# Citrine Informatics

The data analytics platform for the physical world

## The Latest from Citrine

Summit on Data and Analytics for Materials Research

31 October 2016

# Our Mission is Simple

**Add as much value
to your work as possible,
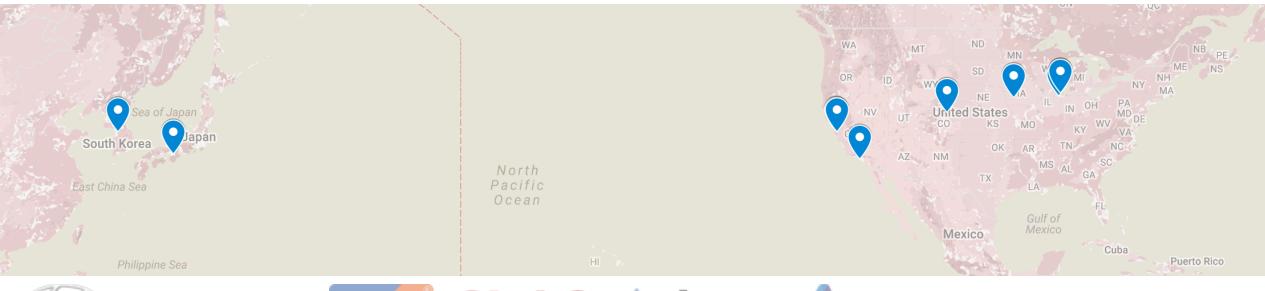immediately,
using data**

# Keys to Industrial Relevance

**UBIQUITY**

**EASE OF USE**

**OBVIOUS ROI**

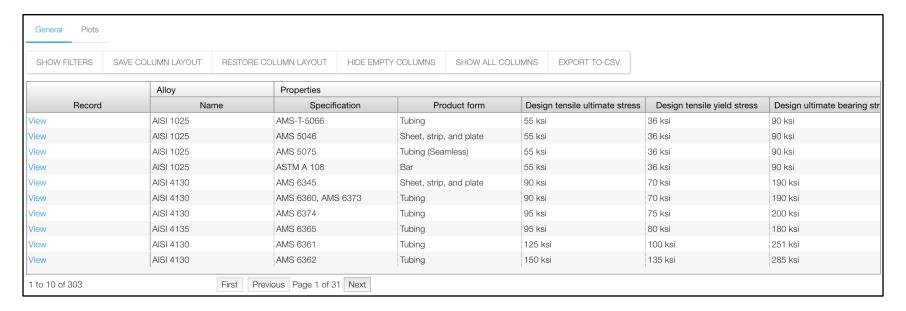# Citrine is the community cloud for materials data, predictive models, & post-processing

+ **<u>All</u> relevant data** in one place, unified from databases, research groups, papers

+ **Predictive AI**, physics-based simulations, and post-processing tools seamlessly integrated with the data

+ **Vibrant ecosystem** of researchers and developers
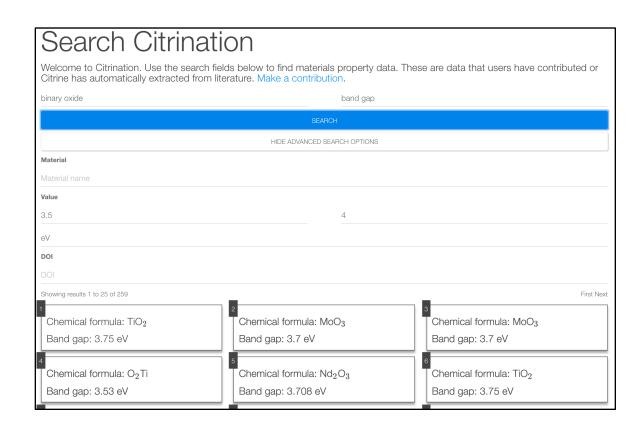
# All Relevant Data

17m+ free data records as pif's on citrination.com (& API)

ASM and MMPDS are now official data partners, providing premium data to the platform; 6 free NIST SRD's & much more

| | General | Plots | | | | |
|---|---|---|---|---|---|---|

| SHOW FILTERS | SAVE COLUMN LAYOUT | RESTORE COLUMN LAYOUT | HIDE EMPTY COLUMNS | SHOW ALL COLUMNS | EXPORT TO CSV |
|---|---|---|---|---|---|

| | Alloy | Properties | | | | |
|---|---|---|---|---|---|---|
| Record | Name | Specification | Product form | Design tensile ultimate stress | Design tensile yield stress | Design ultimate bearing str |
| View | AISI 1025 | AMS-T-5066 | Tubing | 55 ksi | 36 ksi | 90 ksi |
| View | AISI 1025 | AMS 5046 | Sheet, strip, and plate | 55 ksi | 36 ksi | 90 ksi |
| View | AISI 1025 | AMS 5075 | Tubing (Seamless) | 55 ksi | 36 ksi | 90 ksi |
| View | AISI 1025 | ASTM A 108 | Bar | 55 ksi | 36 ksi | 90 ksi |
| View | AISI 4130 | AMS 6345 | Sheet, strip, and plate | 90 ksi | 70 ksi | 190 ksi |
| View | AISI 4130 | AMS 6360, AMS 6373 | Tubing | 90 ksi | 70 ksi | 190 ksi |
| View | AISI 4130 | AMS 6374 | Tubing | 95 ksi | 75 ksi | 200 ksi |
| View | AISI 4135 | AMS 6365 | Tubing | 95 ksi | 80 ksi | 180 ksi |
| View | AISI 4130 | AMS 6361 | Tubing | 125 ksi | 100 ksi | 251 ksi |
| View | AISI 4130 | AMS 6362 | Tubing | 150 ksi | 135 ksi | 285 ksi |

1 to 10 of 303    First   Previous   Page 1 of 31   Next

# Graphical & API (Semantic) Search



"Show me binary oxides with band gap between 3.5 and 4 eV"

# Open Data Matters

# Open Data Matters



JOM
August 2016, Volume 68, Issue 8, pp 2116–2125

**Semi-Supervised Approach to Phase Identification from Combinatorial Sample Diffraction Patterns**

Authors     Authors and affiliations

Jonathan Kenneth Bunn, Jianjun Hu, Jason R. Hattrick-Simpers

"In the current implementation, SS-AutoPhase (semi-supervised AutoPhase) was used to phase map 278 diffractograms from a FeGaPd "open-data" combinatorial thin-film library.**[citation for Citrination]**
…
In this study, the open FeGaPd structural data not only allowed for the validation of SS-AutoPhase, but also it enabled a **new materials discovery from data produced >10 years ago**. By making these data open, the value of the data to the materials community was increased."

# Value of Data Scale in Practice



Initial dataset too small for signal →Larger training set via Citrine platform **Predictive model drove real-world discovery**

# The Citrine Predictive Approach

Start with known physical and chemical relationships

*(priors = DFT ground states, CALPHAD simulations, design rules…)*

then

fit remaining variance to reality (huge quantities of relevant measurements) with machine learning

# Predictive Artificial Intelligence for Materials

Collaboration with Computherm to demonstrate benefits of CALPHAD data in training AI to predict Al alloy mech properties

**AI without CALPHAD**
*RMSE = 82 MPa*



**AI with CALPHAD**
*RMSE = 61 MPa*

# Dataset Visualization



**Scatterplot of UCSB thermoelectrics dataset**
Gaultois *et al.*, *Chem Mater* **25** (2013)

# Dataset Visualization



**Citrine platform recreates visuals from the paper interactively**

Gaultois *et al*., *Chem Mater* **25** (2013)

# Dataset Visualization

**Dynamic Ashby plot of commercial 3D printing materials**

# Uncertainty Quantification

## All Models Have Error Bars



## Predictions are Distributions

# Feature Selection & Importance

**Magpie feature set**
bitbucket.org/wolverton/magpie
doi:10.1038/npjcompumats.2016.28

We are working with the informatics community to build a comprehensive library of *all published features*

Important Features

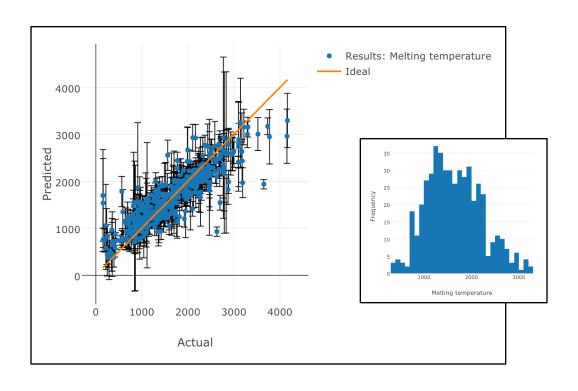| Seebeck coefficient | |
|---|---|
| CHEMICAL_FORMULA_ElectronAffinity_l1 | 0.11953138134990215 |
| CHEMICAL_FORMULA_NsUnfilled_l1 | 0.10335721226261824 |
| CHEMICAL_FORMULA_NUnfilled_l1 | 0.09780109721022519 |
| CHEMICAL_FORMULA_NsValence_l1 | 0.08118081419616913 |
| CHEMICAL_FORMULA_GSestFCClatcnt_l1 | 0.07888443644268245 |
| CHEMICAL_FORMULA_ICSDVolume_l1 | 0.07696848738961315 |
| CHEMICAL_FORMULA_Row_l1 | 0.07500187458125034 |
| CHEMICAL_FORMULA_MiracleRadius_l1 | 0.06839587008787573 |
| CHEMICAL_FORMULA_GSestBCClatcnt_l1 | 0.06776567820725884 |
| CHEMICAL_FORMULA_BoilingT_l1 | 0.06771500457780066 |
| CHEMICAL_FORMULA_GSvolume_pa_l1 | 0.06425279861122248 |
| CHEMICAL_FORMULA_ShearModulus_l1 | 0.06199032157983999 |
| Temperature | 0.03715502350354161 |

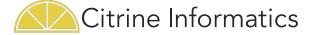Citrine Informatics

# Model Anything!

**NIMS Superconductor Dataset**
*(turns out, superconductors = not easy)*

**NIMS Melting Point Dataset**
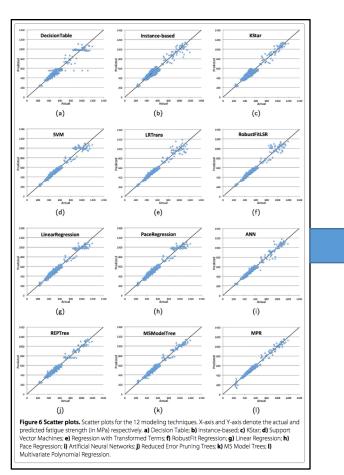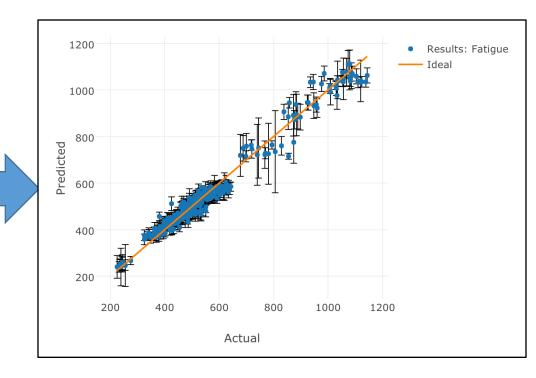*(melting point = much easier)*

# Model Anything!

**Citrine platform creates steel fatigue model from published dataset**
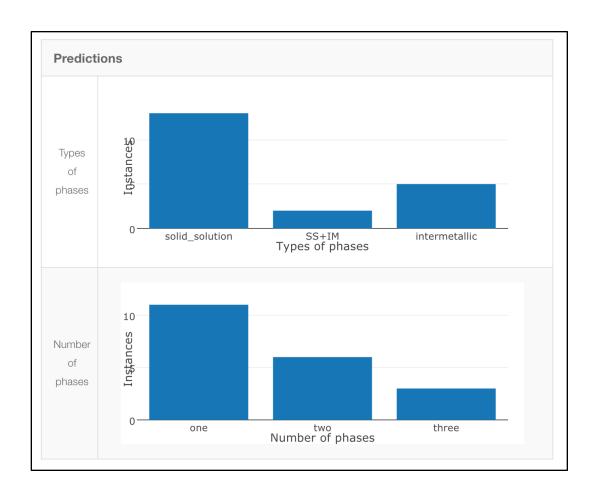Agrawal *et al*., *IMMI* **3** (2014)

# Model Anything!

**Citrine platform trained on
HEA phase stability database**
D Miracle & O Senkov, *Acta Mater* 2016
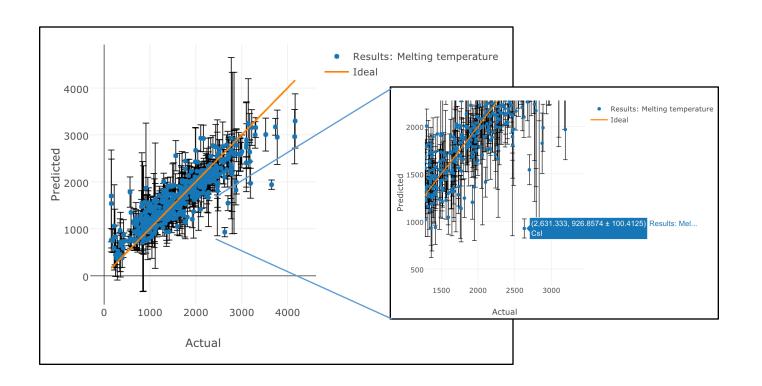
Ex: MoRhRu correctly predicted to be
single-phase SS

# Machine Learning-Assisted Data Curation



NIMS Melting Point Dataset

CsI
Predicted: 927 K
Training: 2631.333 K

1 atm value: 831 K

# Vibrant Ecosystem

Citrine has a new developers' program to enable researchers to publish code that integrates on Citrination

**COMBO Bayesian
Optimization Package**
*K Tsuda, Univ Tokyo / NIMS*

# "Powered by Citrine" Launch

Anchor set of university labs deploying Citrine lab-wide

We are training these users on our API, dataset templates, machine learning templates, PIF data format, and pdf->dataset extraction tools

# Data-Driven Materials Community

Data-Driven Materials Science & Chemistry Newsletter (citrine.io/ddms-newsletter) has >200 weekly readers

*"Your new research highlights are great. There's nothing else out there like this for materials informatics … Particularly when there's a ton of stuff to do in a day, the 1-2 paragraphs plus a figure is a perfect length to start off the day with a hit of research." –a reader*

# Citrine Business Model

Free platform (data & apps) available to everyone

Users of the free platform allow Citrine's algorithms to learn from their data (*Gmail model=monetizing data, not users*)

Industrial users pay for data privacy, while tapping the insights of the free platform

Some premium platform content (e.g., commercial databases)

# Sustainability

Citrine's team of 15+ spends $mm/year to create a scalable, secure, extensible, supported materials data infrastructure for thousands of users—this is not fast, easy, cheap, or temporary
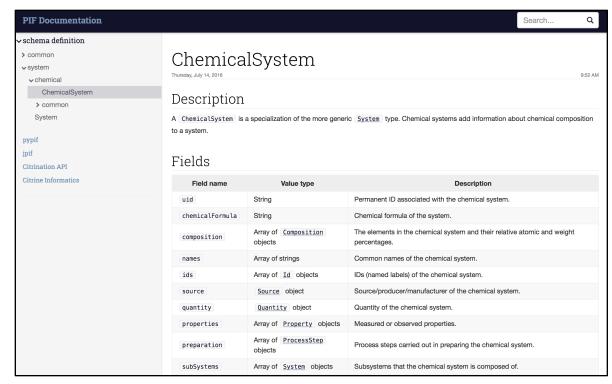
Things we build, track, or have:
- Uptime
- Performance
- Feature velocity
- Security
- Support
- Quality assurance
- Decades of enterprise s/w engineering experience

# Citrine Does Not Lock Users In

Our data structure (pif) is completely open-source JSON: you can export all of your data out of Citrine and back it up elsewhere

We want users using us because they love our platform, not because their data are trapped



**citrine.io/pif**
(also see *MRS Bull* article on pif)

# Let's Create Community Infrastructure

Lots of groups working on roughly the same core web platform features and data plumbing

How can Citrine make it easier for you to build on top of or integrate with our core platform capabilities?

**"Let Citrine handle the IT so you can focus on science"**