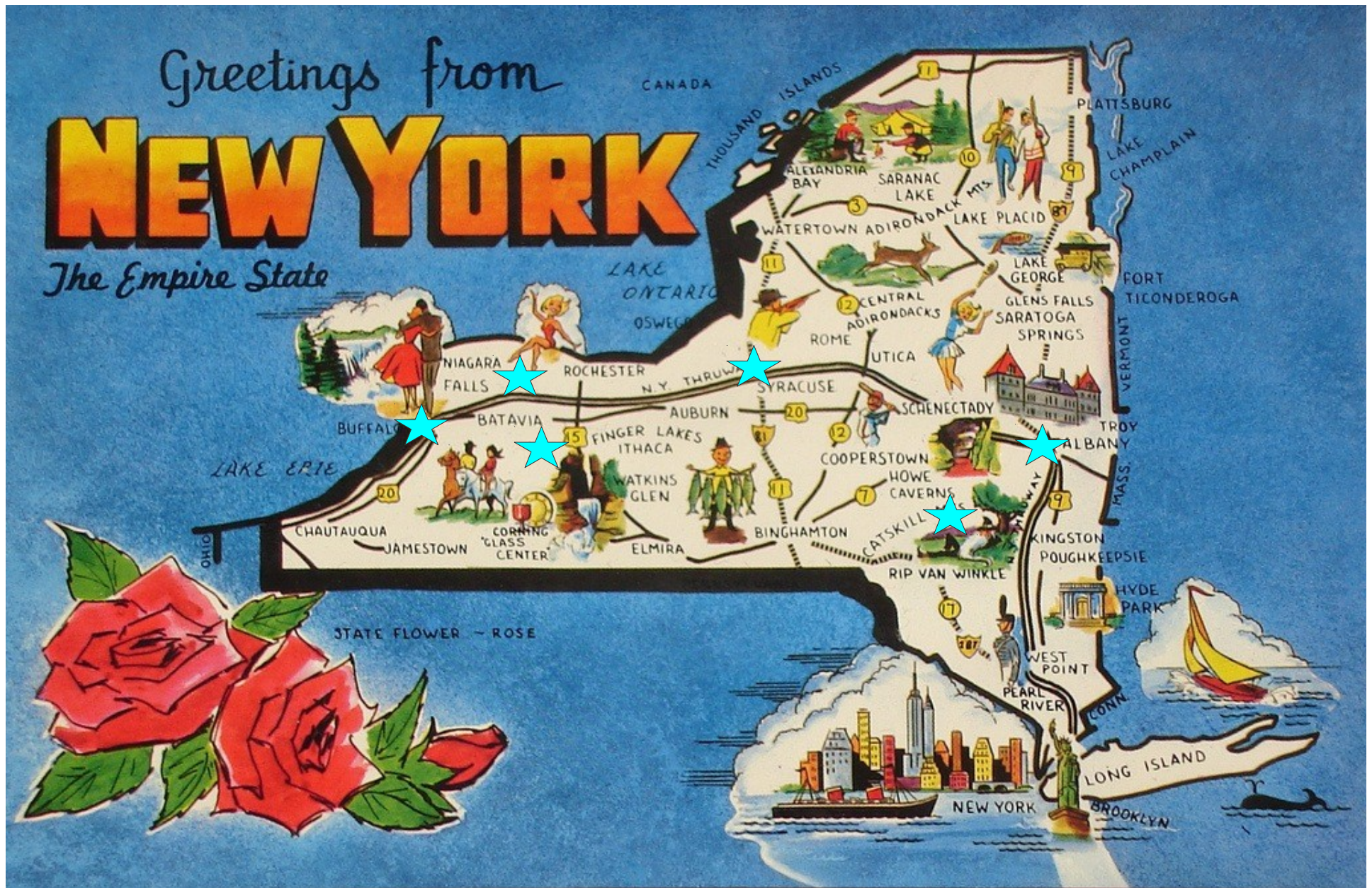
Preliminary shad model output

Penobscot River, no dams



Shad model: spawner abundance







SUNY Buffalo/CCR

- Steve Gallo
- Tom Furlani

SUNY Geneseo

- Kirk M. Anne

SUNY Albany

- Ming Ying
- Yizhen Chen

SUNY Oneonta

- James Greenberg
- Gregory Fulkerson
- Brett Heindl
- Achim Koeddermann
- Brian M. Lowe
- Bill Wilkerson
- Dan Stich

SUNY Upstate Medical

- Nicholas Stam