# Leveraging Public Clouds for DOE Environmental Streaming Data
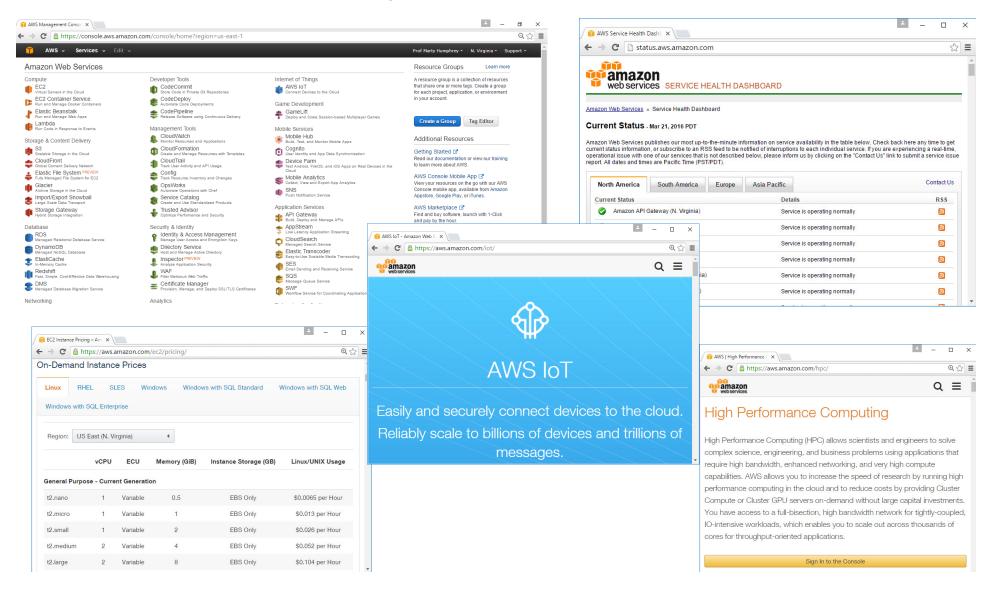
Marty Humphrey

Dept of Computer Science

University of Virginia


Jon Goodall

Dept of Civil and Environmental Engineering
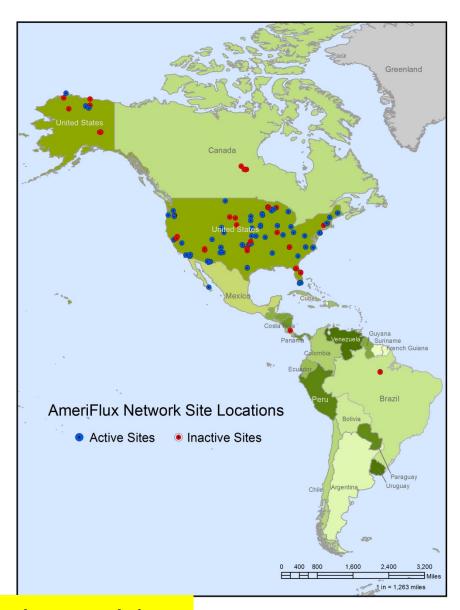
University of Virginia

# Public Clouds should be utilized MORE by Scientists!

# Many DOE applications emerging with environmental streaming data

- AmeriFlux
- NGEE Tropics
- Drone-based sensors
- Environmental monitors in cities
- Traffic sensors
- Etc.

# AmeriFlux, circa 2012



AmeriFlux Network Site Locations

- Active Sites
- Inactive Sites

# Science objectives

- Quantify exchange of carbon, water and energy between terrestrial ecosystems and the atmosphere across a range of vegetation types, disturbance histories, and climatic conditions.

- Understand processes governing the terrestrial carbon cycle and linkages with the water, energy and nitrogen cycles.

- Produce a high-quality data base and synthesize observations across the network.



Courtesy Davis et al '11

# Core measurements

- Fluxes of $CO_2$, water vapor, and sensible heat flux via eddy covariance.

- Radiative fluxes and micrometeorological conditions.

- Biophysical characterization of sites (e.g. vegetation age and type, nutrient status, carbon pool sizes, soil type).



Courtesy Davis et al '11

# AmeriFlux and Streaming Data

- Wind (direction and speed) and trace gas concentrations (mostly $CO_2$ and $H_2O$, but also $CH_4$, NO, $NO_2$, $N_2O$, and others) are measured and stored usually at 10Hz

- Separate mechanism from "data uploads"
  - Currently only tower-driven SCP (for "high-frequency data")
  - Currently only archival in nature
  - 35 configured; 10 active

# AWS IOT

- AWS Lambda: lightweight event-driven programming
- AWS Kinesis: real-time, scalable streaming data sink
- AWS S3: scalable, reliable object store
- AWS DynamoDB: managed noSQL service
- Etc.
- *Plus any open-source projects as needed*
  - *Note to Twitter: please open-source Heron (!)*

- Example: Intel Edison-based rain sensors/gauges (UVa)

# Issues

- How much streaming data is "too much" for public clouds?
- Single custom-build device (e.g., "AmeriFlux AWS IOT device") or integration with existing infrastructure?
- How much info needed for researcher to use site's streaming data?
- How to balance "site ownership" of streaming data vs. real-time nature of the data?
- *Large-scale software design, deployment, and management*