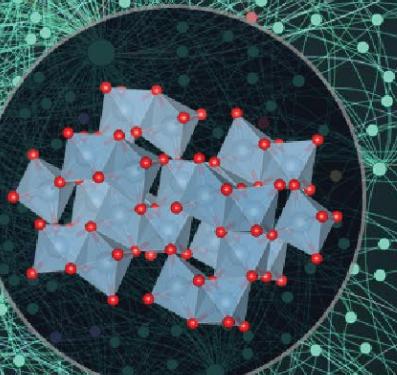
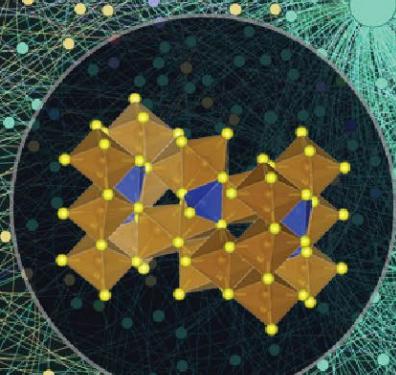
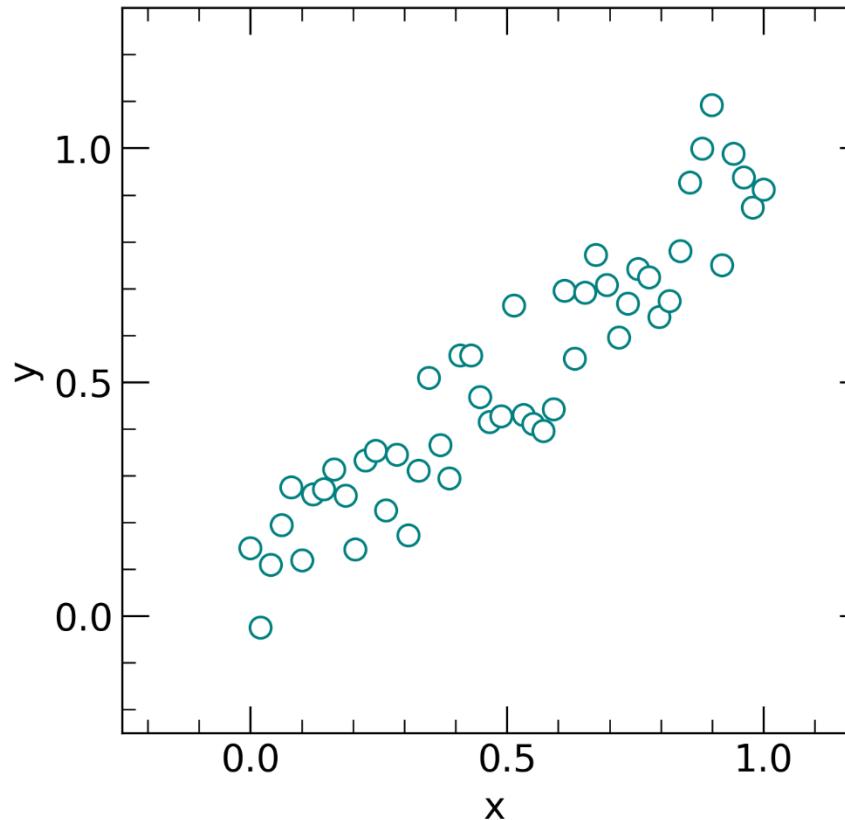


Extrapolation in Materials Informatics

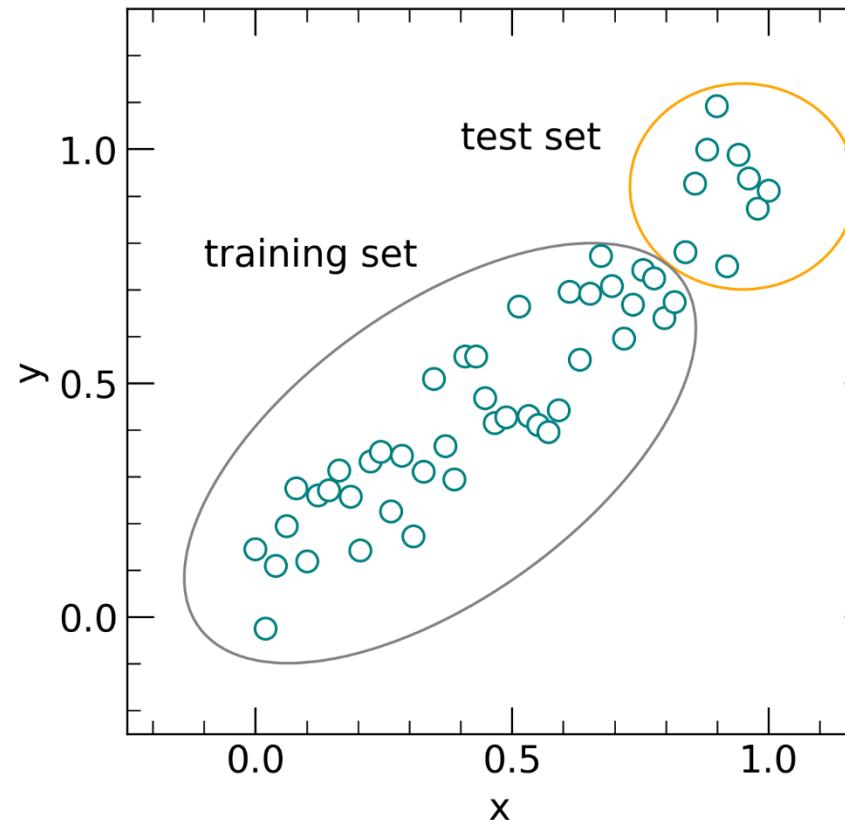


Most machine learning approaches fail with extrapolation

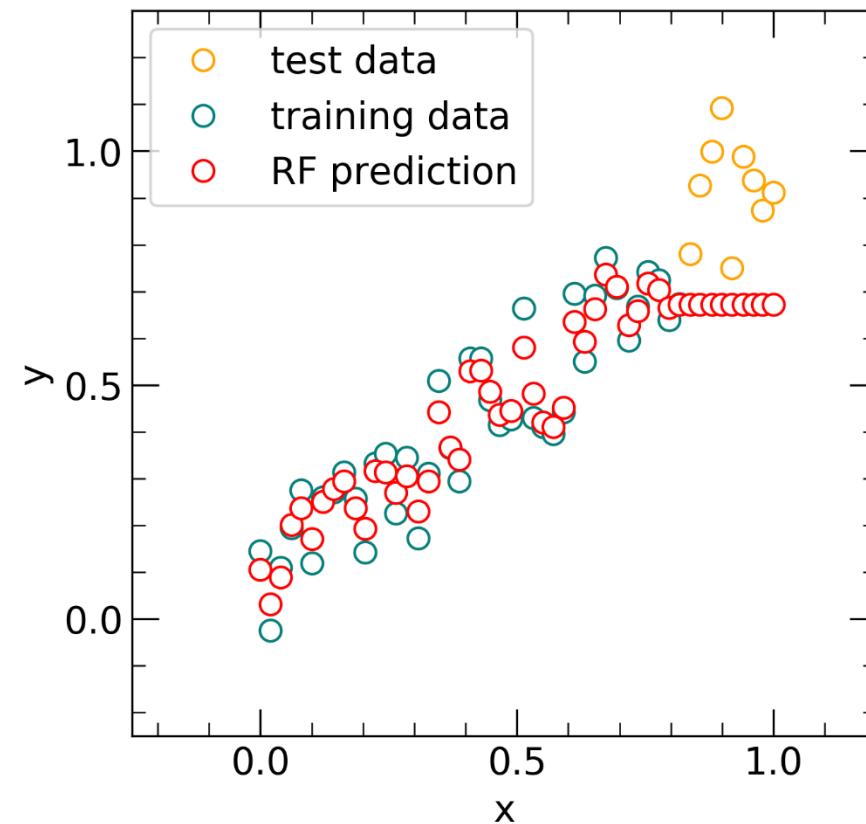
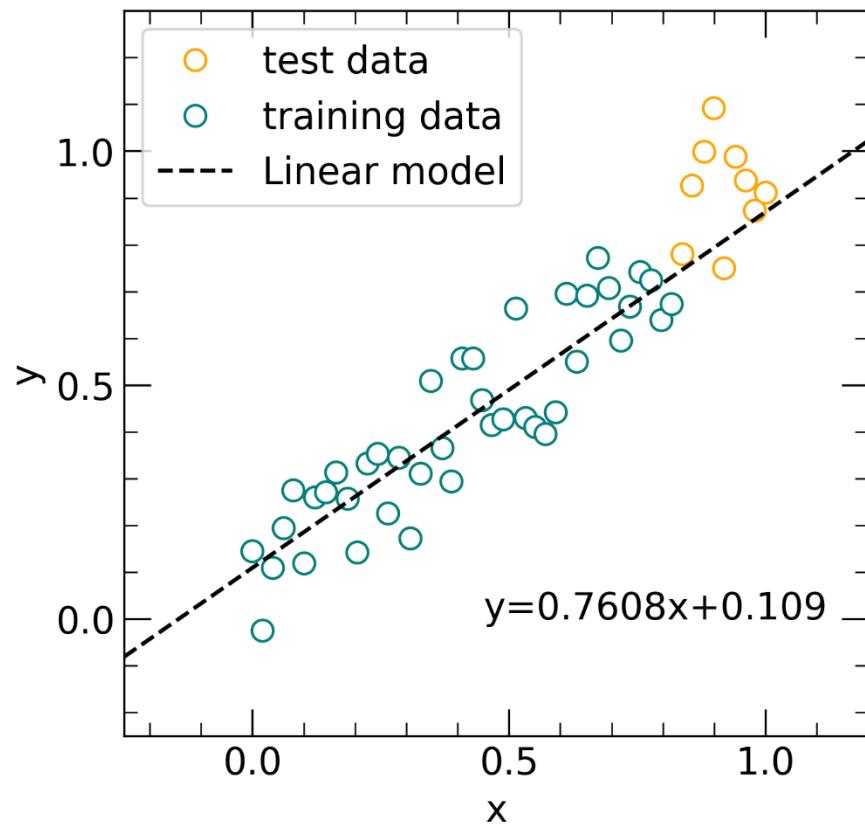


Most machine learning approaches fail with extrapolation

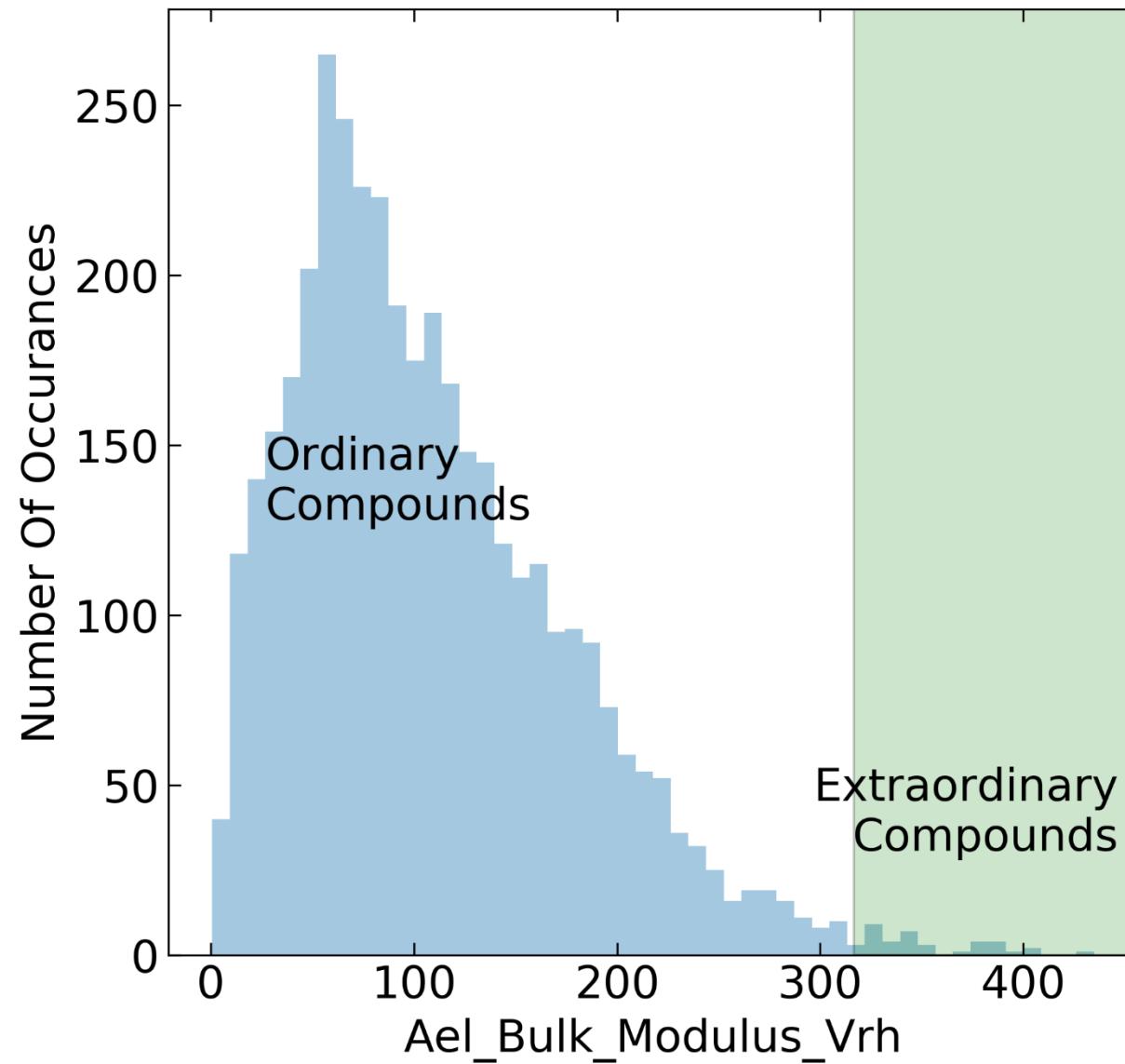
Will we be able to predict the validation set if these are extrapolations?



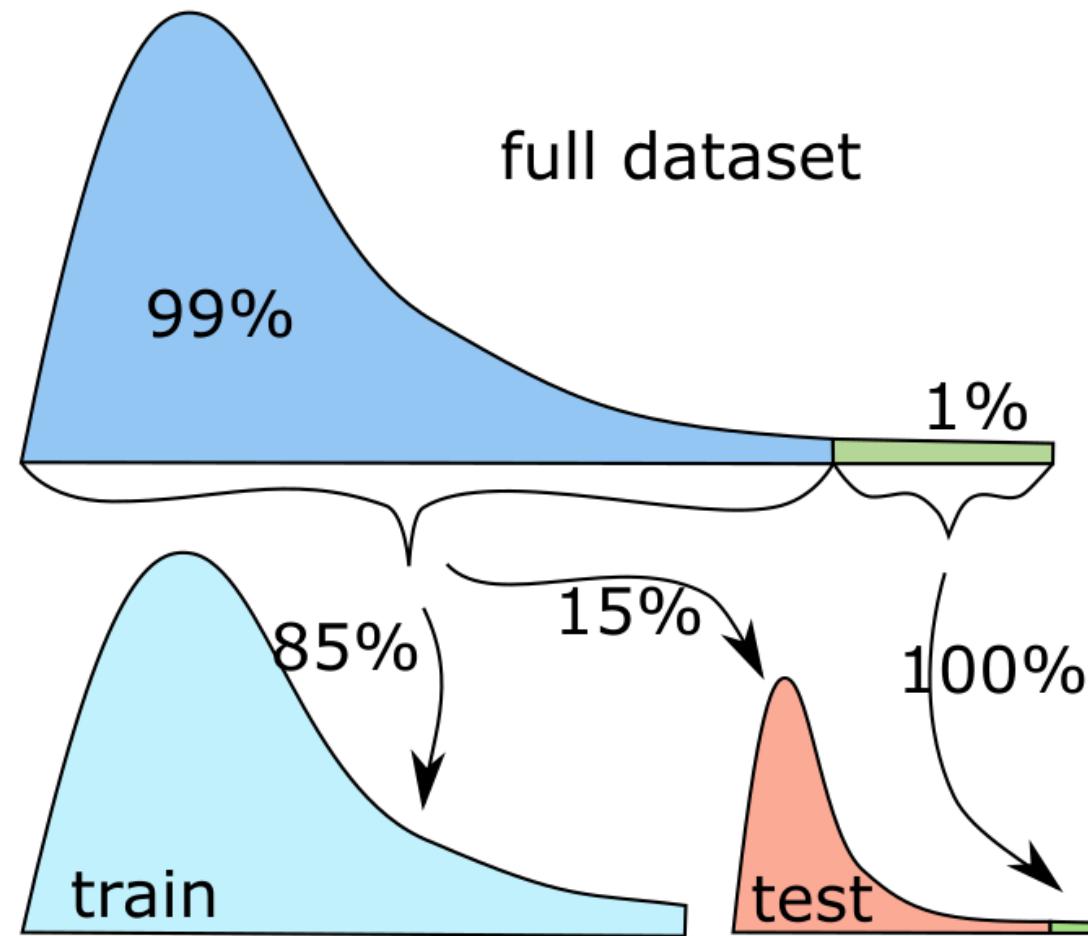
Most machine learning approaches fail with extrapolation



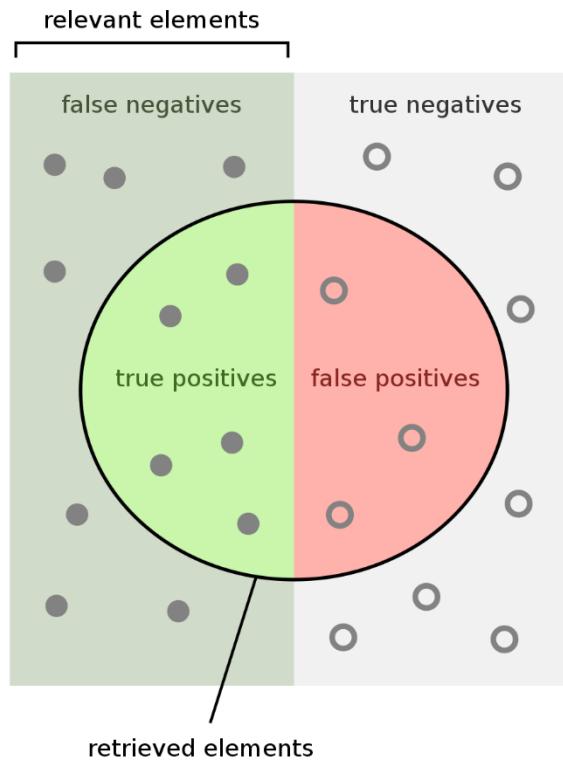
Will machine learning be able to identify extraordinary compounds?



We can test to see how well extrapolation works



We first need to define precision vs recall



How many retrieved items are relevant?

Precision = $\frac{\text{green circle}}{\text{red and green circles}}$

How many relevant items are retrieved?

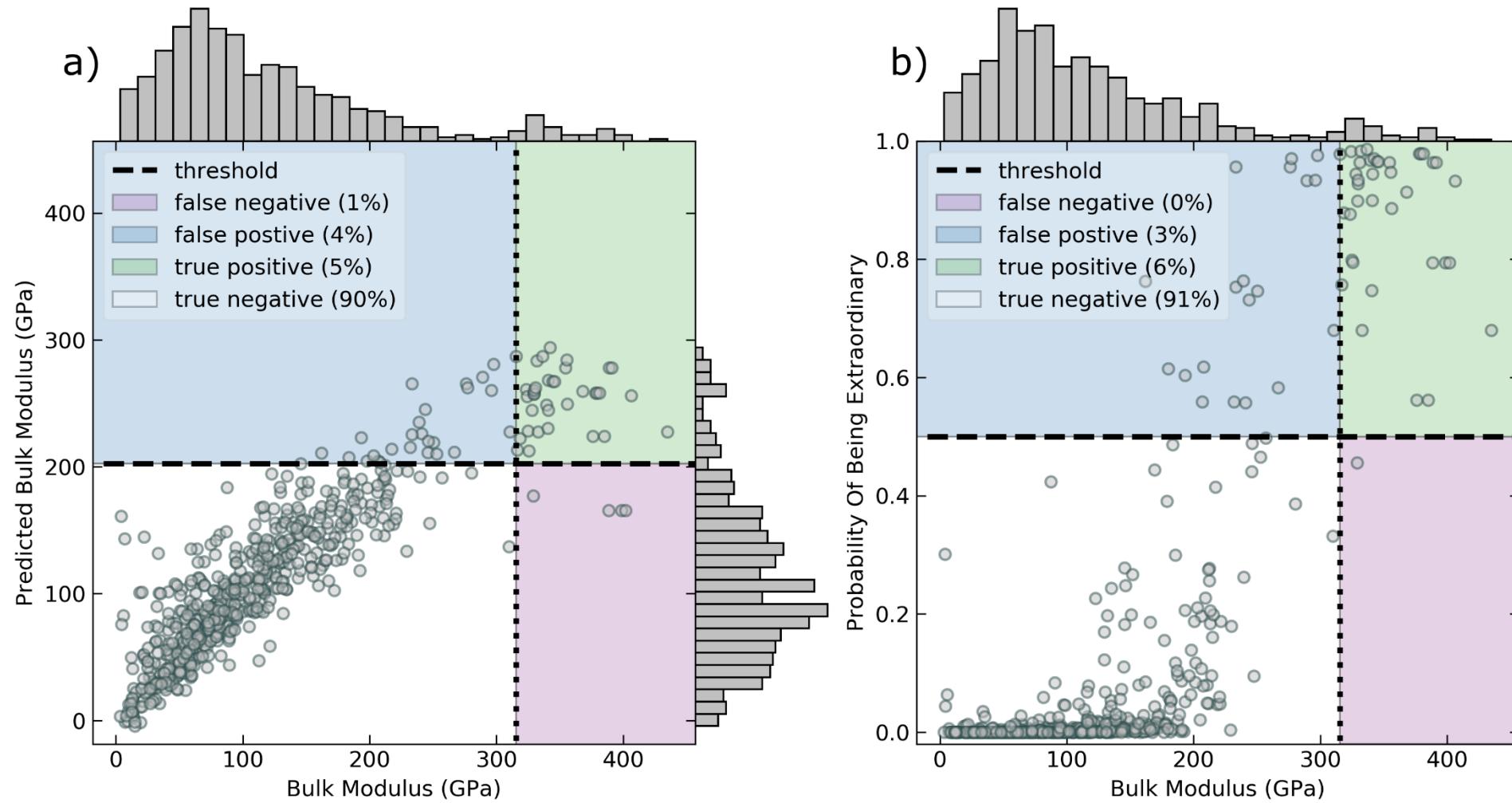
Recall = $\frac{\text{green circle}}{\text{green and blue circles}}$

We need to account for both precision and recall!

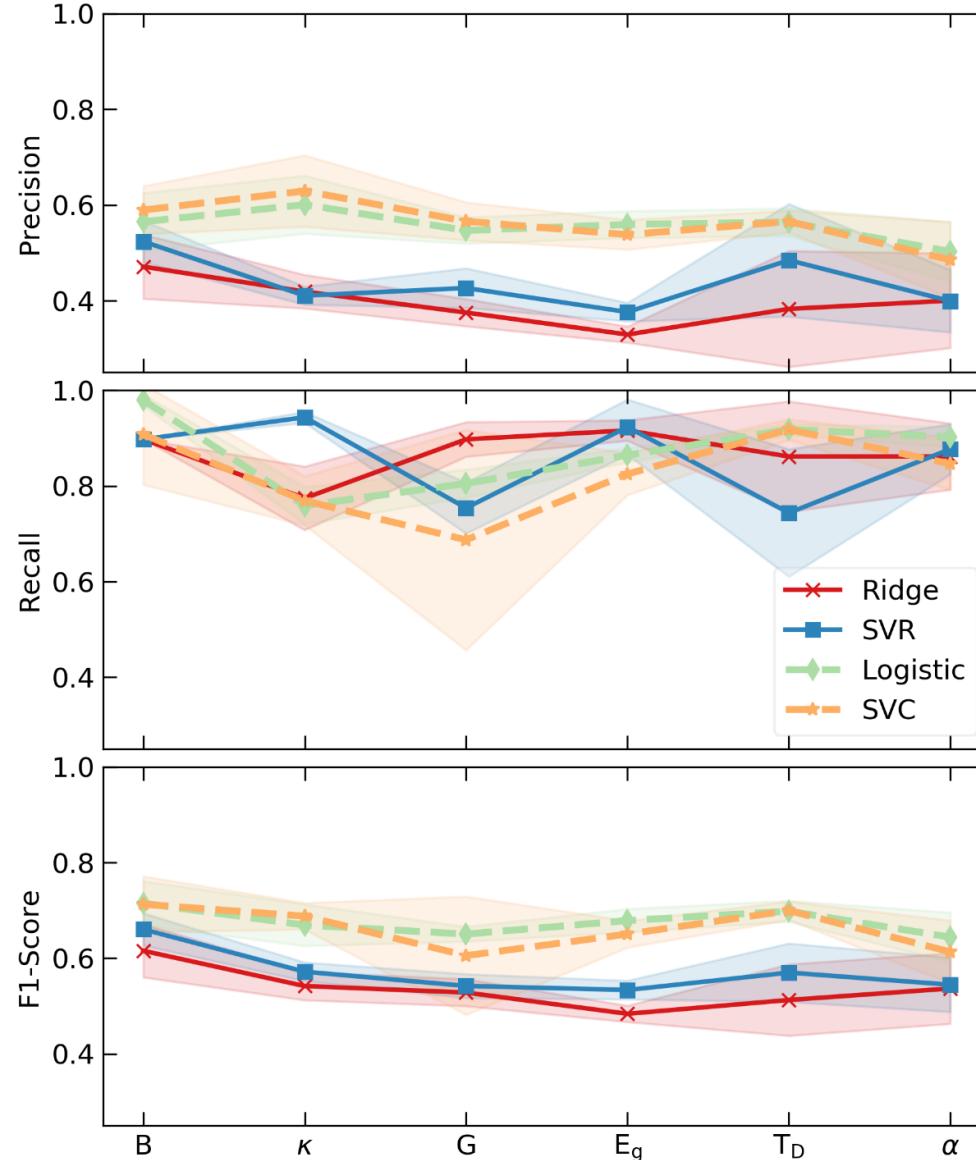
If we take their harmonic mean, we get F1 score

$$F_1 = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

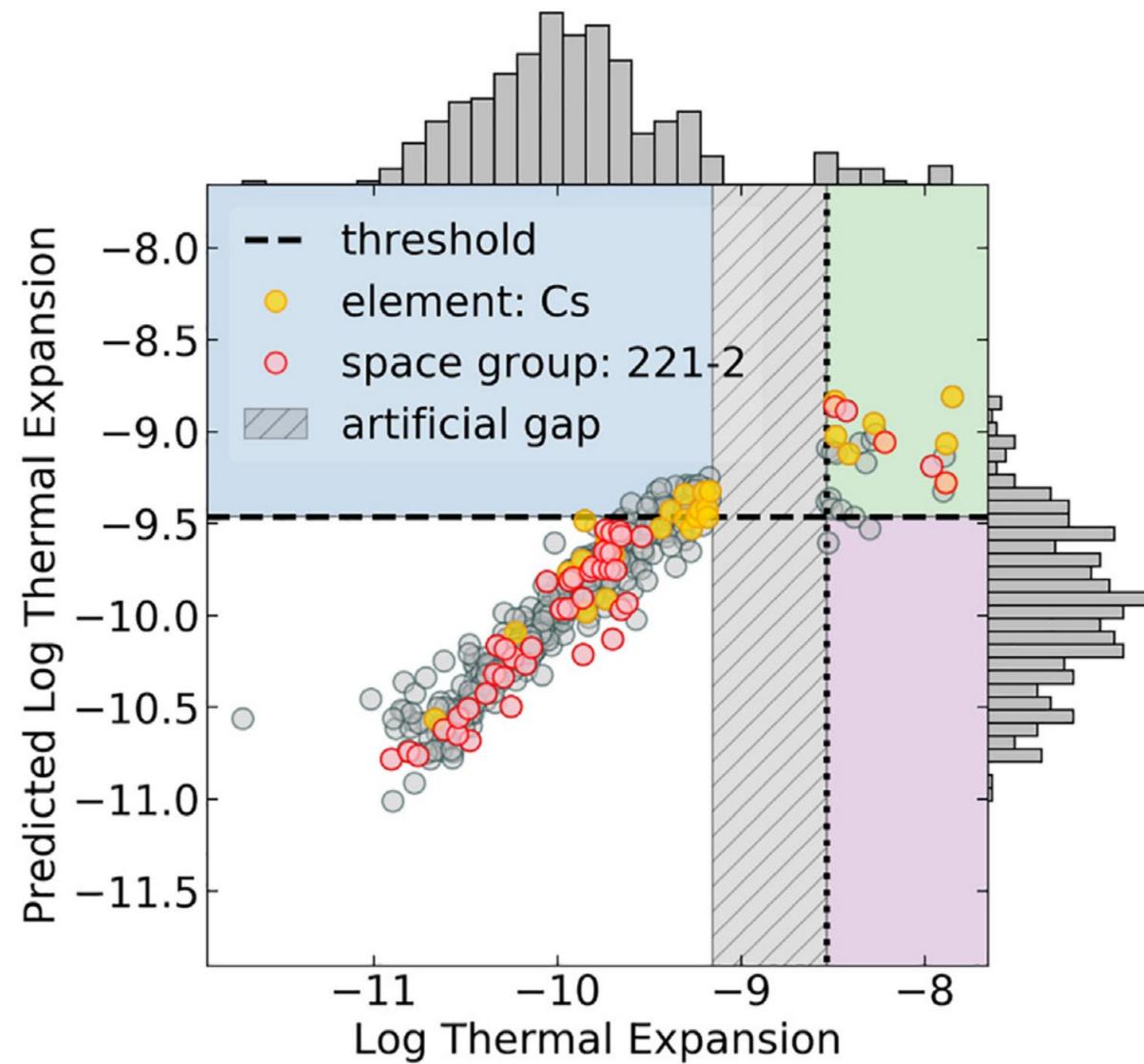
Precision and recall both suggest extrapolation is possible



Classification outperforms regression on almost all properties!



Extrapolation goes beyond target values!



Extrapolation goes beyond target values!

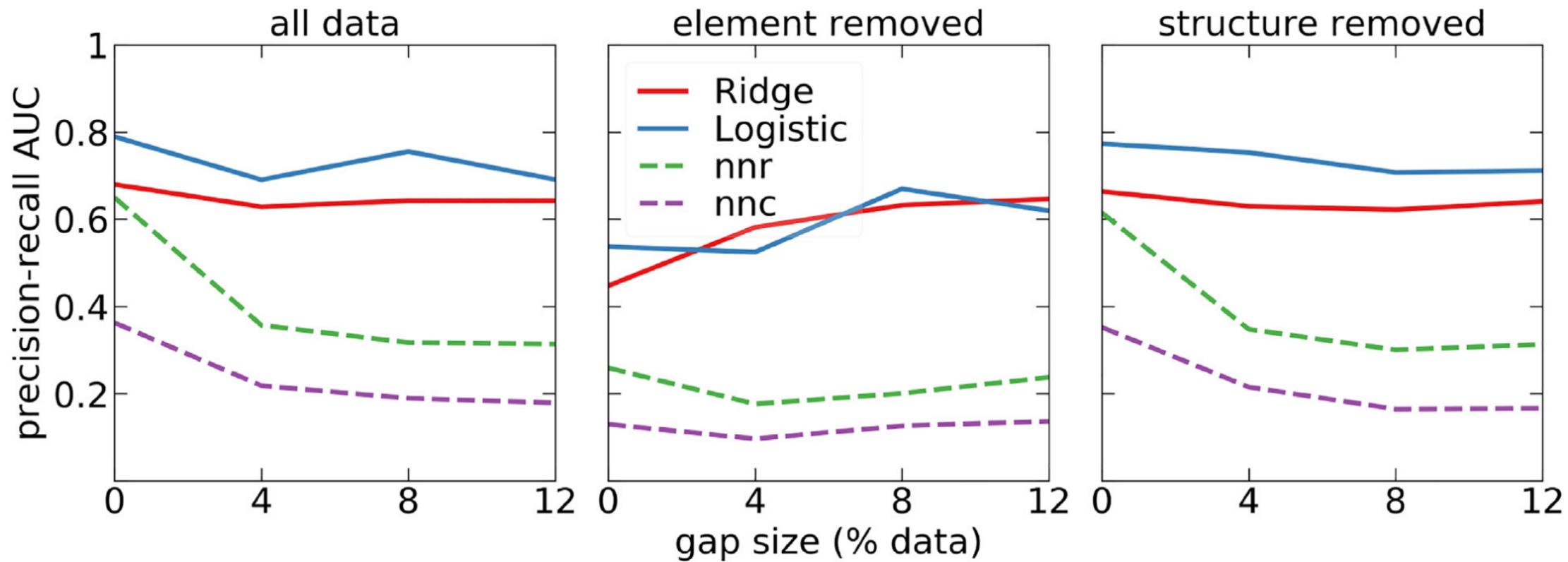


Fig. 5. The precision-recall Area-under-the-curve (AUC) is computed as the average across all properties. The three extrapolation cases are considered for all gap sizes using four different model types, ridge regression (Ridge), logistic regression (Logistic), nearest-neighbor regression (nnr), and nearest-neighbor classification (nnc). (*left*) The entire training data is used (except the gap), (*center*) the most prevalent element is removed from training, and (*right*) the most prevalent structure is removed from training.

Compare DFT to experiment to ask if I.I.D. is valid?

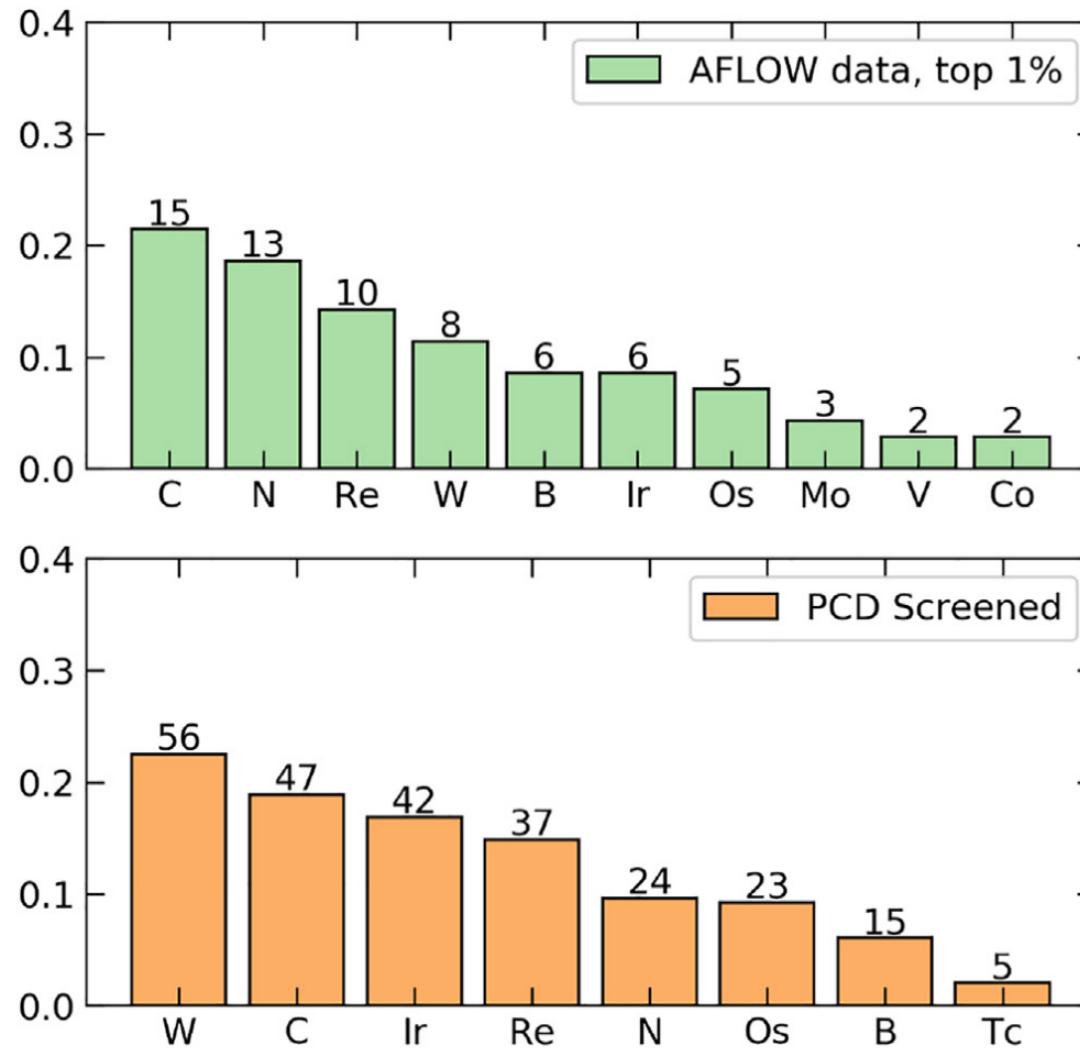


Fig. 6. The ranked elemental prevalence of compounds labeled extraordinary in the original AFLOW data versus those screened from the PCD.

Clustering

