# Mining Massive Datasets from the Allen Telescope Array

Proposal for Stanford CS341 students to use IBM Cloud services to analyze massive datasets from the Allen Telescope Array at Hat Creek Radio Observatory

# IBM Spark@SETI Backgrounder

❑ The SETI Institute operates the Allen Telescope Array (ATA) to observe star systems for radio signals which may provide evidence of extraterrestrial intelligence.

❑ IBM collaborating with the SETI Institute to use IBM Apache Spark services to analyze radio signal data from the ATA.

❑ Astronomers and data scientists from around the world are using the Spark@SETI environment to analyze millions of signal events collected from the ATA over 10 years.

❑ Apache Spark with Jupyter Notebooks are proving to be highly effective in enabling experimental approaches to the analysis of the SETI Institute's archive of signal data.

❑ IBM is offering to make the entire Spark@SETI environment available to Stanford CS341 students during the Spring 2016 quarter.
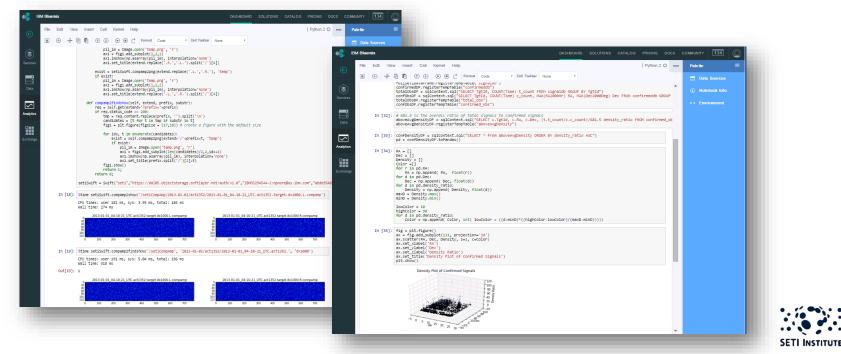
# Spark@SETI for CS341

❑ Benefits of this proposal include:

   ❑ Simple access to the Apache Spark platform to help expand the scope of CS341 projects beyond MapReduce computational methologies.

   ❑ Use of the IBM Spark service at no charge, accessible on the IBM Cloud using a standard web browser.

   ❑ Unlimited access to the ATA data archives, stored on the IBM cloud.

# CS341 Students and Spark@SETI

Students using Spark@SETI for their CS341 projects will receive:

- ❑ A free account on the IBM BlueMix preconfigured with Spark services and direct connections to the ATA data repository.

- ❑ A "kick start" portfolio of Jupyter Notebooks that provide all the foundational Python code to query the ATA databases and read the recorded binary signals for analysis.

- ❑ Access to a 200 million row relational database of ATA signal events

- ❑ Access to approximately 15 million binary "complex amplitude" files that store the raw signal data from the ATA at the moment that a signal was detected

- ❑ Support from the Spark@SETI team for issues relating to platform usage (e.g. dropped Python kernels)

IBM

SETI INSTITUTE

# CS341 Project Goals

❑ CS341 projects using Spark@SETI should not restrict project goals strictly to identifying a signal of interest.

❑ Much of the ongoing work in the Spark@SETI initiative is focused on the inverse problem: that of improving signal classification methods to eliminate signals that are human radio frequency interference (RFI).

❑ CS341 students may produce novel and valuable analytic results, which will be added to the Spark@SETI repository of Spark notebooks for use by astronomers and data scientist from around the world.

❑ Example: Igor Nikitin at the Fraunhofer Institute for Algorithms and Scientific Computing used ATA data to demonstrate how to apply the Radon transform to greatly improve the signal-to-noise ratio of narrow-band signals, thereby permitting the detection and classification of many new signals.

# CS341 Spark@SETI - Examples

Examples of potential CS341 student projects include:

❑ Supervised machine learning of spectrogram images to classify signals according to known categories  (e.g. radar interference, narrow-band zero drift, aircraft RFI).

❑ Analysis of the 200 million signal event database for targets that show unusual consistency of signal features (e.g. consistent power) over long periods of time. Outliers could be reviewed by the radio telescope operations staff for potential re-observation.

❑ Extracting scalar features from signal recordings with PCA and/or MDA results that indicate these features will help to spread and segment signals in useful ways for further analysis by Spark@SETI scientists.

IBM

SETI INSTITUTE