

# The MAterials Simulation Toolkit for Machine Learning (MAST-ML): Automating Development and Evaluation of Machine Learning Models for Materials Property Prediction

**Ryan Jacobs, Tam Mayeshiba, Ben Afflerbach, Dane Morgan**  
*(University of Wisconsin – Madison, WI USA)*

**Luke Miles, Max Williams, Matthew Turner, Raphael Finkel**  
*(University of Kentucky, Lexington, KY USA)*

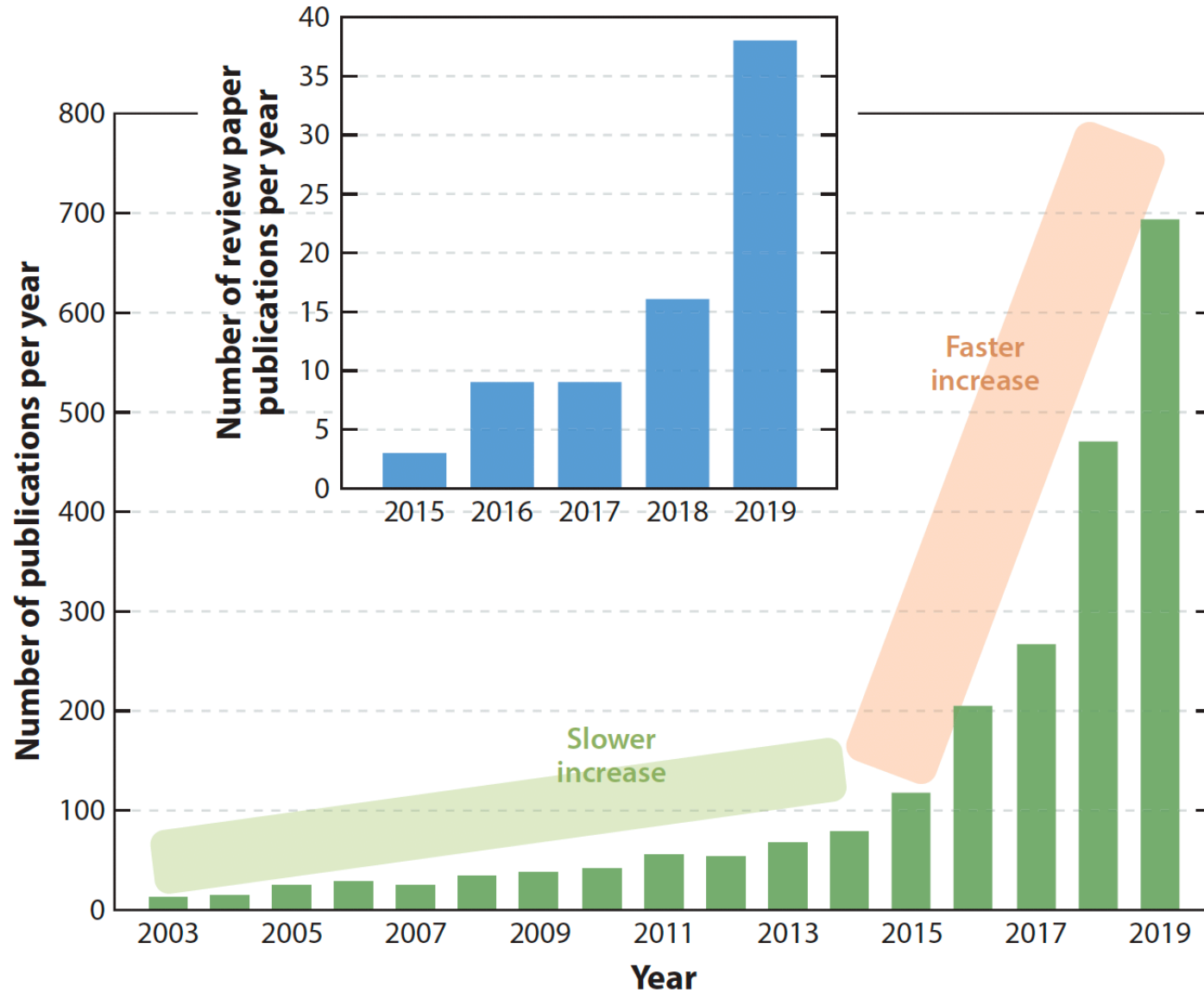
*Most Recent Skunkworks MASTML members:*  
**Avery Chan, Hock Lye Lee, Min Yi Lin**

<https://github.com/uw-cmg/MAST-ML>

**NanoHub ML Workshop**  
**5/19/2021**

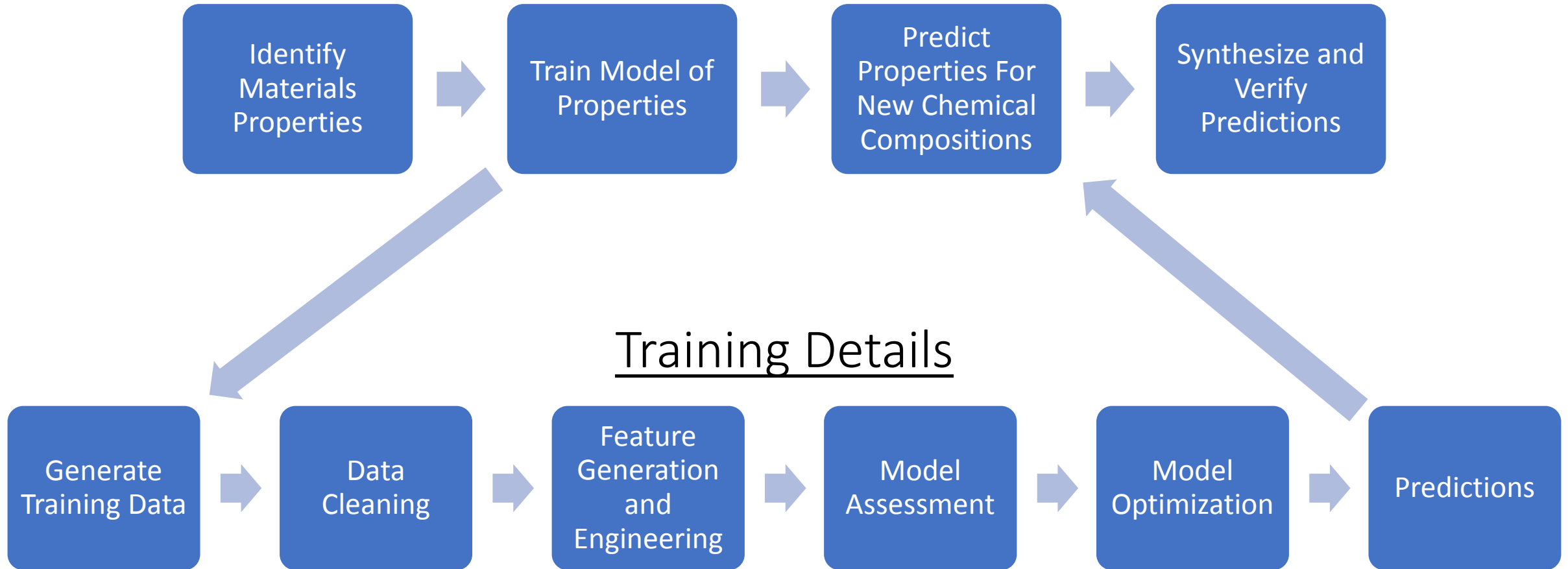


# Machine learning in Materials Science is Exploding



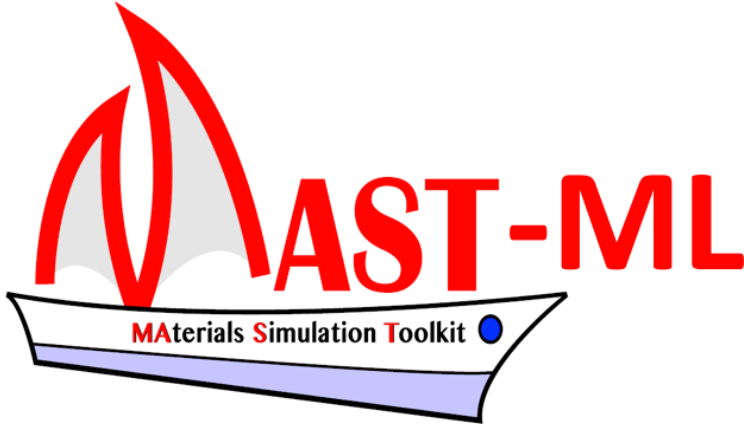
# A Basic Materials Design Workflow

---



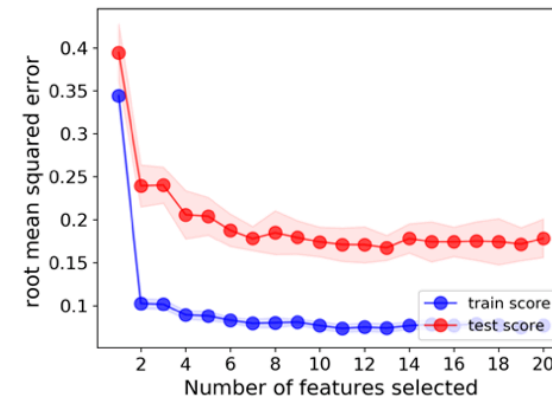
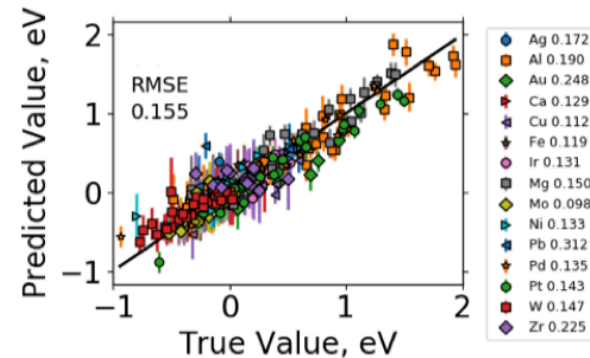
# What is MAST-ML?

MAST-ML is an open-source Python package designed to broaden and accelerate the use of machine learning in materials science research, particularly for non-experts.

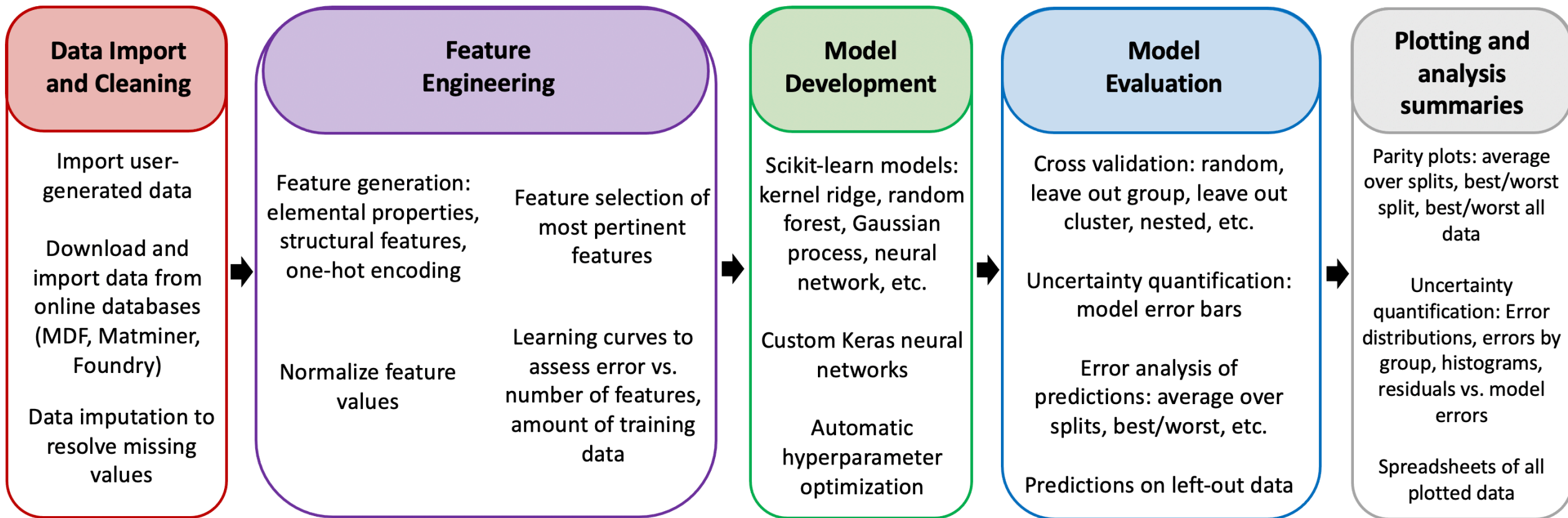


Automated machine learning tools for materials informatics research (MAST-ML)

<https://github.com/uw-cmg/MAST-ML>

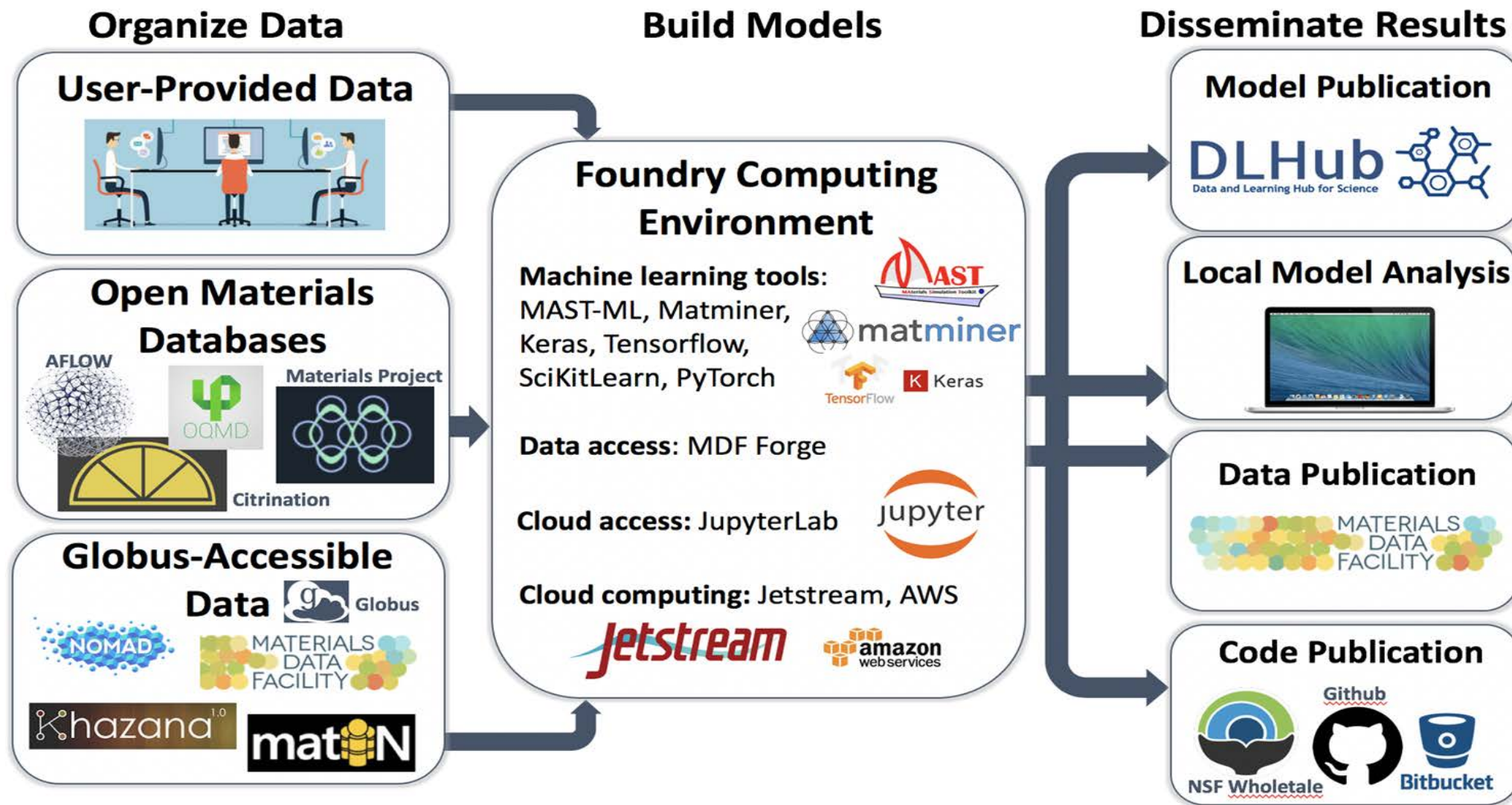


# MAST-ML automates the supervised learning workflow

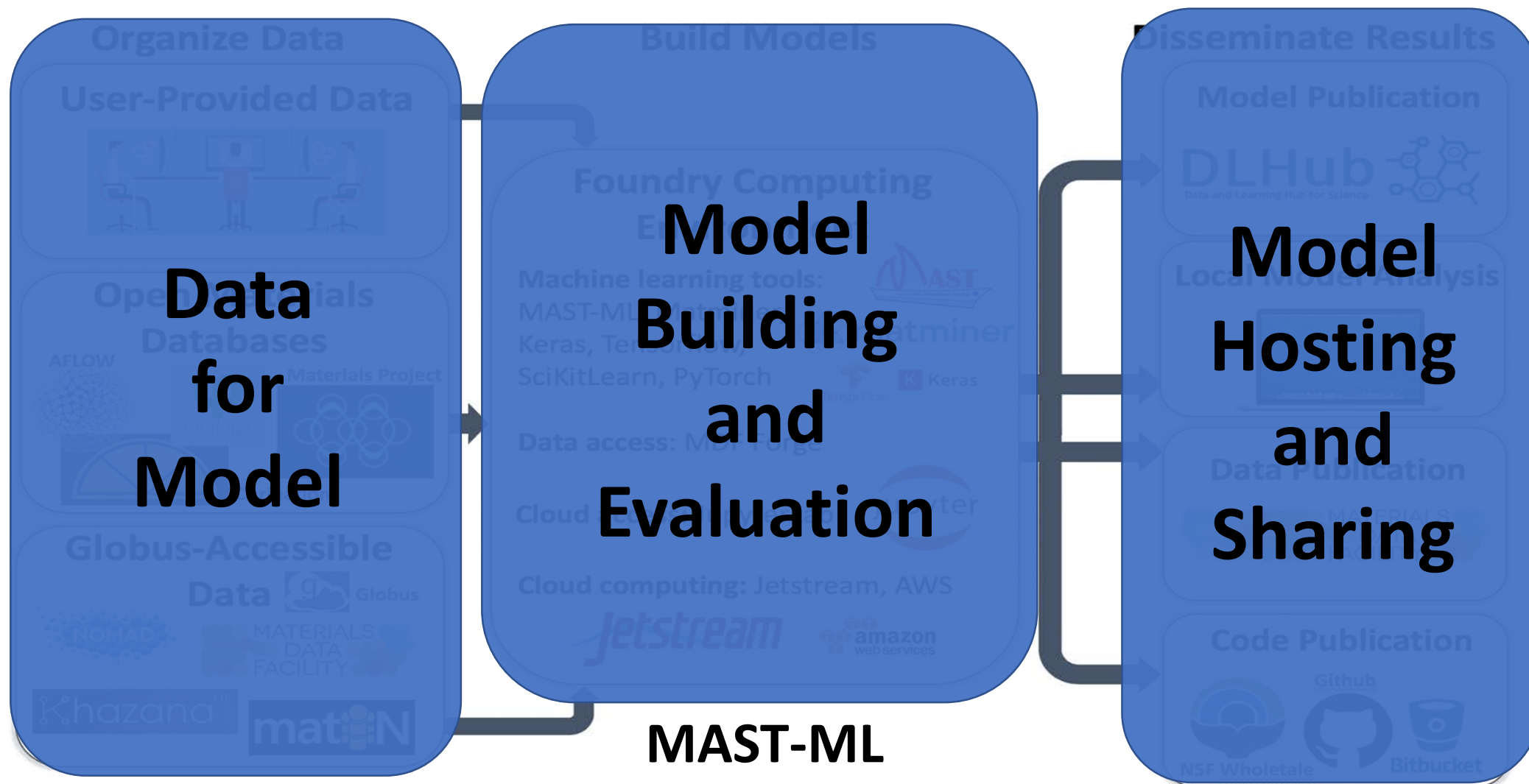


- MAST-ML supports the full library of scikit-learn modules, and can be used to construct neural networks with Keras (based on tensorflow)
- MAST-ML allows for the simultaneous execution of an arbitrary combination of data preprocessing, feature generation/selection, model types and model evaluation metrics

# (NSF CSSI) Machine Learning Materials Innovation Infrastructure



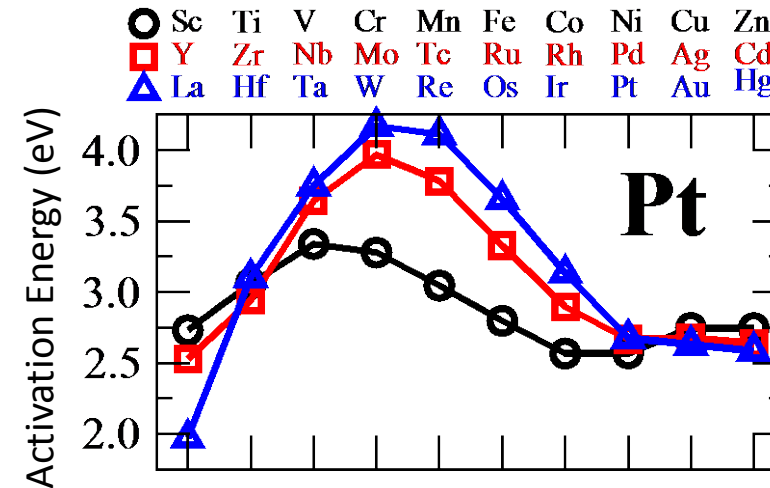
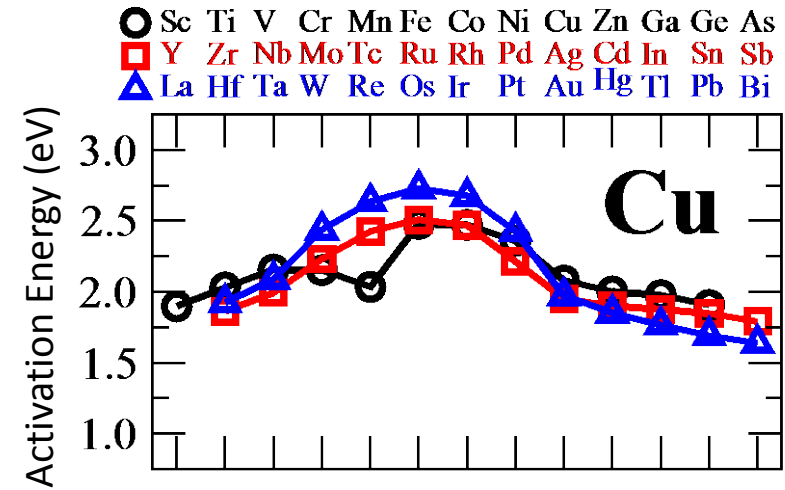
# (NSF CSSI) Machine Learning Materials Innovation Infrastructure



Model building, evaluation, and key connections  
between data and model dissemination

# Test Problem: Impurity Diffusion Database

- Diffusion of dilute impurity X in host H. We have DFT calculations of 440 values, but want ~4,000. [1, 2]
- Assume Y= Activation energies measured relative to host, X= Host descriptors, Impurity descriptors. Find  $Y=F(X)$ .
- Descriptors = elemental properties like melting temperature, bulk modulus, electronegativity, ... and their ratios, differences, etc. (MAGPIE set)[3]
- F is determined using standard machine learning regression methods (e.g., Gaussian Process Regression (Gaussian Kernel) (GPR), Random Forest (RF), neural network).
- Fit F with calculated data (15 hosts, 440 M-X pairs)

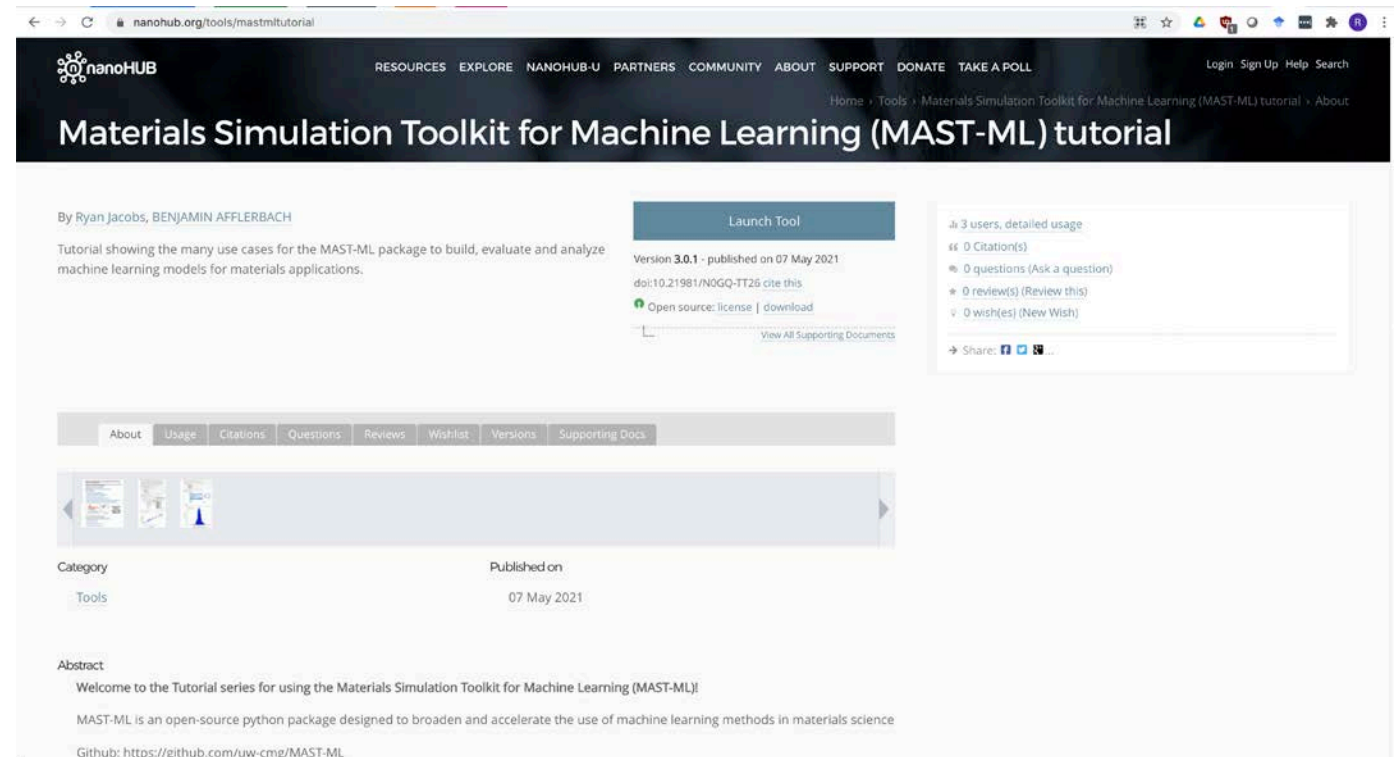


<http://diffusiondata.materialshub.org/>



# Getting Started with the MAST-ML tutorial on NanoHub

- Link to Tool:  
<https://nanohub.org/tools/mastmltutorial>
- Select “Launch Tool”
- A Jupyter notebook environment will open (may take a minute)
- Click on cell and run with Shift+return
- Data will be saved to local directory, see next slides for how to download results



The screenshot displays the NanoHub website interface for the MAST-ML tutorial. The page title is "Materials Simulation Toolkit for Machine Learning (MAST-ML) tutorial". It is authored by Ryan Jacobs and Benjamin Afflerbach. The page includes a "Launch Tool" button, version information (3.0.1, published 07 May 2021), and a DOI (10.21981/NOGQ-TT26). There are also links for "Open source: license" and "download". A sidebar on the right shows user statistics: 3 users, 0 citations, 0 questions, 0 reviews, and 0 wishes. The main content area has tabs for "About", "Usage", "Citations", "Questions", "Reviews", "Wishlist", "Versions", and "Supporting Docs". Below the tabs, there is a category "Tools" and a "Published on" date of "07 May 2021". The abstract section begins with "Welcome to the Tutorial series for using the Materials Simulation Toolkit for Machine Learning (MAST-ML)!" and describes MAST-ML as an open-source Python package for machine learning in materials science. The GitHub link is provided as <https://github.com/uw-cmg/MAST-ML>.