



Cyberinfrastructure for Al@Edge Computing

The Sage Grande Testbed

Nicola Ferrier

Northwestern Argonne Institute for Science and Engineering

Deputy Director of Sage Grande Project



6th GLOBAL RESEARCH PLATFORM WORKSHOP, Sept 15, 2025



Sage Grande: An Open Artificial Intelligence Testbed Supporting Edge Computing and Intelligent Sensing

Goals

Training and Educating the Al Workforce



First generation Sage nodes deployed 2019-2024 (MSRI-1 #1935984 2019-2024)



Provide academic community Al resources to:

- Develop new Al algorithms and technologies
- Advance techniques to enhance AI safety & privacy
- Explore next-generation Al-enabled infrastructure for real-time monitoring of wildfire, flooding, and drought to support early warning and mitigation for farms
- Deploy 300 new state-of-the-art Al-enabled platforms across the USA, and in every state
- Train students to use Al-enabled cyberinfrastructure with hackathons, workshops, and new curricula



Partnerships with American Companies

Institutional Lead: Northwestern Univ.

Partners: Univ. of California, San Diego; Univ. of Illinois Chicago, Colorado State, Univ. of Utah, Univ. of Hawaii, Univ. of Oregon, George Mason Univ.,



Sage Grande Award: #2436842

Start Date: 2/15/2025

\$25.6M over 5 years

Field Proven

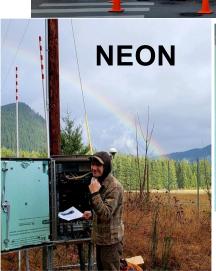
sagecontinuum.org













Almost 10 years of experience NSF MSRI-1 Pilot 2019-2024

Building on NSF Array of Things (2016-2018)



Al@Edge platform brings together sensors, actuators, and computing.



software



BYOS – bring your own sensors

- Tree sap monitor
- Hawaii rain
- Insect screen
- Lidar (Nantucket and Chicago)

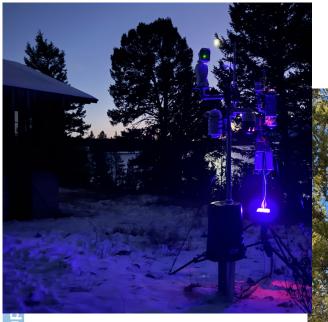




Volcano National Park October 2023



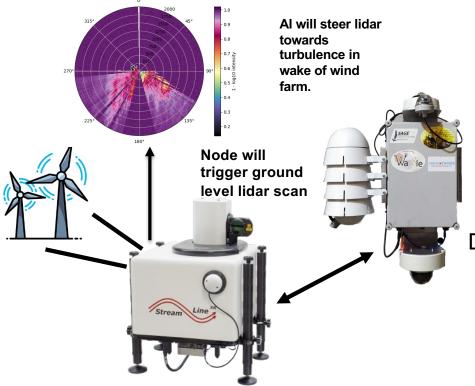
Lorawan connected tree sap monitor

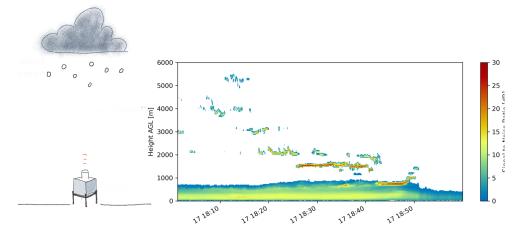






Sage + Doppler Lidar





DL is a great target for edge processing:

- Doppler spectra are not routinely stored
 - · Currently only moments are archived
- Spectra contains rich information that should be processed differently when different phenomena are present



Jackson, R. C., and Coauthors, 2023: ARMing the Edge: Designing Edge Computing—Capable Machine Learning Algorithms to Target ARM Doppler Lidar Processing. *Artif. Intell. Earth Syst.*, **2**, 220062, https://doi.org/10.1175/AIES-D-22-0062.1

A Few Examples of Al@Edge



Wildfire Detection



Jakub Szumny

Biodiversity



Dario Dematties

Urban Safety



Anagha Tiwari

Drought Agriculture



Eugene Kelly

Lidar/Digital Twins



Bobby Jackson

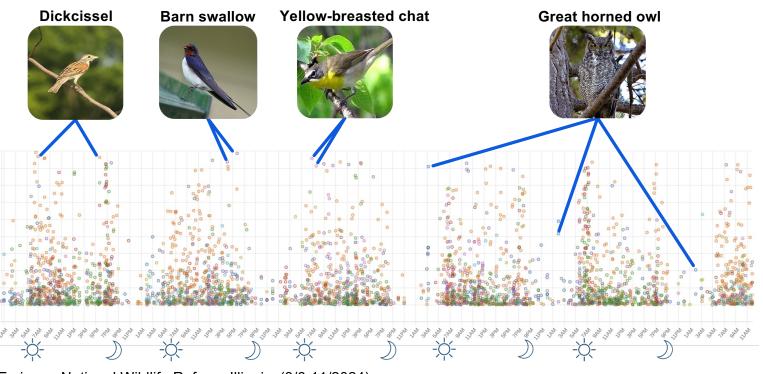
Predict Flooding



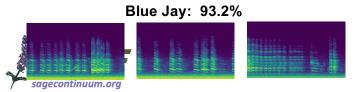


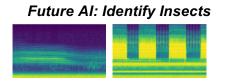
Seongha Park

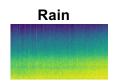




Emiquon National Wildlife Refuge, Illinois. (3/6-11/2024)





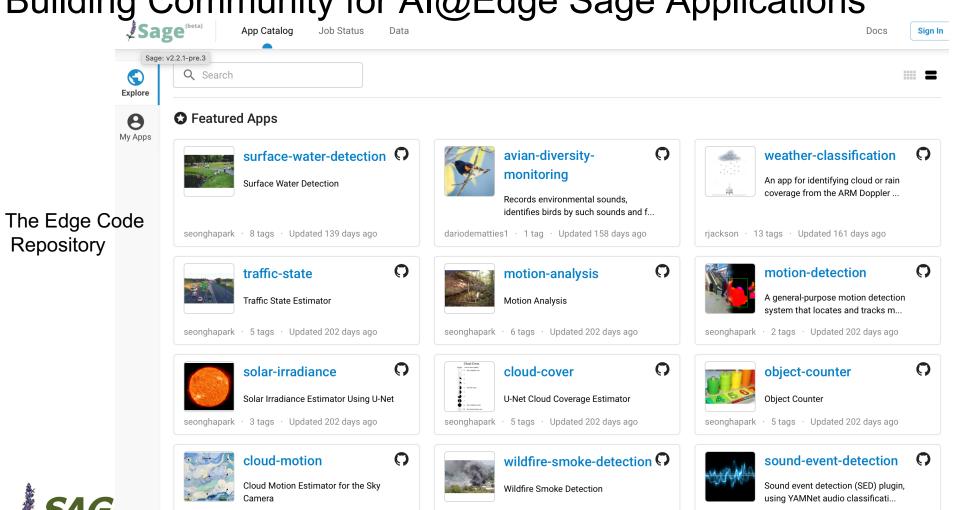








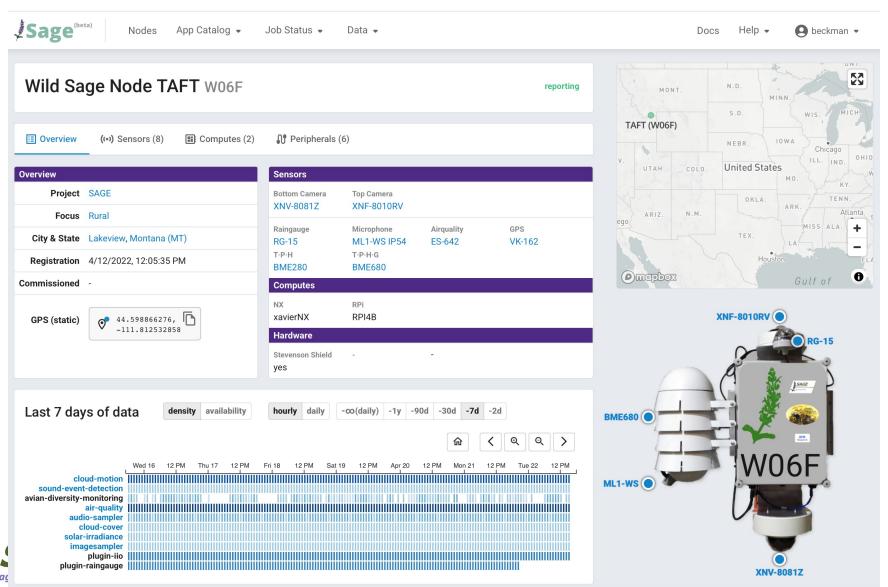
Building Community for Al@Edge Sage Applications



seonghapark · 2 tags · Updated 259 days ago

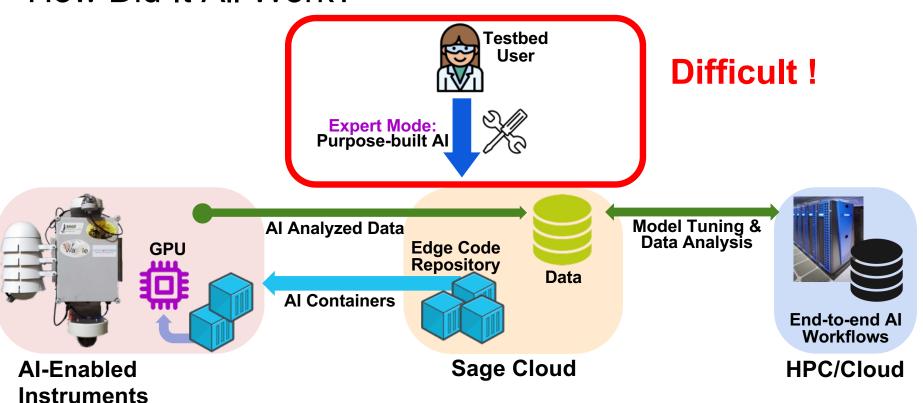
bhupendraraut · 2 tags · Updated 202 days ago

dariodematties · 1 tag · Updated 278 days ago



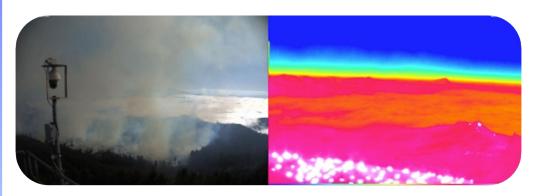


How Did It All Work?





Testbed Advance AI -- Guide what we can imagine next





National NSF Cyberinfrastructure to Support AI Research

Computer Science:

- Al inference, federated / self-supervised learning
- LLM benchmarking and correctness at edge
- Cybersecurity, privacy, real-time data analysis
- Adaptive power management
- Programming model for continuum computing
- Resource management, streaming data

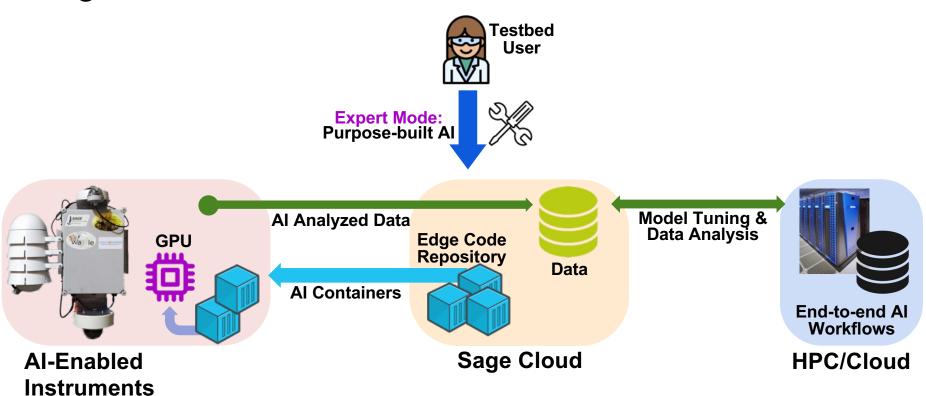
Domain Science:

- Developing next-gen Al-enabled instruments
- LLMs for in-situ data analysis
- Domain-specific tuned AI models for edge
- Understanding drought, wildfire, and flooding
- Real-time forecasts with edge-HPC models
- Adaptive steering of data collection



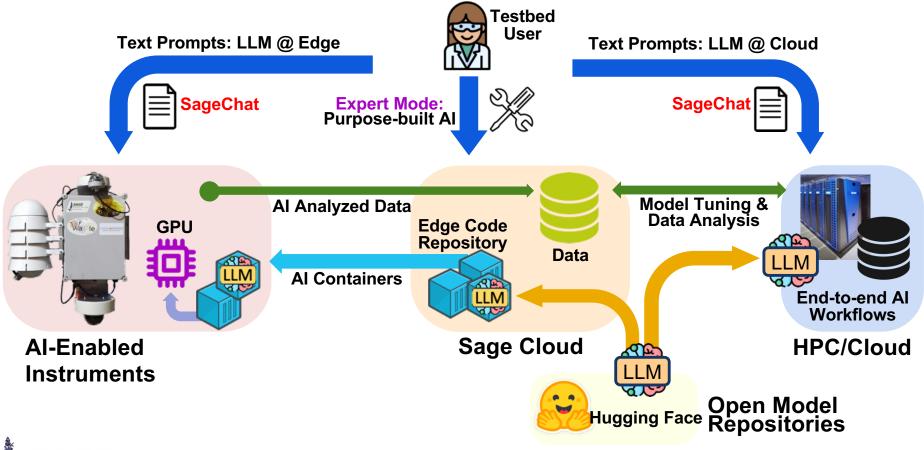


Sage v1: How *Did* It All Work?





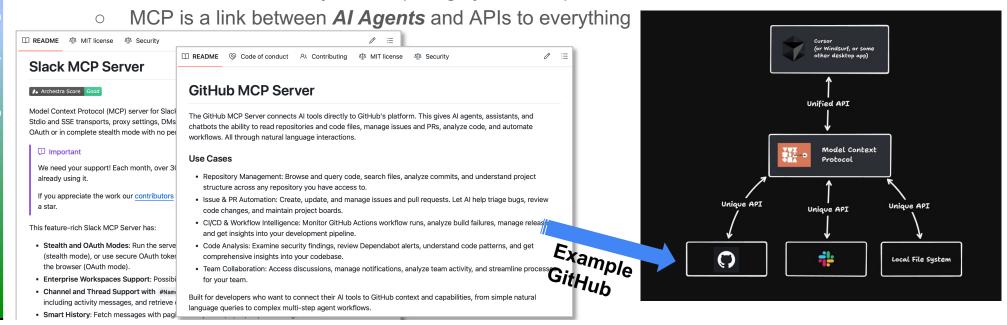
How Does (Will) It All Work?





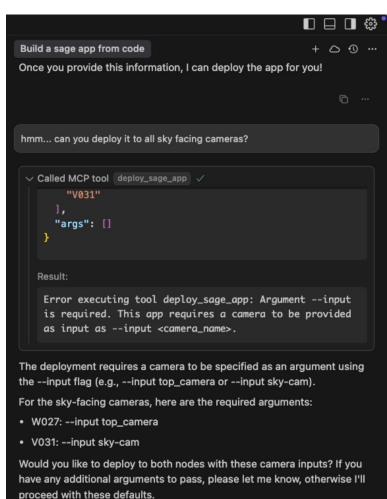
Natural Language to interact with Sage Grande Testbed: SageChat Prototype and Future Directions

- Goal: Harness LLMs --- provide natural language interfaces and expand user community for Al-enabled cyberinfrastructure and edge computing.
- Anthropic's "Model Context Protocol" (MCP) is open standard: Nov 2024
 - Provides standard way for computing systems to publish interfaces that AI can use



Sage Grande has a Prototype MCP!!!





What does the SGT MCP Provide Users? A new way to write programs and control Sage

Some Examples:

© Core Capabilities

1. Sensor Data Exploration

Discover what's available:

- "Show me all available Sage nodes and their sensors" "What environmental data is available from the last hour?" "List nodes in Chicago with temperature sensors"
- Query specific data:

"Get temperature readings from node W "Show me rainfall data from all nodes Find existing plugins: "What's the current air quality in na

3. Job Submission & Management

Submit jobs easily:

- "Deploy a cloud cover detection job to n "Run audio sampling on all nodes in Hawa
- "Start a multi-plugin ML suite on prairi

Monitor and manage:

- "Check the status of job 12345"
- "Show me recent data from my flowering p
- "Remove job 67890 from the scheduler"

Research & Data Analysis

Exploratory Data Analysis:

- "I'm studying pollinator activity. What Sage data would be m
- "Show me temperature trends across different ecosystems in t "Find correlations between flowering patterns and weather da

Hypothesis Testing:

"Compare bird activity between urban and rural nodes during "Analyze the relationship between air quality and plant heal "Test if rainfall patterns affect flowering timing in prairi

Publication-Ready Analysis:

- "Generate a statistical summary of biodiversity metrics from "Create a comprehensive environmental report for node W06D o
- "Export data in a format suitable for scientific publication

2. Plugin Discovery & Development

- "Find plugins for monitoring bird sounds"
- "What computer vision plugins are available for plant detection?"
- "Show me plugins that work with PTZ cameras"

Get development guidance:

- "How do I create a custom Sage plugin for detecting flowers?"
- "What's the best approach for real-time audio analysis on edge nodes?"
- "Help me optimize a YOLO model for edge deployment"



SageChat --- The Future

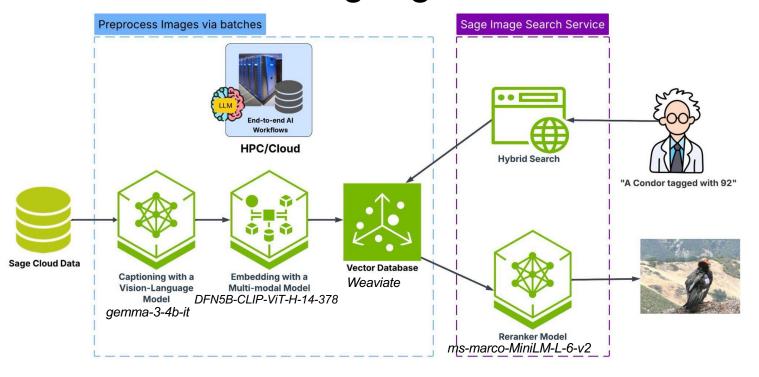
- Tutorial and educational interactions:
 - "Teach me about Al inference at the edge"
- Code and Al model development:
 - "I'd like to improve the performance of my code, either by quantizing the model or using a more efficient run-time system. Provide some step-by-step paths forward"
- Workflows spanning MCPs:
 - "I'd like to link my Al@Edge code with a HPC forecast that predicts drought. Please find the HPC resources and apply for a starter allocation and then run my workflow"
- MCP for the edge: (Agents at the Edge!)
 - I'd like to examine all incoming data at the edge for anomalies and then autonomously modify the data sampling protocol and redirect cameras and instruments to better collect data"

Resources:

- https://github.com/waggle-sensor/sage-mcp/blob/main/docs/EXAMPLES.md
- https://github.com/waggle-sensor/sage-mcp/tree/main



Process Archived Data with Al Use Natural Language Search Interface



(as of June, 2025)

Preprocess Computation:

- 10 sec/image.
- Sage archive: ~33M images
- 15K images/day
- ~382 weeks of H100 (unparallelized).
- scales by Images/Nodes

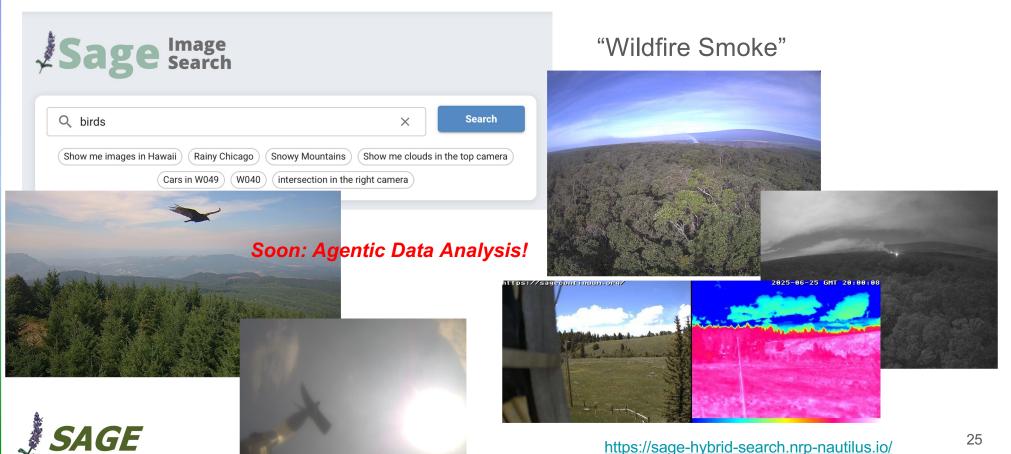
Sage Image Search Service Computation:

- 20 sec/query for top 25 images.
- Sage has ~303 users
- If all users made one query, it would take ~1.68 hours of H100 to process all. (not yet parallelized).
- · Scales by Users.



Al to Make Data Useful and Broadly Accessible to Scientific Communities

Early Prototype that shows the promise of AI for data insights



Integrating Sage Grande with other NSF Resources: NRP/NDP, Pelican/OSDF, PATh (HT Condor)

