

**SISSO
(SURE-
INDEPENDENCE
SCREENING AND
SPARSIFYING
OPERATOR)**





OVERVIEW OF SISSO

- Regressor technique which yields predictive surrogate models expressed as analytical formula.
- Used in conjunction with compressed sensing techniques to filter features.(OMP)
- Overcomes limitation of smaller datasets.



MAIN PARAMETERS

- OMP(Orthogonal Matching Pursuit) – Selects features of data set to be used based off their coefficients. Number of features can be specified.
- Operation set – Specifies possible operations to be used. E.g "(+)(*)(^2)(^3)(^-1)(cos)(sin)".
- Descriptor dimension/rung/complexity – Specifies number of features, composed operations and operations used.
- SISSO then selects model which yields lowest RMSE after specifying parameters.

OMP DEMO

- Applied to Mxenes dataset.
- Input:

```
OMP(descript_dim=3, #number of descriptors to filter out
    data=df_to_ML,
    data_fraction_for_test=0.2,
    split_random_seed=1,
    save_path=SAVE_PATH #set to 0 if don't want to save
)
```

- Results:

	Descriptors	Coefficients
0	charge	373.1032887
1	c	-0.03389305751
2	a	0.6661875896

* All done *

Process finished with exit code 0

- 3 most correlated descriptors are found with their respective coefficients.

SISSO DEMO

- Applied to Mxenes dataset.

- Input:

```
SISSO(OMP_dim=0, #0 if no OMP use, else state dimension to chose via omp
      data=df_to_ML,
      CLEAN_RUN_DIR=False, #False to have SISSO eqn. Need to be false for OMP use too
      optree_depth=2, # rung of the feature space to be constructed
      op_set="(+)(*)(^2)(^3)(^-1)(cos)(sin)", # "(+)(*)(^2)(^3)(^-1)(cos)(sin)"
      descriptor_dim=3, # number of descriptors used, i.e number of columns composed
      maxi_complexity=3, # max feature complexity (number of operators in a feature)
      data_fraction_for_test=0.2,
      split_random_seed=4,
      save_path=SAVE_PATH, #set to 0 if don't want to save
      )
```

- Results:

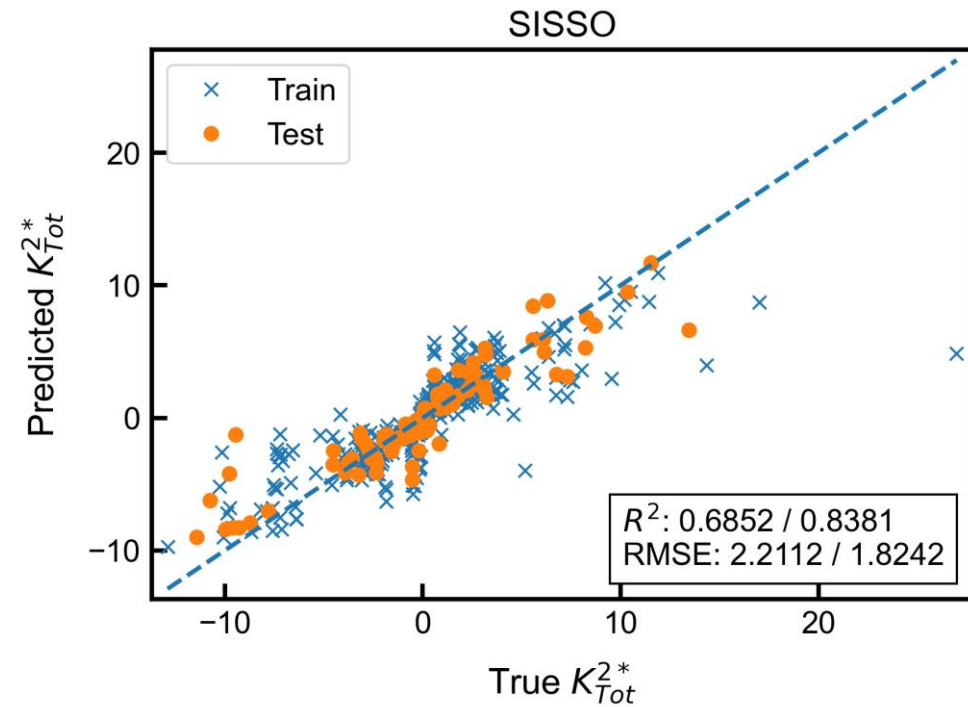
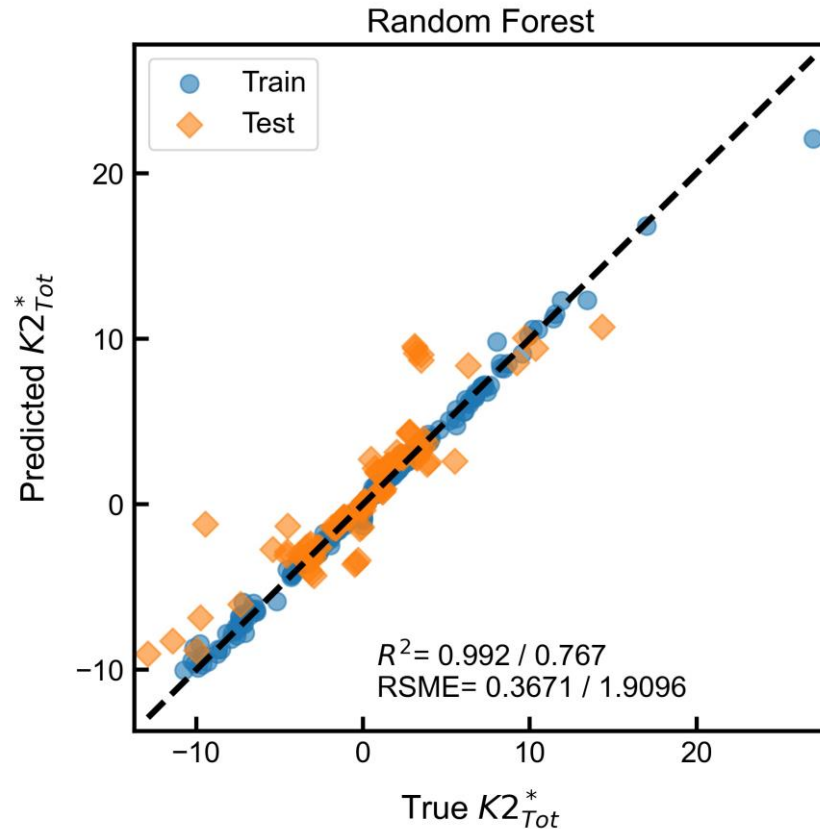
```
SISSO w/o OMP:
+24689514(sin(charge)*(atomic_states*atomic_tot_elec)) -24688847(atomic_states*(charge*atomic_tot_elec)) +3839(sin(a)*(charge*atomic_states))
Dimension: 3
RMSE: 2.2112/1.8242
R^2: 0.6852/0.8381
*****
*           All done           *
*****
```

- Analytic equation is obtained with dimension, RMSE and R-squared values.

COMPARISON WITH GENERAL ML

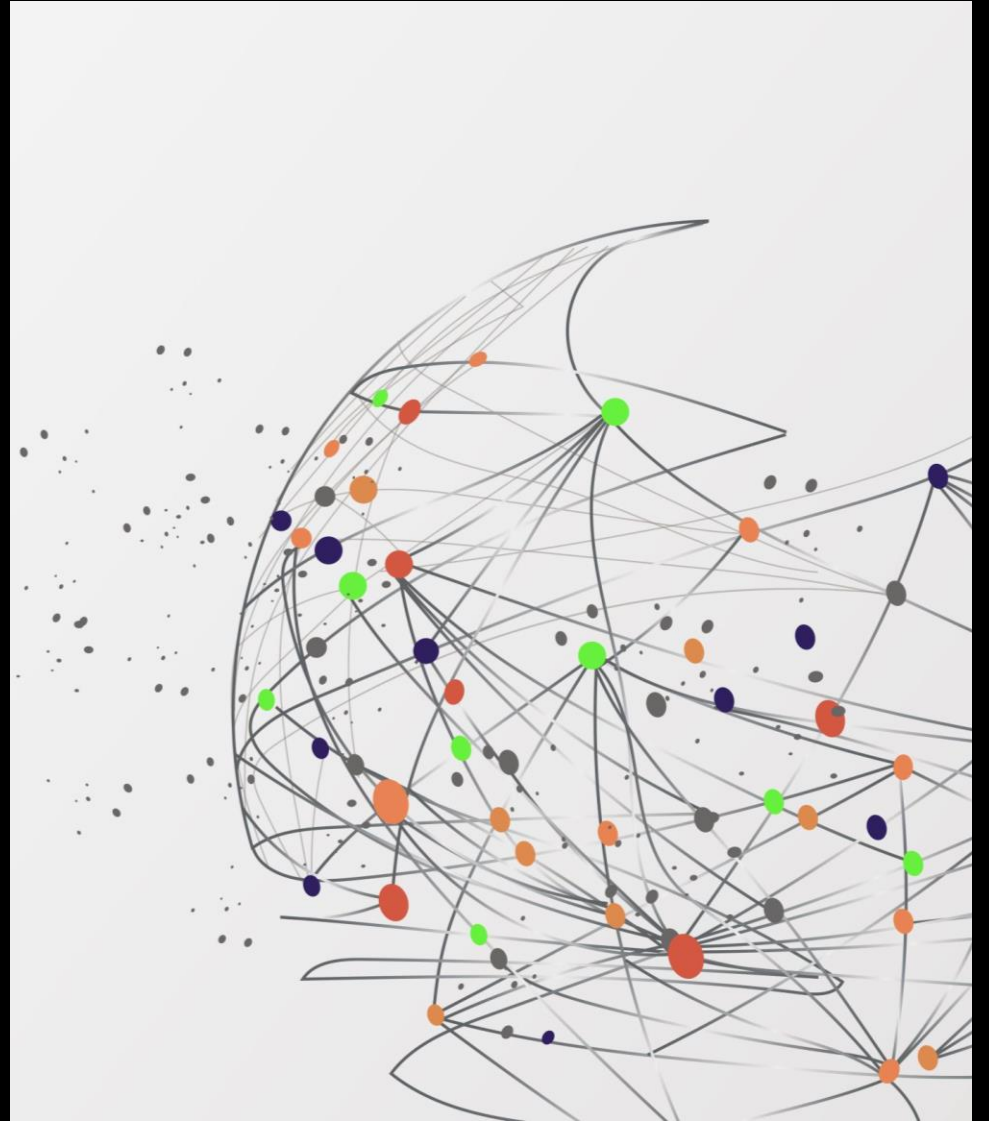
- Applied to Mxenes dataset.

- Plots:



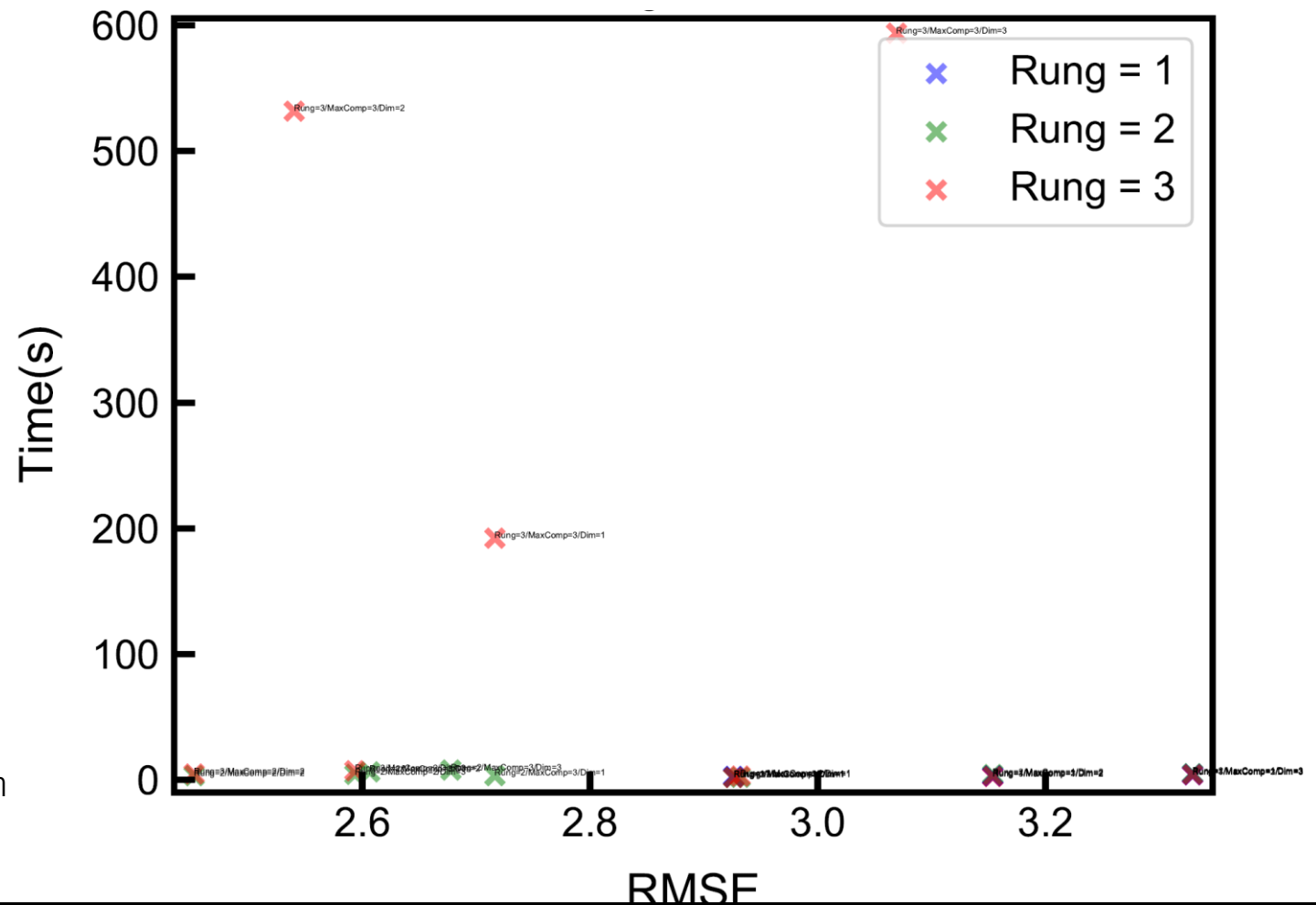
- Performance is comparable, while attaining an analytic equation.

SISSO PERFORMANCE AND RESULTS ANALYSIS



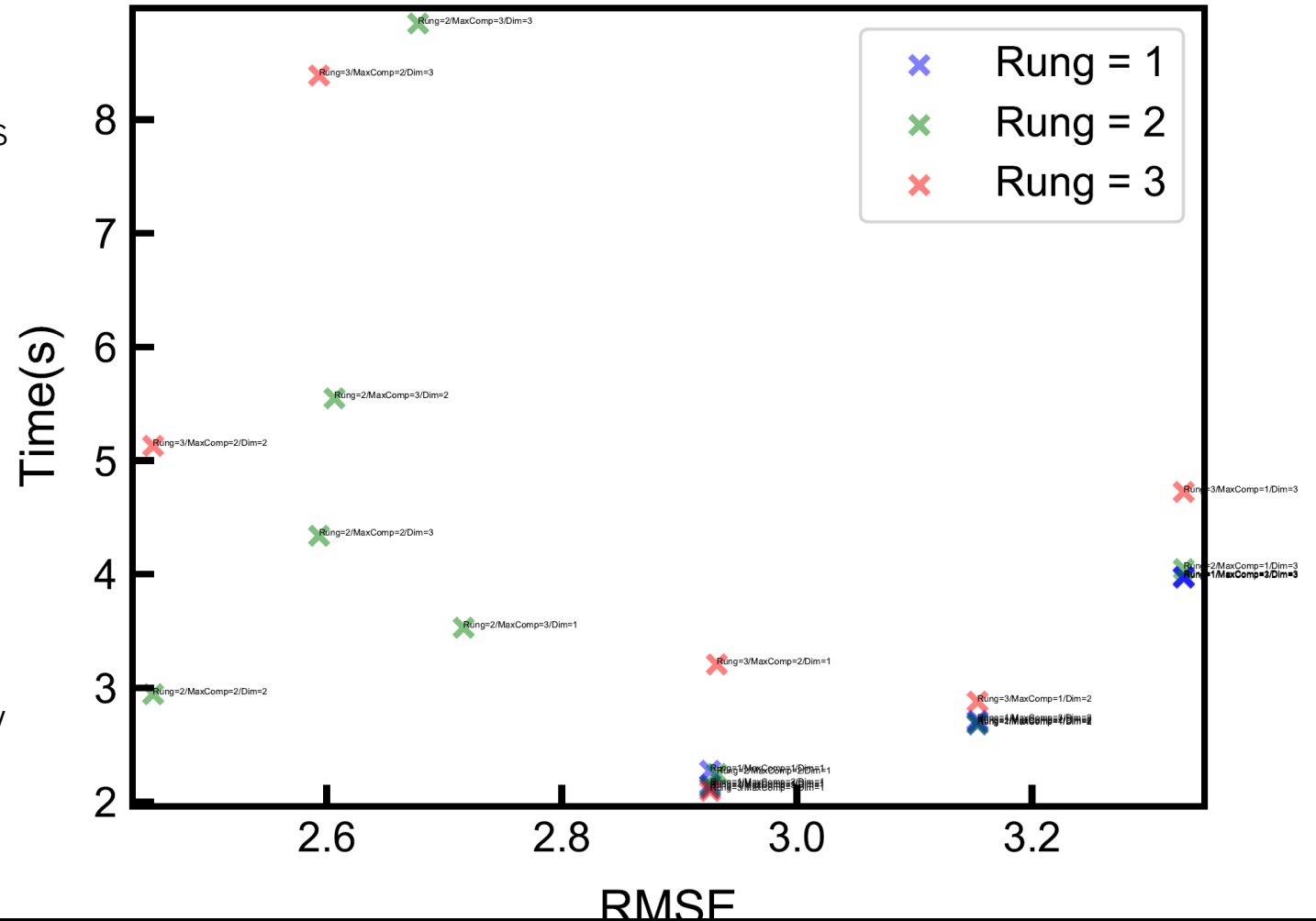
PERFORMANCE PLOT

- Applied to Mxenes dataset with target variable K_{Tot}^{2*}
- Points with varied parameters are plotted. (Rung, Maximum complexity, Dimension)
- Observed that rung = 3 (red) is not necessarily more optimal for accuracy or time as points are distributed evenly. (potentially overfitting)
- Observed that rung = 1 (blue) is less accurate as is distributed towards bottom right.



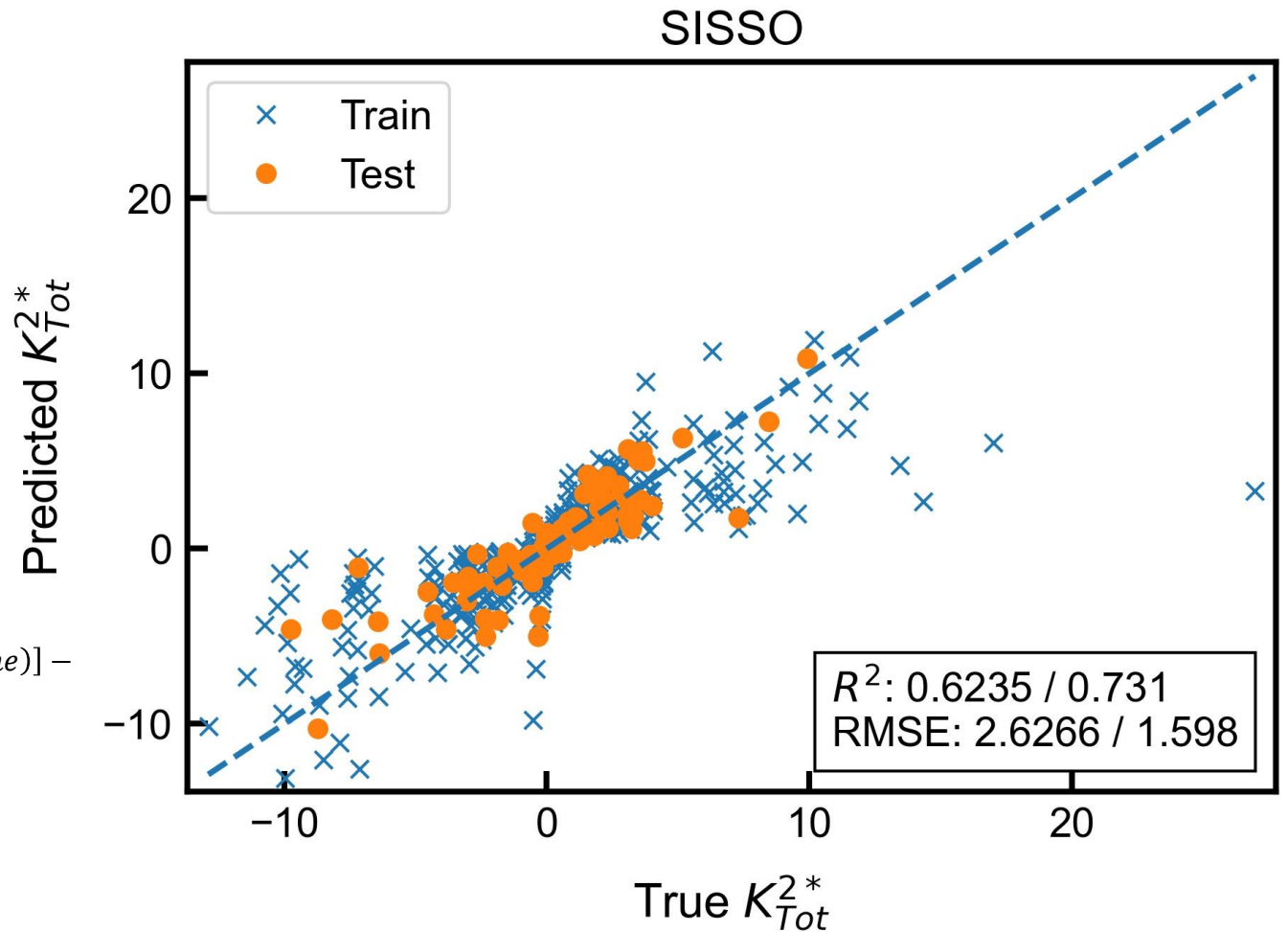
PERFORMANCE PLOT (ZOOMED)

- Zoomed into bottom portion of previous plot
- Observed that rung = 2 (green) is optimal for both time and precision on mxenes dataset.
- Best performance is observed with parameters rung, maximum complexity and dimension = 2. (Bottom left)



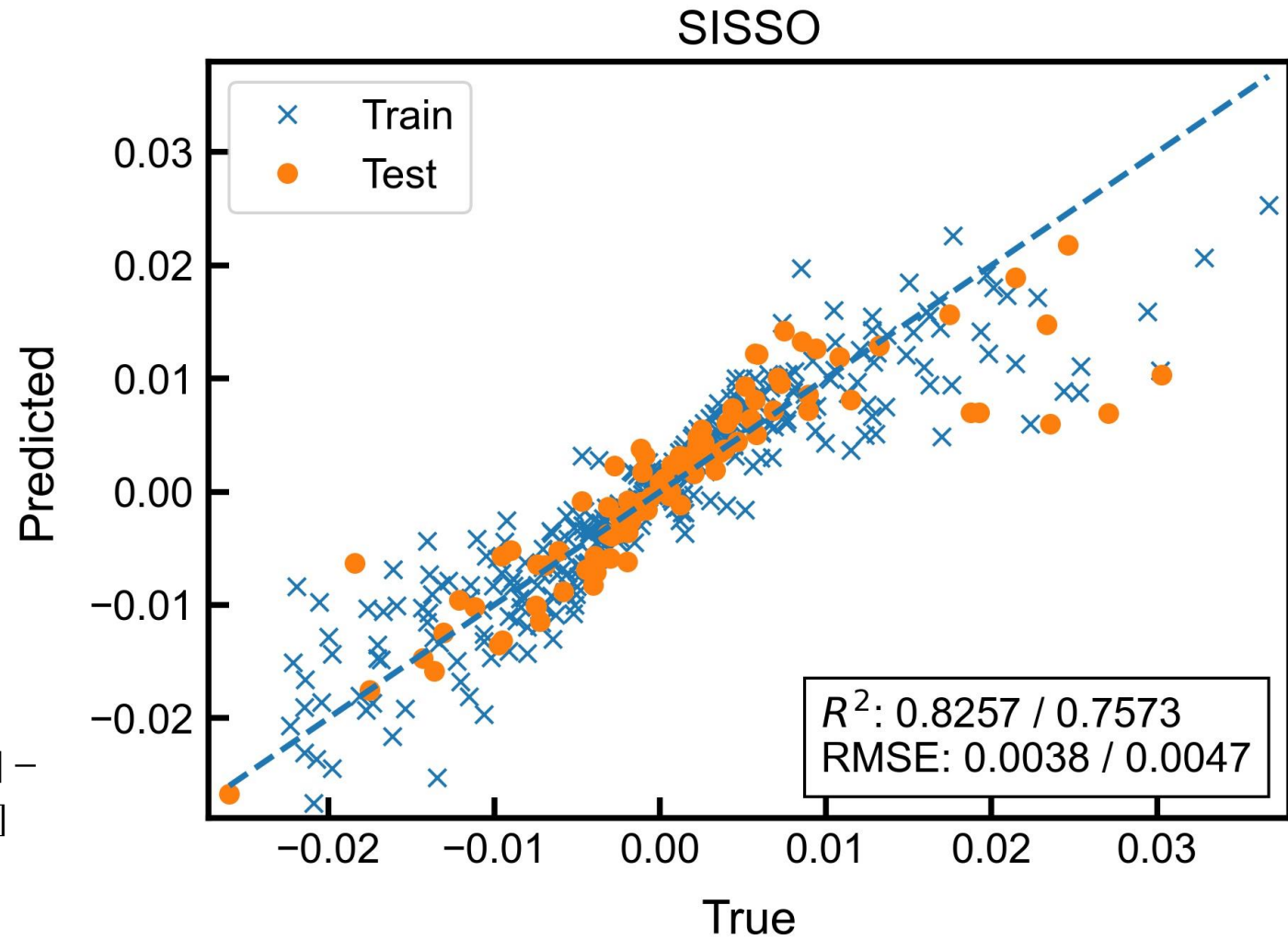
APPLICATION TO MXENES DATA (K_{Tot}^{2*})

- Applied best parameters from analysis (previous slide) to K_{Tot}^{2*} .
- Accuracy is high
- Analytic Equation:
$$6308.041021[atomicstates * (charge * totelecpervolume)] - 8.738623357[(c33 * charge) * atomicstates]$$



APPLICATION TO MXENES DATA(EFFICIENCY)

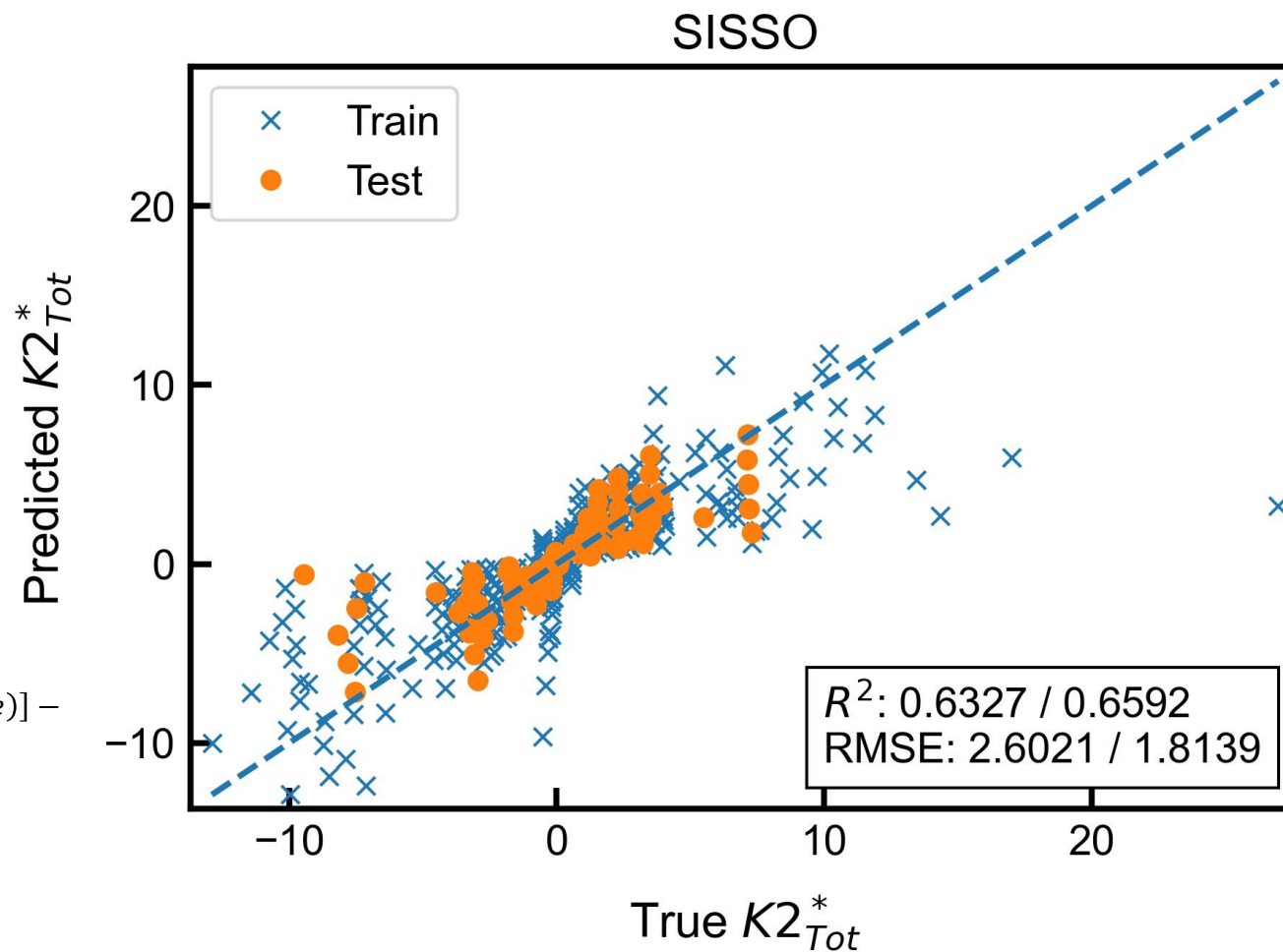
- Applied same parameters to predict effitotsigned. (rung,maxcomp,dim = 2)
- Accuracy is high
- Analytic Equation:
$$0.008881534082[charge * abs(c11 - c33)] - 4.840361695e - 06[c11 * (mass * charge)]$$



APPLICATION TO STRATIFIED MXENES

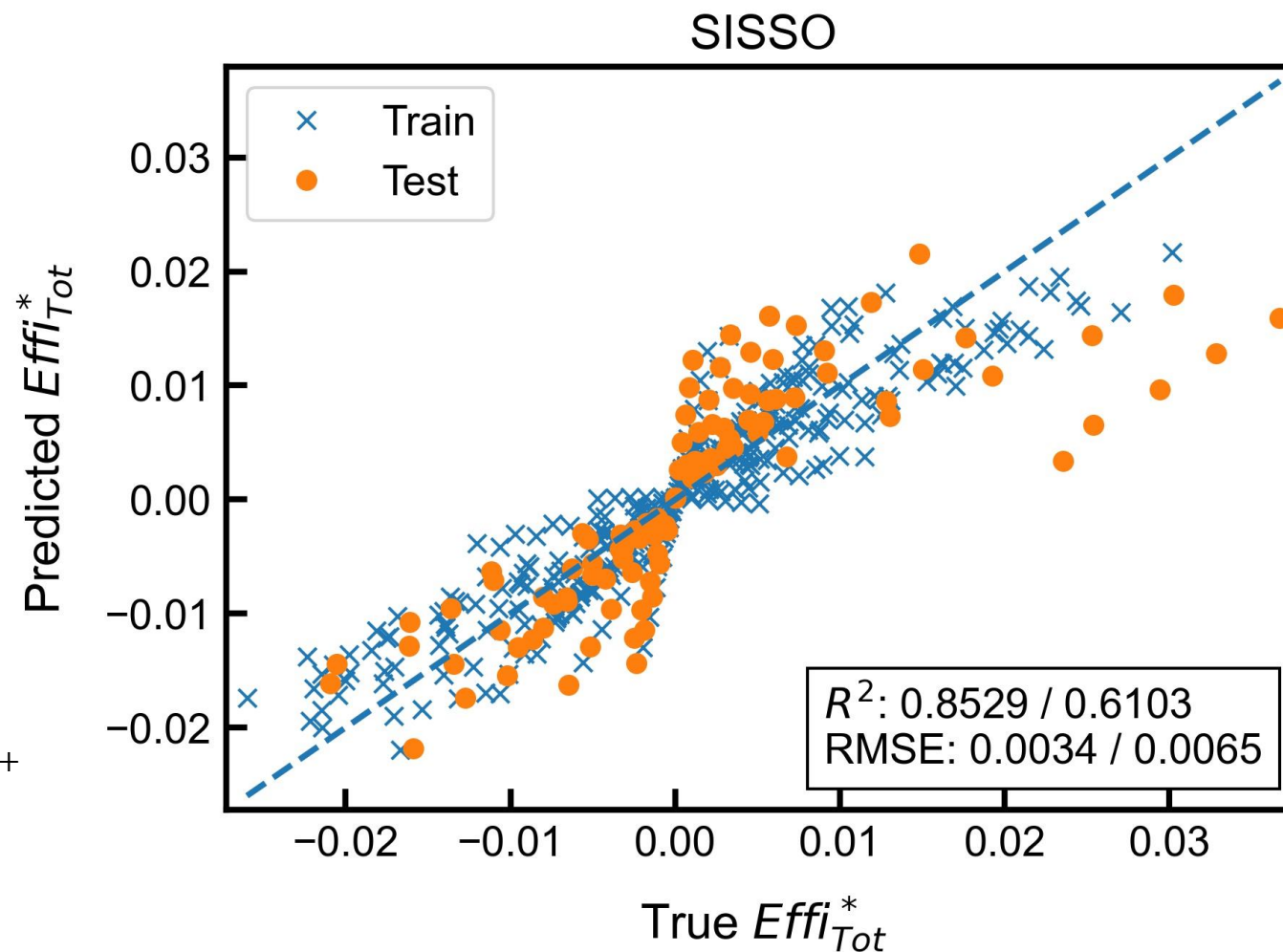
DATA (K_{Tot}^{2*})

- Applied best parameters from analysis to K_{Tot}^{2*} with stratification by material group.
- Accuracy is high
- Analytic Equation:
$$6215.127894[atomicstates * (charge * totelecpervolume)] - 8.597285392[(c33 * charge) * atomicstates]$$

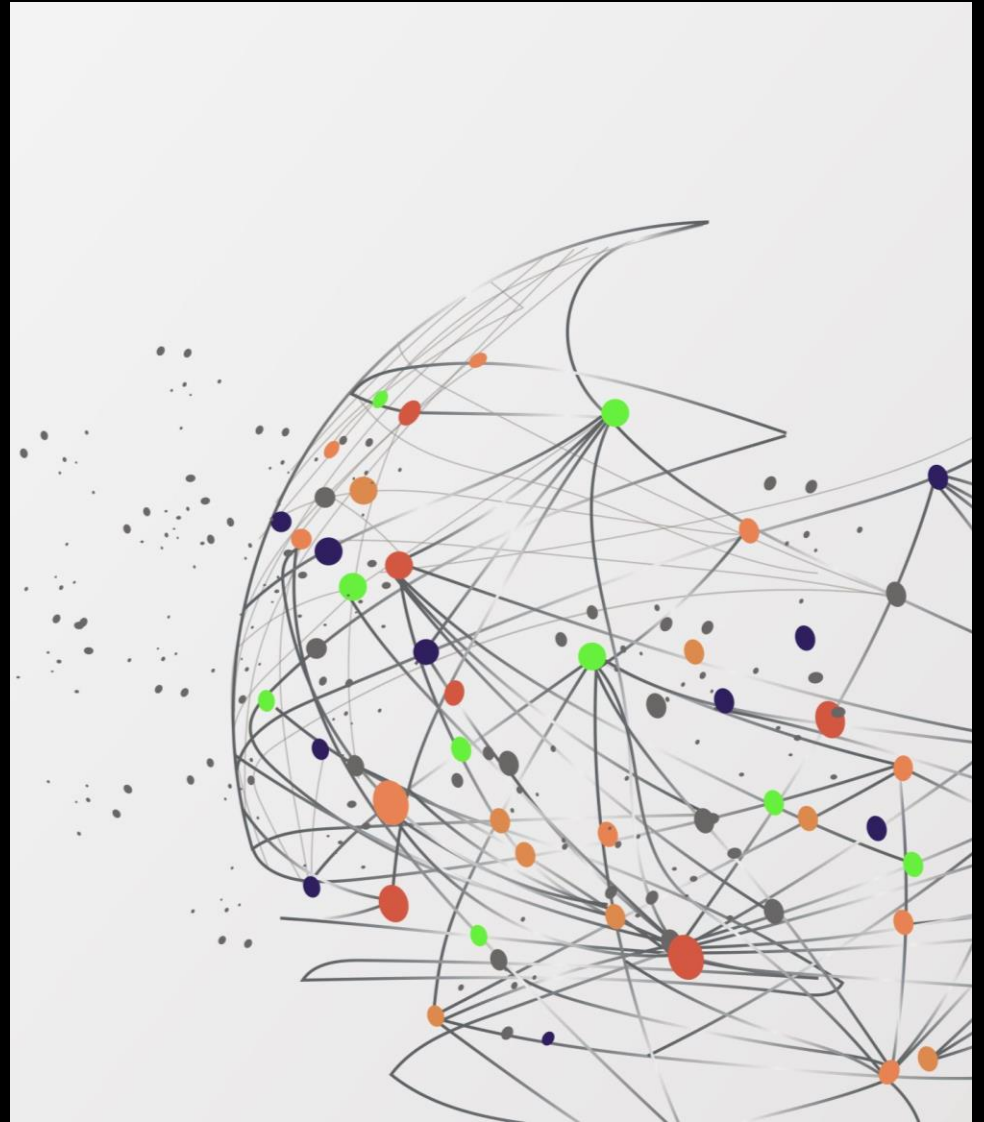


APPLICATION TO STRATIFIED MXENES DATA(EFFICIENCY)

- Applied same parameters to predict effitotsigned with stratification by material group. (rung,maxcomp,dim = 2)
- Accuracy is high
- Analytic Equation:
$$0.3931274926[(charge * atomic_tot_elec) - charge] + 0.1873179569[charge * (atomic_tot_elec - density)]$$



SISSO-HEAS RESULTS ANALYSIS





METHODOLOGY

- Selected target variable. (Elongation, YS, HV, UTS)
- Used OMP on HEAS data to select 5 most correlated features.
- Applied SISSO to the set of selected features, with hyperparameter tuning (range 1 to 3) and 5-fold cross-validation for rung, dimension and maximum complexity. (most frequent optimal combination of parameters: 1,2,1)
- Repeated for HEAS data with only composition features and with only physical features to compare results.

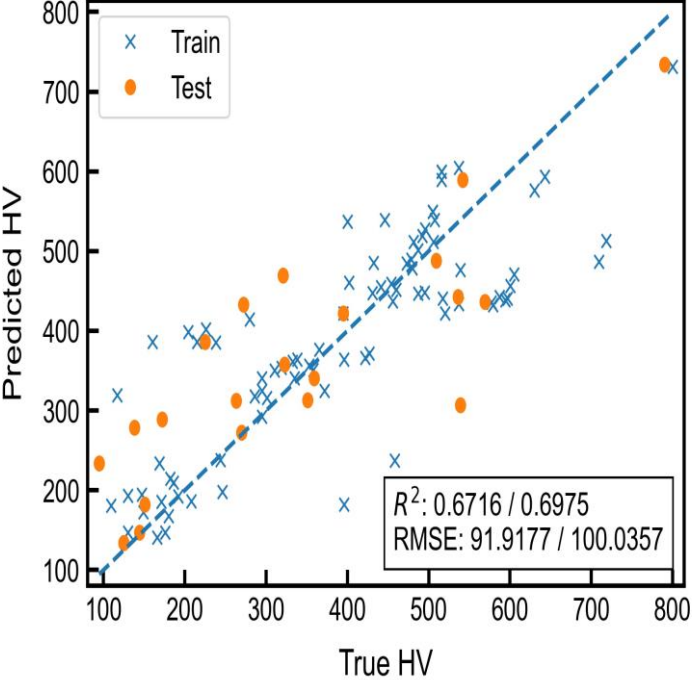
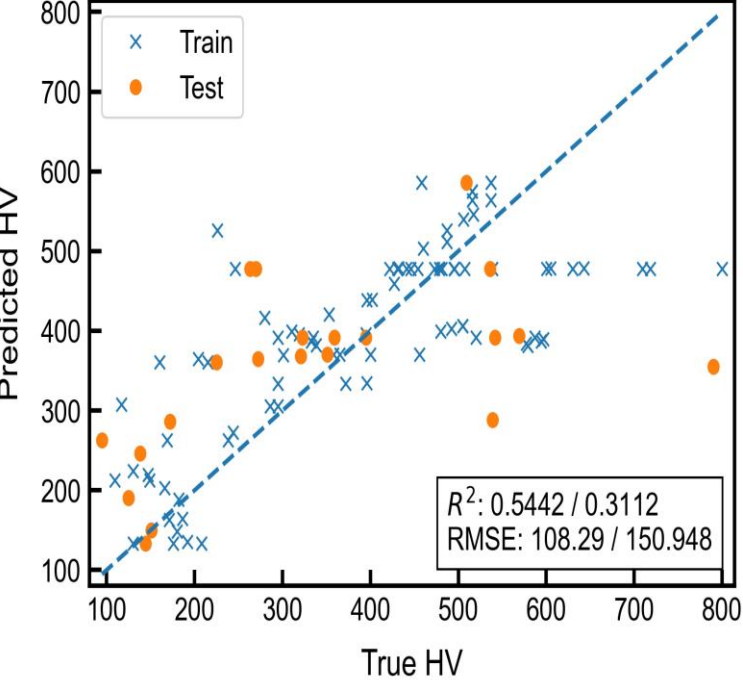
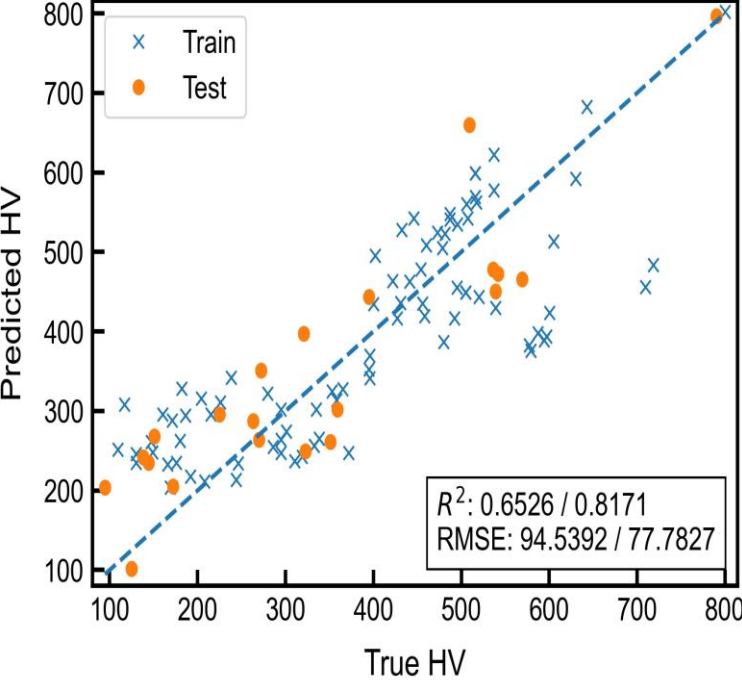
TARGET VARIABLE:
ELONGATION(%)

Dataset	Whole data	Composition data	Physical data
Plot	<p>SISSO</p> <p>Predicted Elongation_(%)</p> <p>True Elongation_(%)</p> <p>R^2: 0.559 / 0.4511 RMSE: 14.2597 / 15.0964</p>	<p>SISSO</p> <p>Predicted Elongation_(%)</p> <p>True Elongation_(%)</p> <p>R^2: 0.4464 / -0.3721 RMSE: 15.9766 / 23.3333</p>	<p>SISSO</p> <p>Predicted Elongation_(%)</p> <p>True Elongation_(%)</p> <p>R^2: 0.5118 / 0.3347 RMSE: 15.0033 / 16.6203</p>
Equation	<p>$-52.97817093(\textit{Shear_modulus_delta} - \textit{Ni})$ $-14.76128983(\textit{Mixing_enthalpy}$ $\ast \textit{Electronegativity_local_mismatch})$</p>	<p>$-55.07104127(\textit{Al} + \textit{Mo}) + 1516.61397(\textit{Co})^3$</p>	<p>$-1.777507734(\textit{Mixing_enthalpy}$ $- \textit{Interant_electrons})$ $-4.300657093(\textit{Interant_electrons}$ $\ast \textit{Shear_modulus_delta})$</p>

Accuracy (worst to best): Composition -> Physical -> Whole

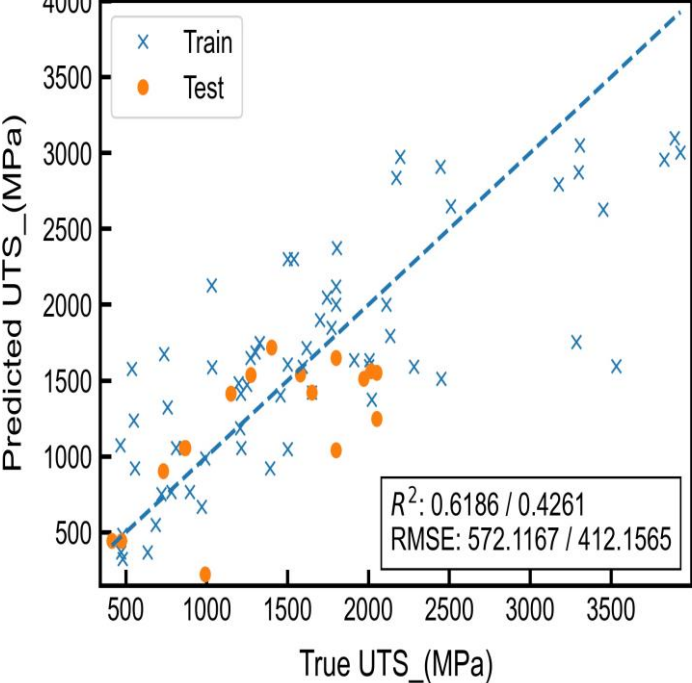
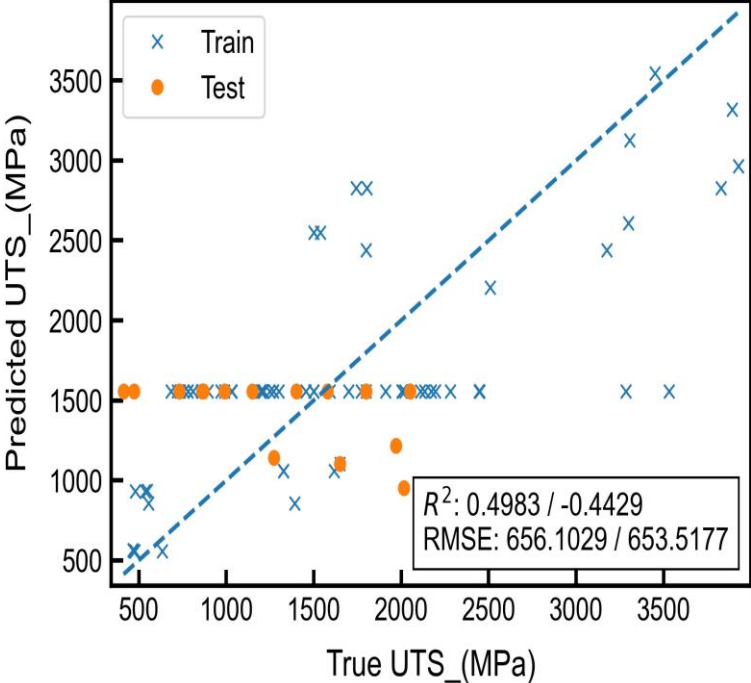
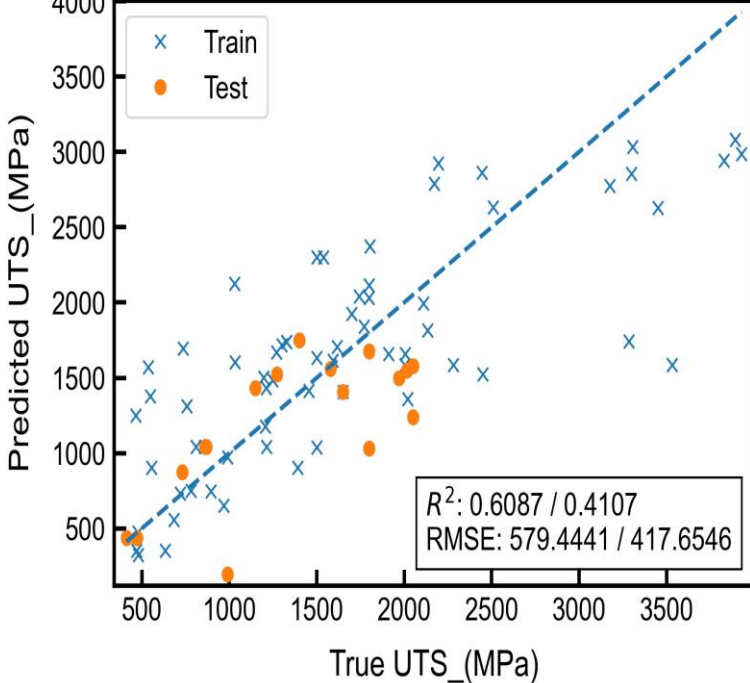
TARGET VARIABLE:

HV

Dataset	Whole data	Composition data	Physical data
Plot	<p>SISSO</p>  <p>Predicted HV</p> <p>True HV</p> <p>R^2: 0.6716 / 0.6975 RMSE: 91.9177 / 100.0357</p>	<p>SISSO</p>  <p>Predicted HV</p> <p>True HV</p> <p>R^2: 0.5442 / 0.3112 RMSE: 108.29 / 150.948</p>	<p>SISSO</p>  <p>Predicted HV</p> <p>True HV</p> <p>R^2: 0.6526 / 0.8171 RMSE: 94.5392 / 77.7827</p>
Equation	$4194.979928 \exp(\text{Electronegativity_delta}) - 766.5006764(\text{Mn} + \text{Shear_modulus_strength_model})$	$-861.5603106(\text{Co} + \text{Mn}) + 431.6943479(\text{Al} - \text{Zr})$	$1754.503458(\text{Electronegativity_delta} - \text{Shear_modulus_strength_model}) + 11.33349587(\text{Mixing_enthalpy} - \text{VEC_mean})$

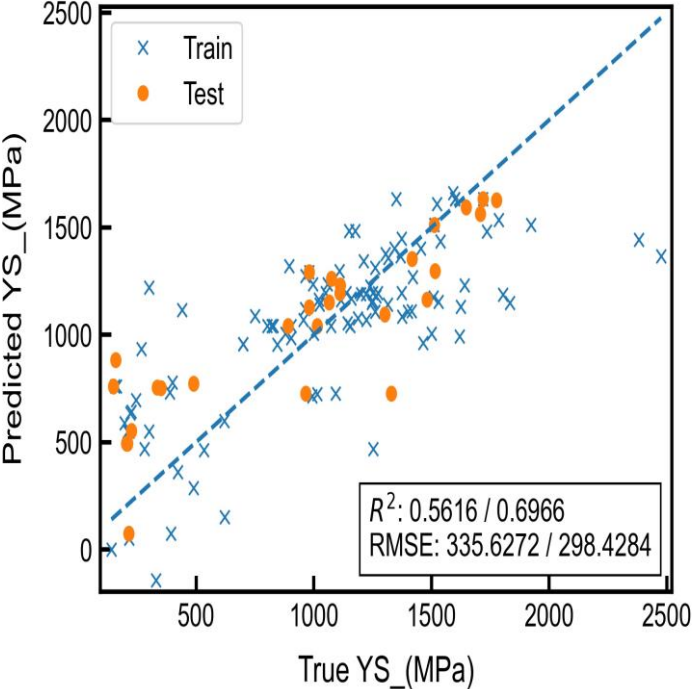
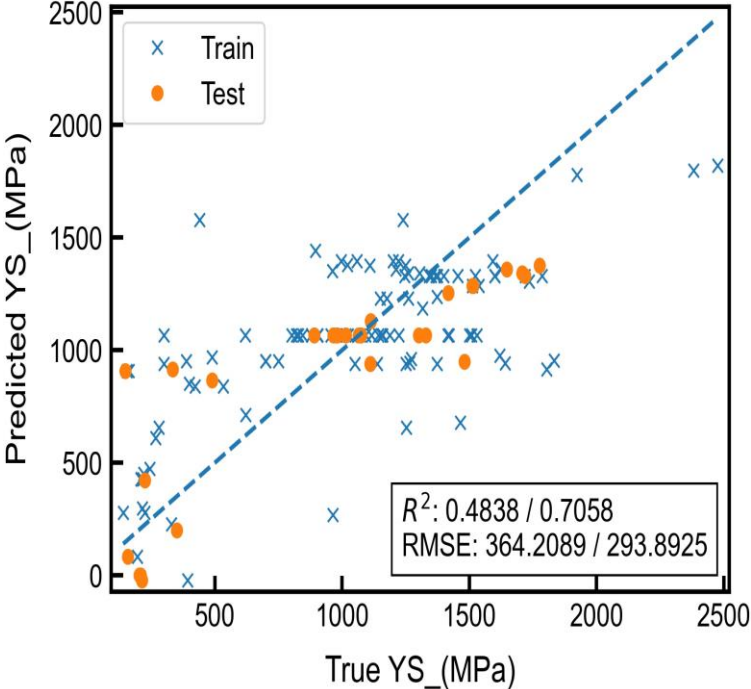
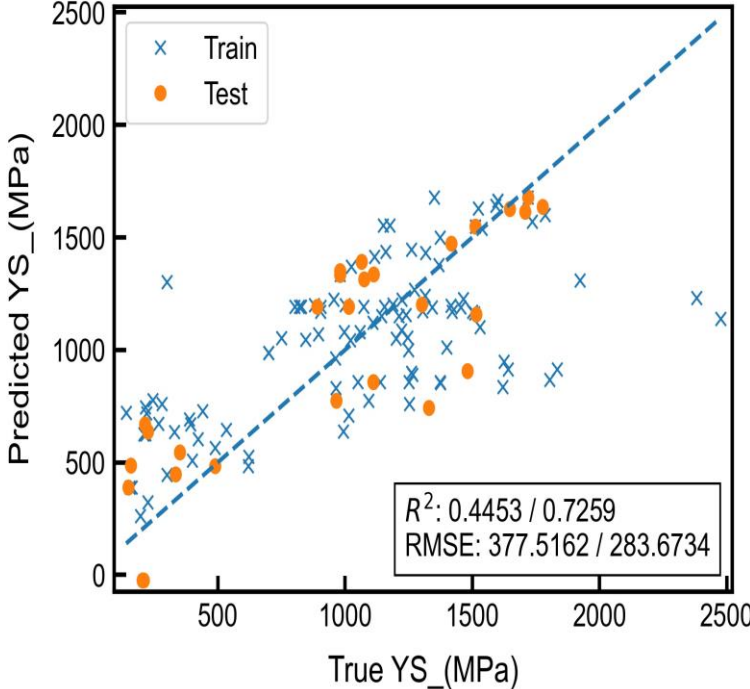
Accuracy (worst to best): Composition -> Whole -> Physical

TARGET VARIABLE:
UTS(MPA)

Dataset	Whole data	Composition data	Physical data
Plot	<div><p>SISSO</p><p>Predicted UTS_(MPa)</p><p>True UTS_(MPa)</p><p>R^2: 0.6186 / 0.4261 RMSE: 572.1167 / 412.1565</p></div>	<div><p>SISSO</p><p>Predicted UTS_(MPa)</p><p>True UTS_(MPa)</p><p>R^2: 0.4983 / -0.4429 RMSE: 656.1029 / 653.5177</p></div>	<div><p>SISSO</p><p>Predicted UTS_(MPa)</p><p>True UTS_(MPa)</p><p>R^2: 0.6087 / 0.4107 RMSE: 579.4441 / 417.6546</p></div>
Equation	$36551.71097(Radii_gamma * Shear_modulus_delta) - 35154.83364(Sn + Shear_modulus_delta)$	$31794.64019(Mo * Zr) - 2499.812945(Cu + Mn)$	$35335.70666(Radii_gamma * Shear_modulus_delta) - 33784.28339(Shear_modulus_delta)$

Accuracy (worst to best): Composition -> Physical = Whole

TARGET VARIABLE:
YS(MPA)

Dataset	Whole data	Composition data	Physical data
Plot	<div><p>SISSO</p><p>Predicted YS_(MPa)</p><p>True YS_(MPa)</p><p>R^2: 0.5616 / 0.6966 RMSE: 335.6272 / 298.4284</p></div>	<div><p>SISSO</p><p>Predicted YS_(MPa)</p><p>True YS_(MPa)</p><p>R^2: 0.4838 / 0.7058 RMSE: 364.2089 / 293.8925</p></div>	<div><p>SISSO</p><p>Predicted YS_(MPa)</p><p>True YS_(MPa)</p><p>R^2: 0.4453 / 0.7259 RMSE: 377.5162 / 283.6734</p></div>
Equation	$13265.55335(\text{Electronegativity_delta} - \text{Mg}) - 2067.858599(\text{Cu} + \text{Mn})$	$-2561.36509(\text{Mn} - \text{Cr}) + 1314.036639 \text{abs}(\text{Co} - \text{Mo}) - 4508.87323 \text{abs}(\text{Co} - \text{Mg})$	$9499.973507(\text{Yang_delta} + \text{Electronegativity_delta}) - 0.6182925161(\text{VEC_mean})^3$

Accuracy (worst to best): Whole -> Composition -> Physical



NOTABLE RESULTS

- For all targets, using only composition data was less accurate than only physical data.
- For targets HV, YS, whole data was less accurate than only physical. This implies that adding composition data as predictors reduced the overall accuracy.
- Only for target Elongation, whole data was more accurate than physical. This implies that composition data helped improve accuracy.
- For target YS, composition features made up a larger proportion of the features used in whole data compared to other targets. (3 comp. and 1 phy. feat. vs 1 comp. and 3 phy. feat.) However, using whole data performed worse than composition or physical data only in this case.

MULTI-TASK SISO





OVERVIEW

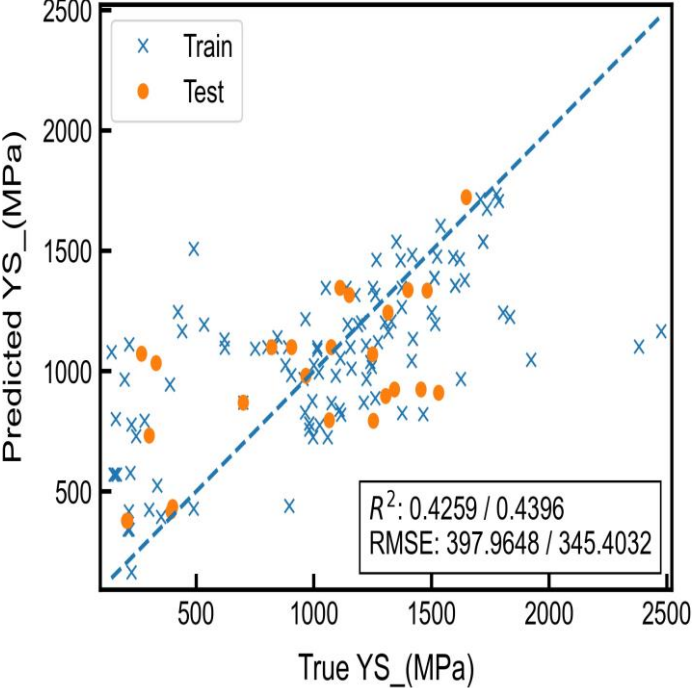
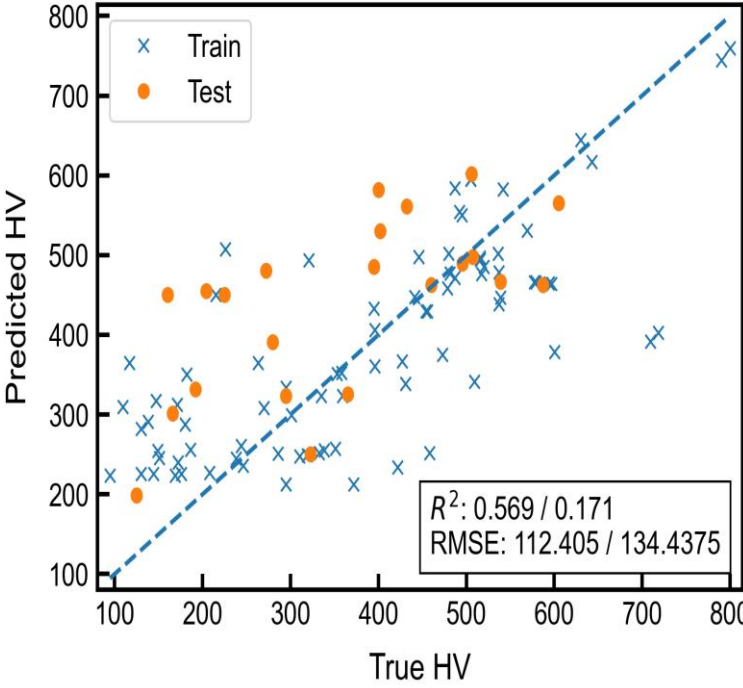
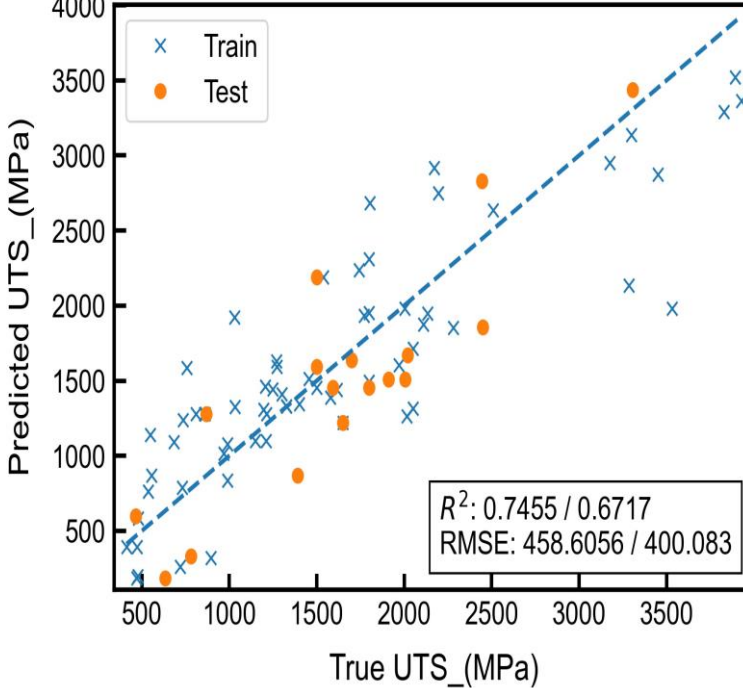
- Takes as inputs, “X” data columns and multiple “Y” target columns.
- For each feature (variable in the equation), single-task SISSO is run for each individual “Y” target column.
- RMSE for that feature is calculated as the average RMSE of each single task run, and the features which minimise RMSE are selected for usage in the equation.
- Generates analytic equation for each target, with a fixed set of features. i.e. All targets share the same features but with different coefficients in their equations.
- Advantages: Predictions are generated for NA values.



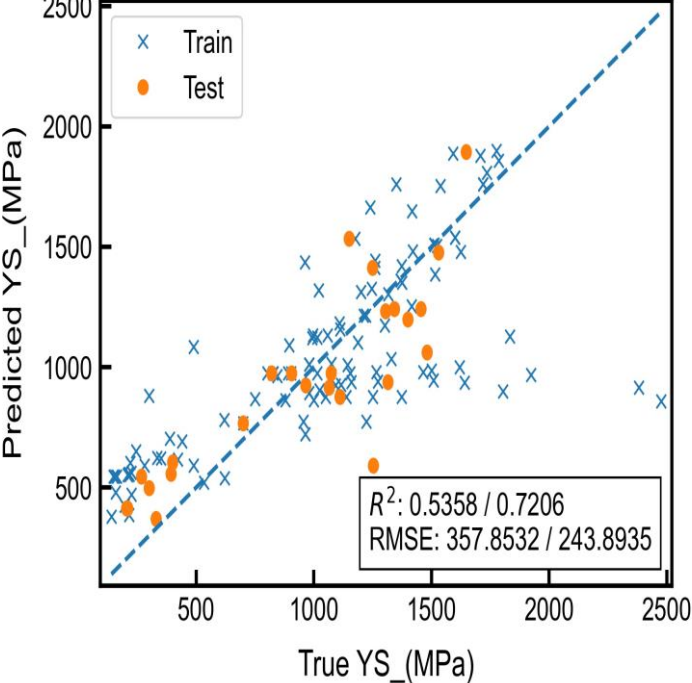
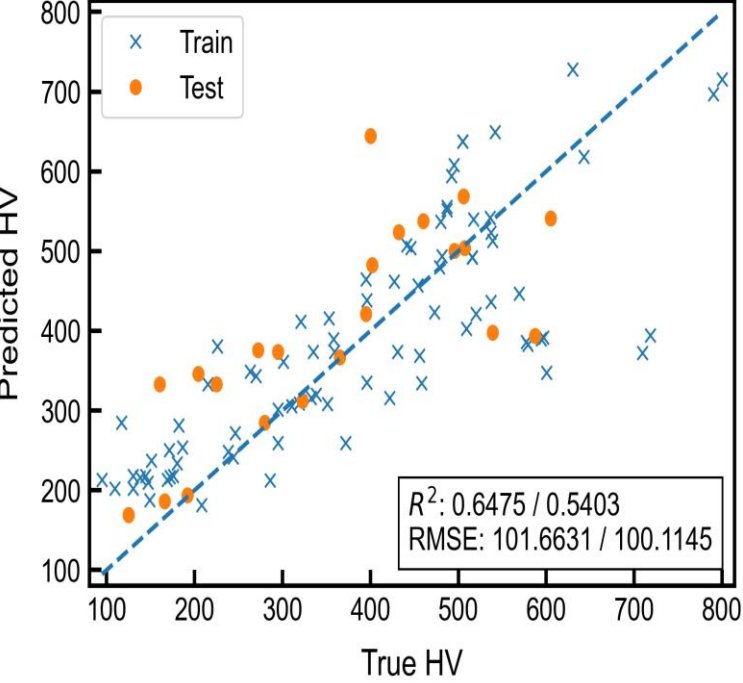
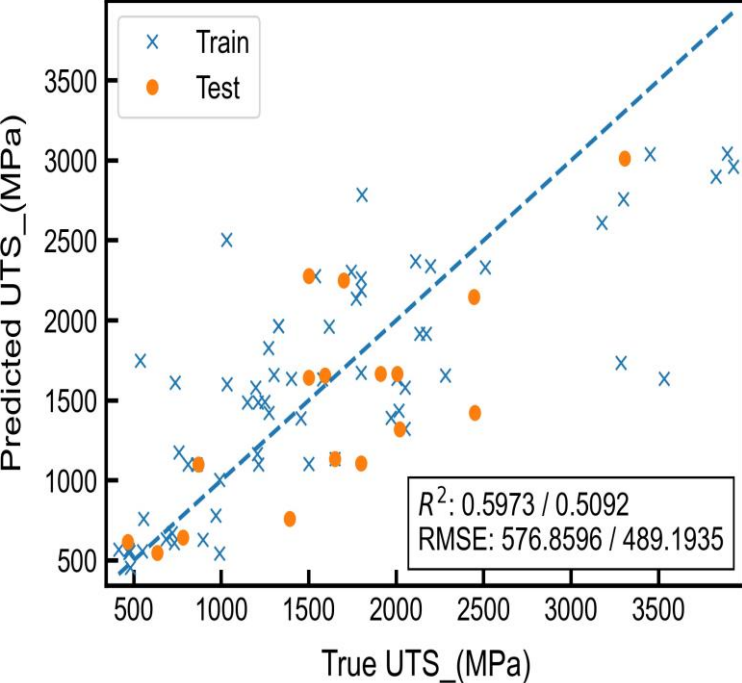
APPLICATION TO HEAS DATASET

- “Y” target variables are input as 'YS (MPa)', 'HV' and 'UTS (MPa)'.
- Performed with only physical data.
- Run with configuration $\text{rung} = 2$, $\text{dimension} = 2$ and complexity 2.
- Run with configuration $\text{rung} = 1$, $\text{dimension} = 2$ and complexity 1.
- Compared with previous results with optimal config: 1,2,1 for config: 1,2,1.

CONFIG:
RUNG=2,DIM=2,COMP=2

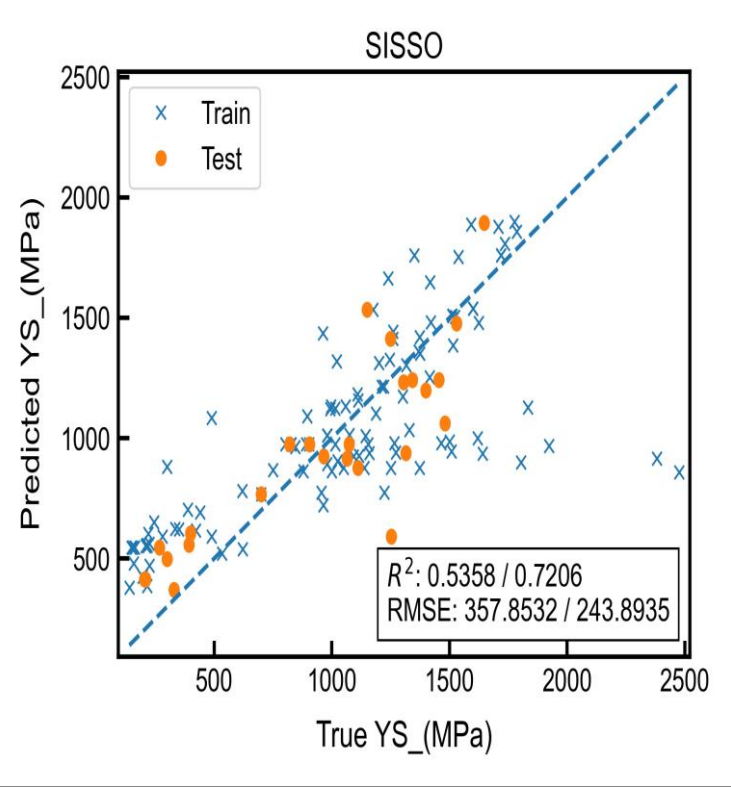
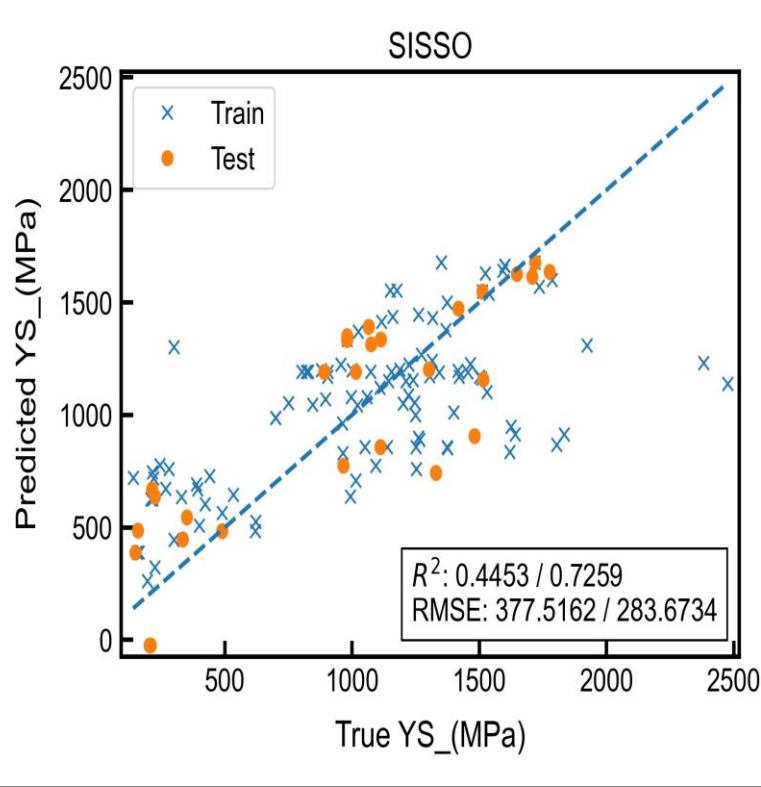
Target	YS (MPa)	HV	UTS (MPa)
Plot	<div><p>SISSO</p><p>Predicted YS_(MPa)</p><p>True YS_(MPa)</p><p>$R^2: 0.4259 / 0.4396$ $RMSE: 397.9648 / 345.4032$</p></div>	<div><p>SISSO</p><p>Predicted HV</p><p>True HV</p><p>$R^2: 0.569 / 0.171$ $RMSE: 112.405 / 134.4375$</p></div>	<div><p>SISSO</p><p>Predicted UTS_(MPa)</p><p>True UTS_(MPa)</p><p>$R^2: 0.7455 / 0.6717$ $RMSE: 458.6056 / 400.083$</p></div>
Equation	<div>$-48827.35389(\log(Radii_gamma) * Shear_modulus_strength_model)$$-41.56337773abs(Interant_d_electrons - (Mixing_enthalpy + Shear_modulus_local_mismatch))$</div>	<div>$-51409.68208(\log(Radii_gamma) * Shear_modulus_strength_model)$$-0.6022698101abs(Interant_d_electrons - (Mixing_enthalpy + Shear_modulus_local_mismatch))$</div>	<div>$-188661.3072(\log(Radii_gamma) * Shear_modulus_strength_model)$$-51.10224241abs(Interant_d_electrons - (Mixing_enthalpy + Shear_modulus_local_mismatch))$</div>

CONFIG:
RUNG=1,DIM=2,COMP=1

Target	YS (MPa)	HV	UTS (MPa)
Plot	<p>SISSO</p>  <p>Predicted YS_(MPa)</p> <p>True YS_(MPa)</p> <p>R^2: 0.5358 / 0.7206 RMSE: 357.8532 / 243.8935</p>	<p>SISSO</p>  <p>Predicted HV</p> <p>True HV</p> <p>R^2: 0.6475 / 0.5403 RMSE: 101.6631 / 100.1145</p>	<p>SISSO</p>  <p>Predicted UTS_(MPa)</p> <p>True UTS_(MPa)</p> <p>R^2: 0.5973 / 0.5092 RMSE: 576.8596 / 489.1935</p>
Equation	<p>$-151253.3373(\text{Yang_delta}$ * $\text{Shear_modulus_strength_model}$) $-127.484471\text{abs}(\text{VEC_mean}$ $- \text{Mean_cohesive_energy})$</p>	<p>$-77254.61692(\text{Yang_delta}$ * $\text{Shear_modulus_strength_model}$) $-30.73113347\text{abs}(\text{VEC_mean}$ $- \text{Mean_cohesive_energy})$</p>	<p>$-348775.8453(\text{Yang_delta}$ * $\text{Shear_modulus_strength_model}$) $-25.31638202\text{abs}(\text{VEC_mean}$ $- \text{Mean_cohesive_energy})$</p>

COMPARISON:

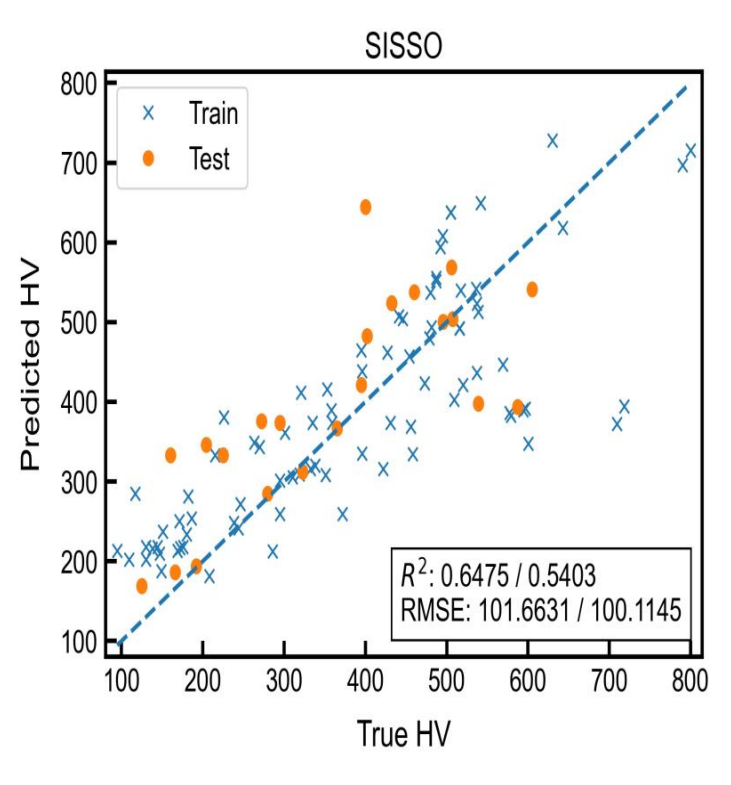
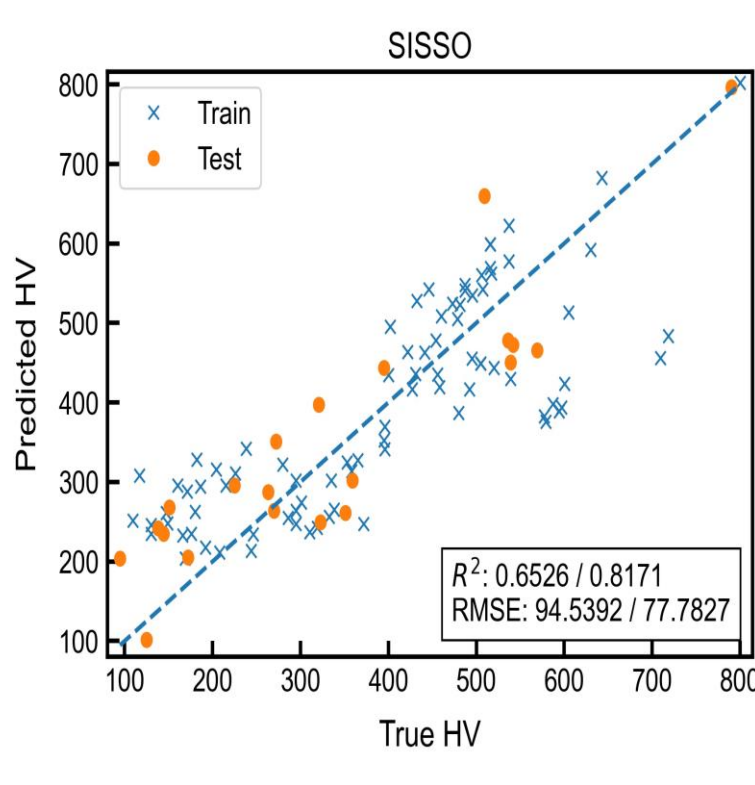
YS(MPa)

SISSO Type	Multi-task	Single-task
Plot		
Equation	$-151253.3373(\text{Yang_delta} * \text{Shear_modulus_strength_model}) - 127.484471 \text{abs}(\text{VEC_mean} - \text{Mean_cohesive_energy})$	$9499.973507(\text{Yang_delta} + \text{Electronegativity_delta}) - 0.6182925161(\text{VEC_mean})^3$

- Configuration: 1,2,1 is used for both.
- Multi-task performs better in comparison to single-task.
- YS data has 56 nan values out of 193. i.e. 29% nan values.

COMPARISON:

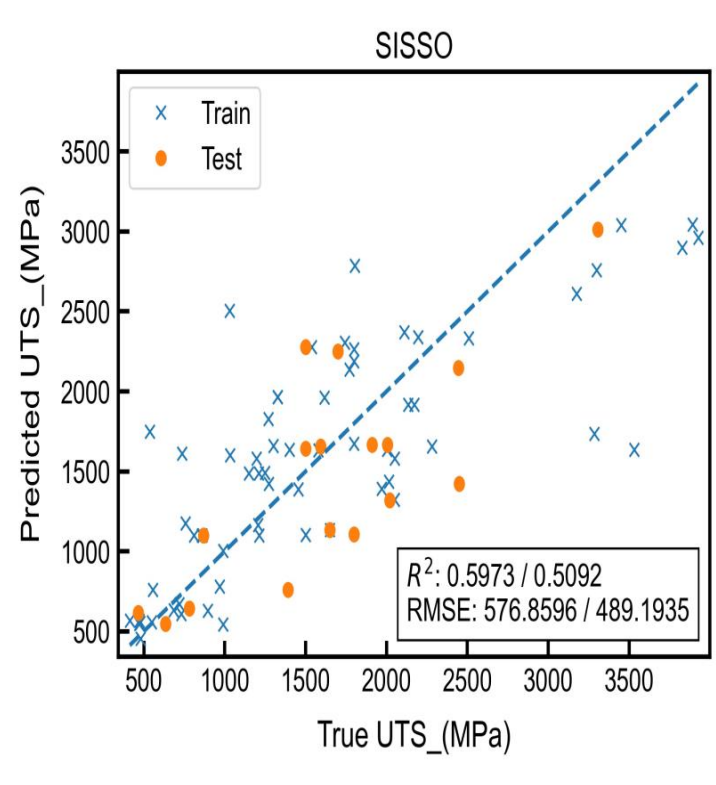
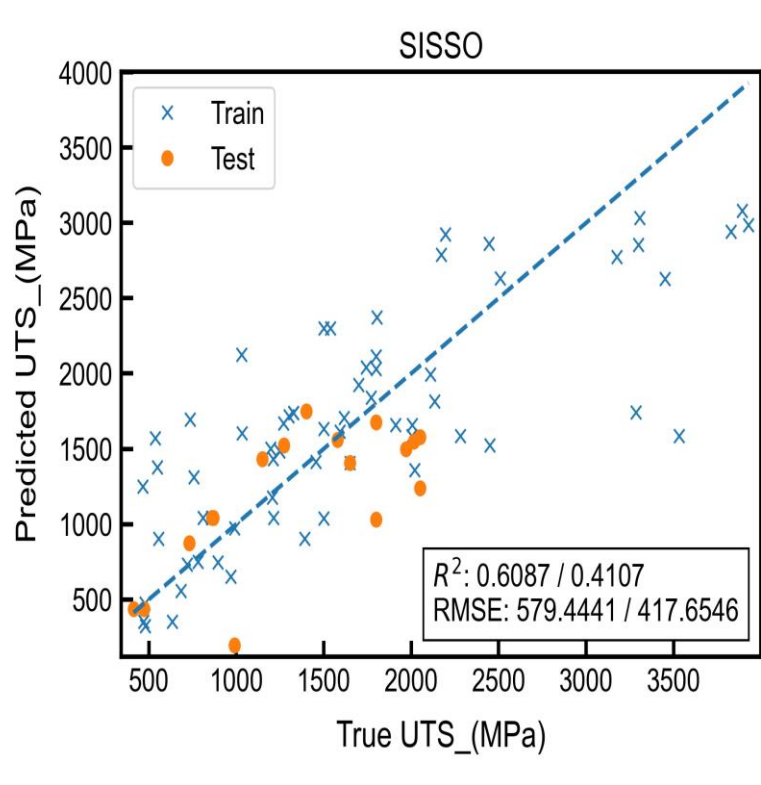
HV

SISSO Type	Multi-task	Single-task
Plot		
Equation	$-77254.61692(\text{Yang_delta} * \text{Shear_modulus_strength_model}) - 30.73113347 \text{abs}(\text{VEC_mean} - \text{Mean_cohesive_energy})$	$1754.503458(\text{Electronegativity_delta} - \text{Shear_modulus_strength_model}) + 11.33349587(\text{Mixing_enthalpy} - \text{VEC_mean})$

- Configuration: 1,2,1 is used for both.
- Multi-task performs worse in comparison to single-task.
- HV data has 86 nan values out of 193 entries. i.e. 44.5% nan values.

COMPARISON:

UTS(MPa)

SISSO Type	Multi-task	Single-task
Plot		
Equation	$-348775.8453(\text{Yang_delta} * \text{Shear_modulus_strength_model}) - 25.31638202 \text{abs}(\text{VEC_mean} - \text{Mean_cohesive_energy})$	$35335.70666(\text{Radii_gamma} * \text{Shear_modulus_delta}) - 33784.28339(\text{Shear_modulus_delta})$

- Configuration: 1,2,1 is used for both.
- Multi-task performs worse in comparison to single-task.
- UTS data has 111 nan values out of 193 entries. i.e. 57.5% nan values.



NOTABLE RESULTS

- It can be inferred that multi-task performance is close to single-task.
- It appears that columns with fewer missing entries benefit more from multi-task SISSO.
- Similarly, it appears that columns with more missing entries lose accuracy from using multi-task SISSO. It could be that the larger target datasets corrupt the accuracy by holding a larger influence on the selection of features.
- It is to be noted for the comparison that multi-task had been run without OMP (with multiple targets, it might not be possible to specify correlation in the given software) and without cross-validation (due to software limitations). Having both would further help validate the hypotheses above as single-task had been run with both.