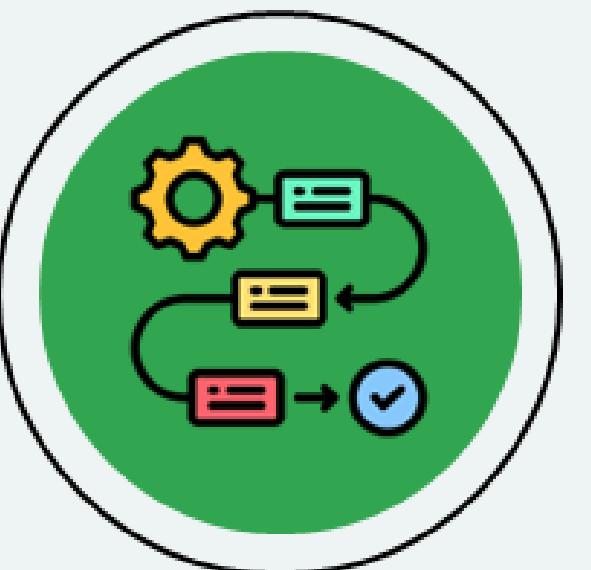


DATA mining
“mELTING pOINT”



AGENDA



**Problem
Motivation
Research Question**

**Data Exploration
Feature
Engineering**

**Methodology
Models
Results**

**Interpretation
Key Learnings
Future Work**



**Problem
Motivation
Research Question**

JOHN HEDENGREN · COMMUNITY PREDICTION COMPETITION · 15 DAYS TO GO

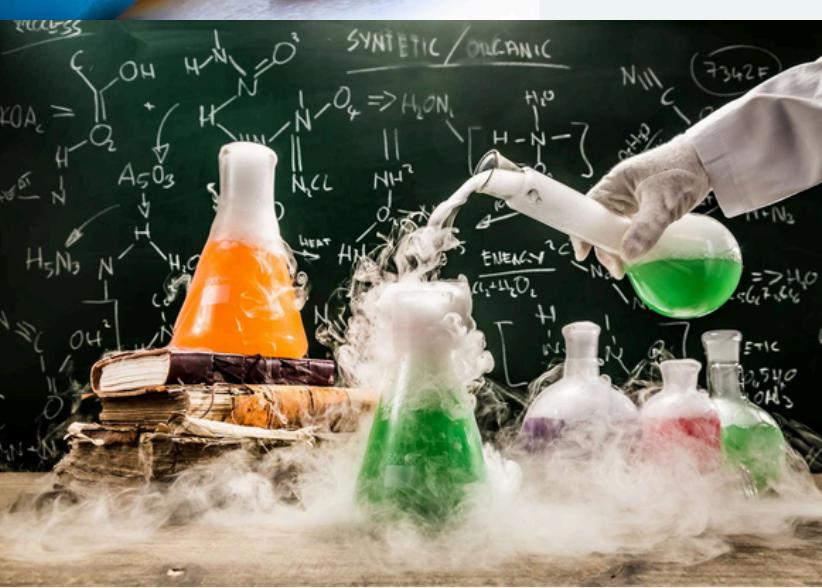
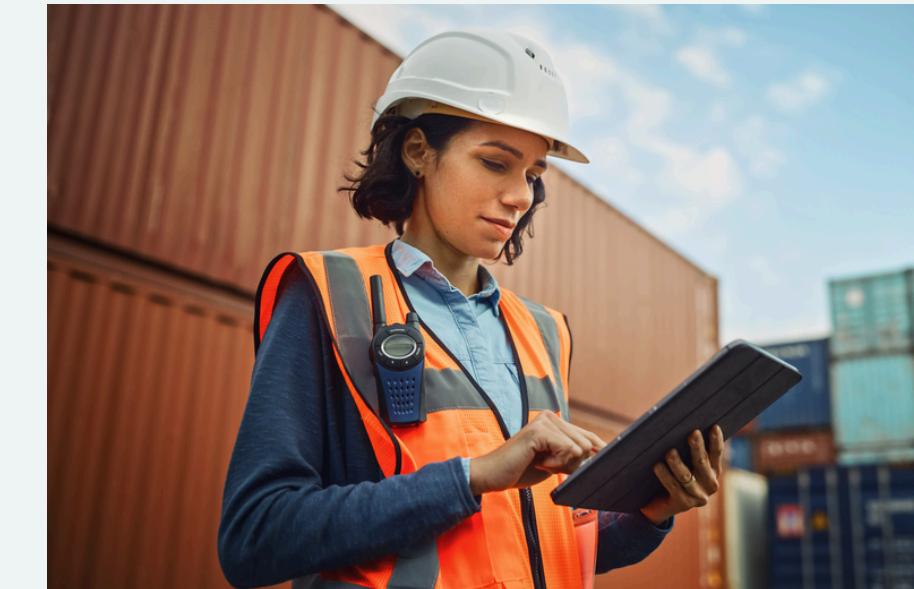
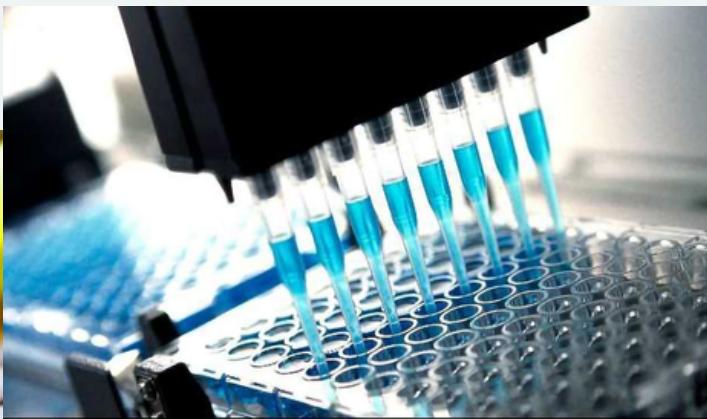
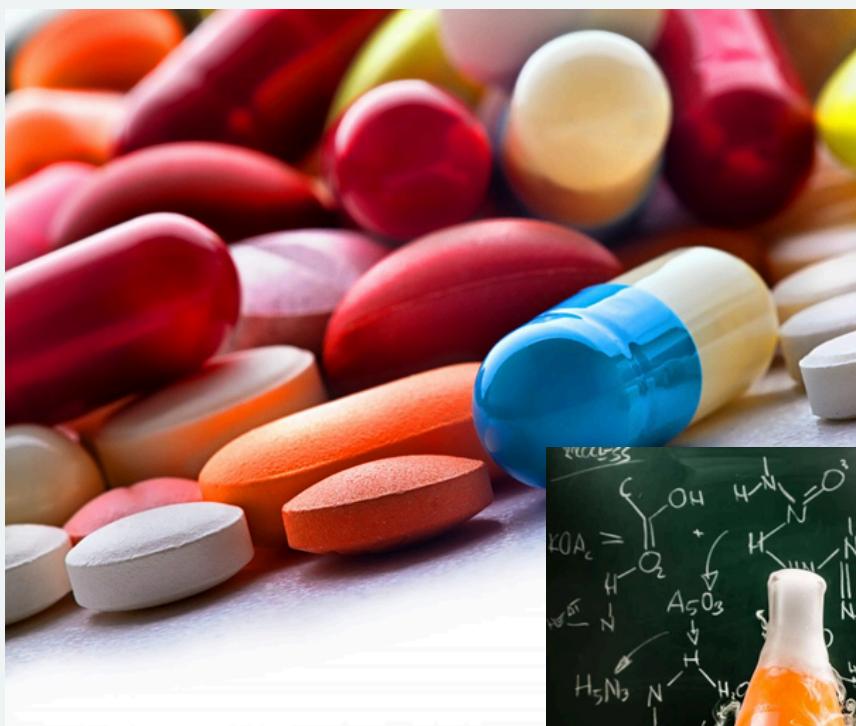
[Submit Prediction](#)

...

CONTEXT & RELEVANCE

Thermophysical Property: Melting Point

Your goal is to build ML models that predict melting point in temperature units of Kelvin for organic compounds given molecular descriptors.



PROBLEM

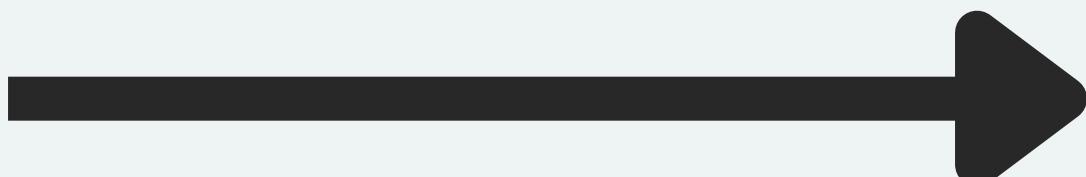


GOAL



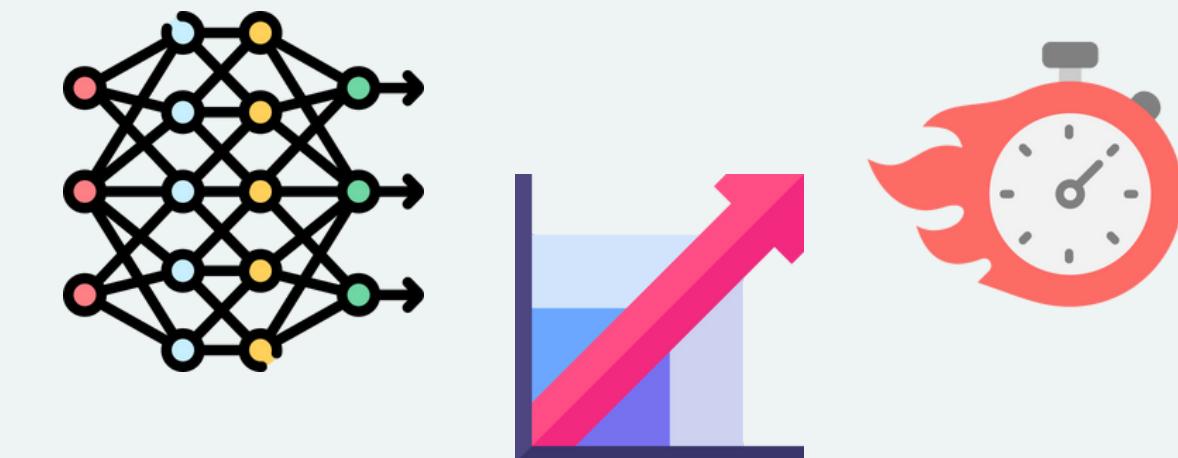
Experimental Determination

- Laboratory measurements are expensive
- Experiments require significant time
- Large chemical space cannot be fully tested



Machine Learning Based Prediction

- Predict MT without physical experiments
- Fast and scalable evaluation
- Enables early stage screening of compounds



CENTRAL RESEARCH QUESTION

“

To what extent can the melting point of organic molecules be predicted from molecular structure using machine learning methods?”

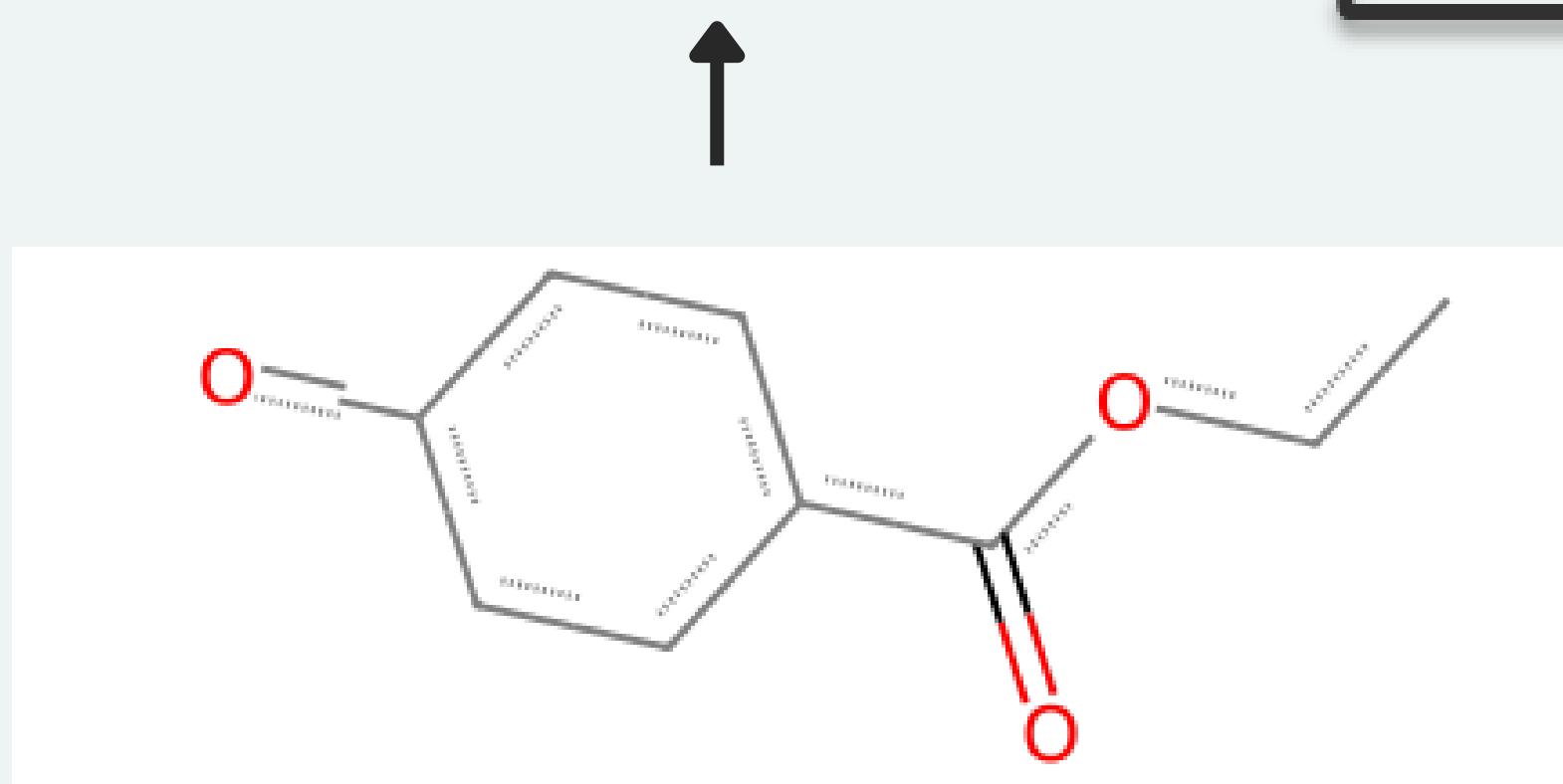




**Data Exploration
Feature
Engineering**

WHAT ARE WE TRYING TO PREDICT?

ID	SMILES	Tm	Group 1	Group 2	Group 15	Group 30	...	Group 424
2175	CCOC(=O)c1ccc(O)cc1	213.15	0	1	0	2	...	0



Two dimensional representation of an organic molecule

Groups	Best Pattern	Correlation
Group 1	CH3	(0.915379)
Group 2	CH2	(0.923681)
Group 15	Aromatic_C	(0.897967)
Group 30	Phenol_OH	(0.985509)
Group 45	COO	(0.269556)
Group 424	Aromatic C	(0.278946)

DATA SET

Files

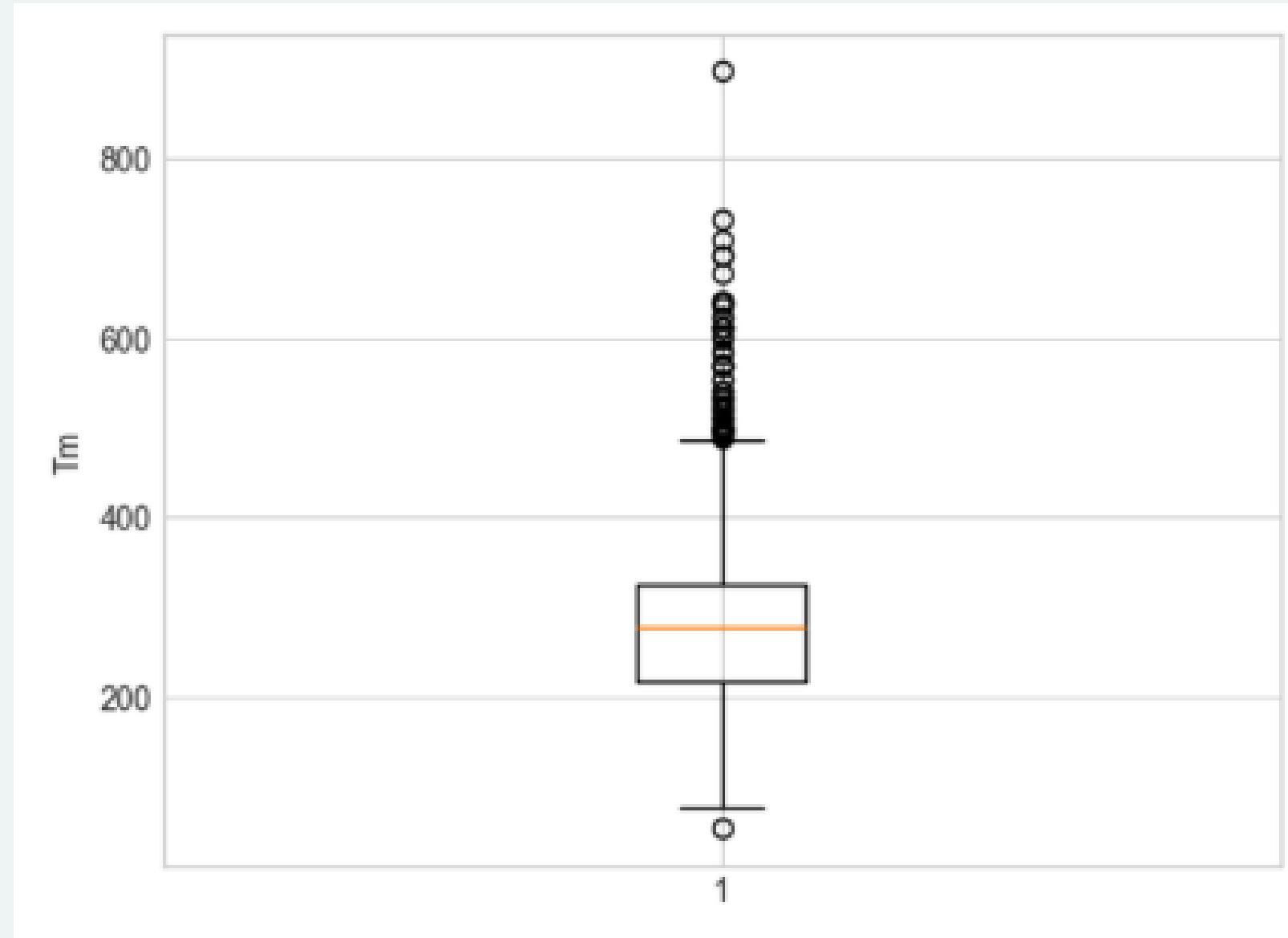
- train.csv : Features + target (Tm)
 - test.csv : Features only, no target
 - sample_submission.csv : Template with columns [id,Tm]

Columns

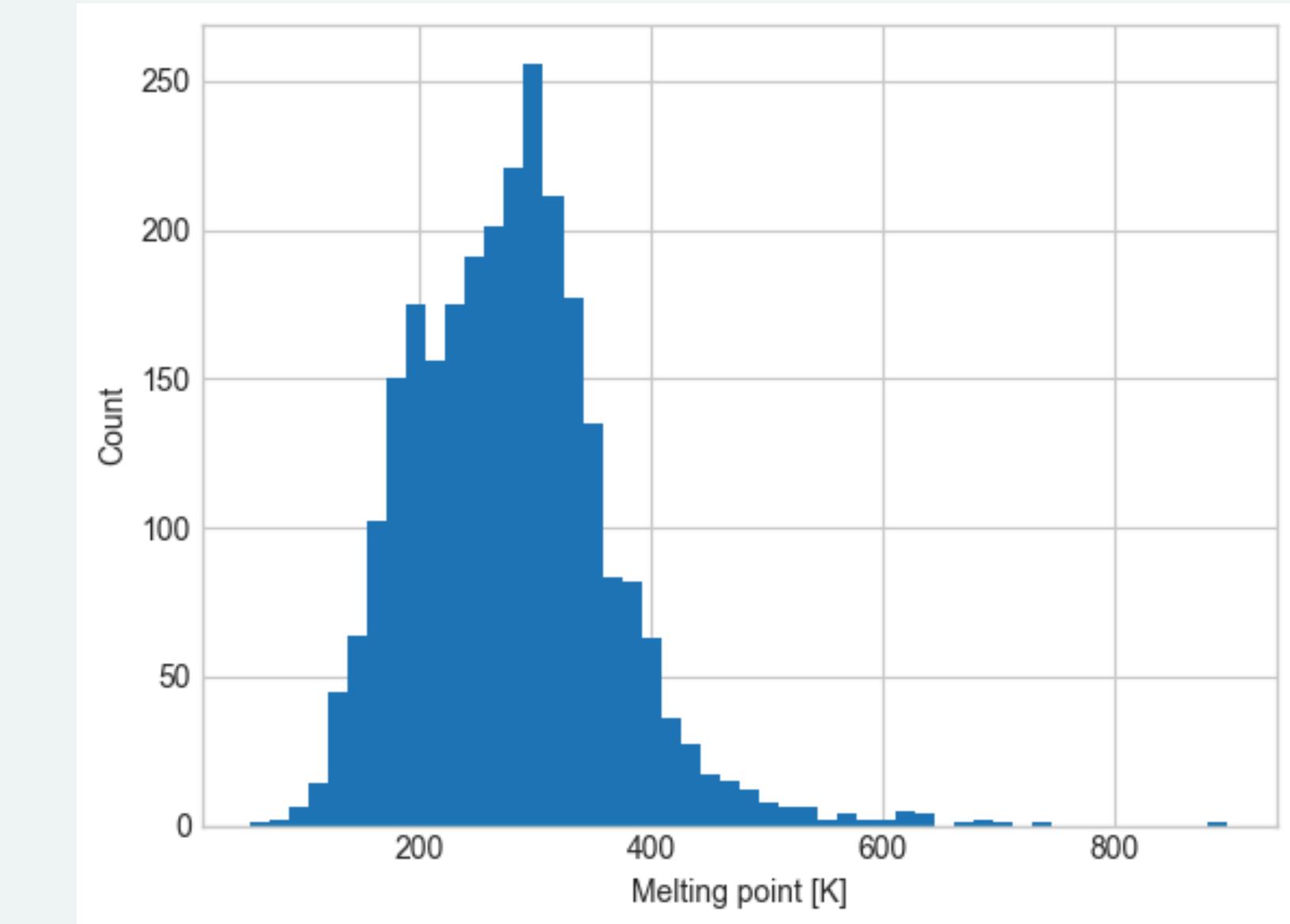
- id : unique ID
 - SMILES : molecular string
 - Group 1..N : descriptor features
 - Tm : melting point (Kelvin) [train only]

DATA EXPLORATION

Distribution of the Target Variable Tm



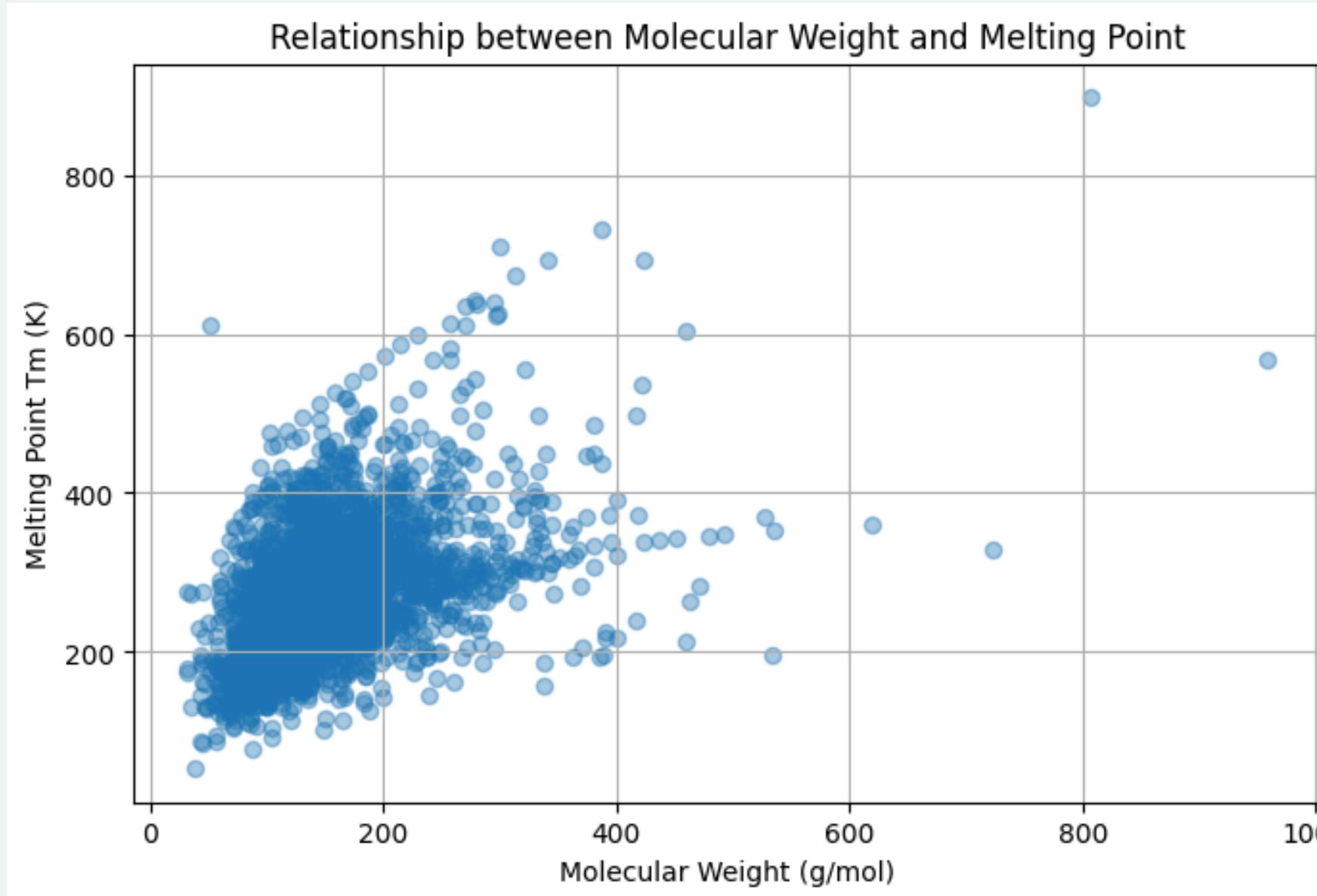
- Max $T_m = 897.15\text{K}$
- Min $T_m = 53.54\text{K}$



- right skewed distribution
- The mean $T_m \approx 278\text{ K}$

FEATURE ENGINEERING: WEIGHT

We used the SMILES to derive the mol mass of the molecules and added this as a new feature

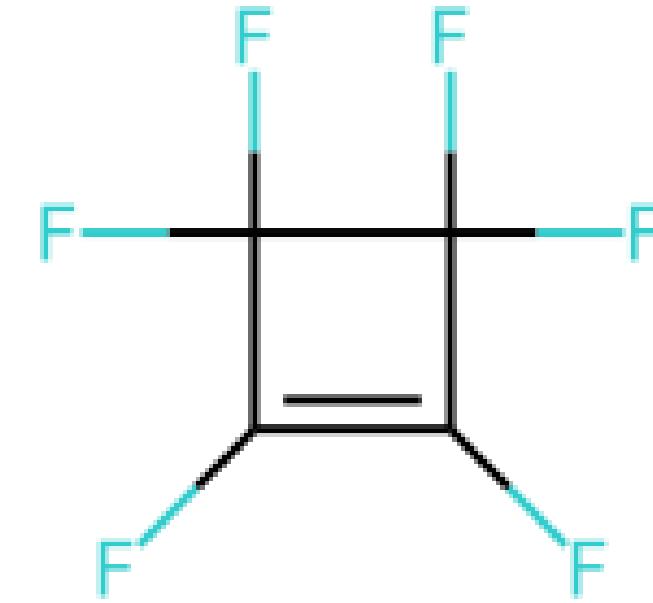


index	SMILES	MolWeight	Tm
0	FC1=C(F)C(F)(F)C1(F)F	162.032	213.15
1	c1cccc2c(c1)ccc3Nc4cccc4c23	217.271	407.15
2	CCN1C(C)=Nc2ccc cc12	160.220	324.15
3	CC#CC(=O)O	84.074	351.15
4	CCCC(S)C	118.245	126.15

FEATURE ENGINEERING: STRUCTURE (1)

Example with one Tuple

SMILE: FC1=C(F)C(F)(F)C1(F)F



A chemist would recognize:

- Symmetry
- No H-Bridges
- negative electric charge
- ...

Expectation: High melting point

The Groups would look like this:

Group	Value
Aromatic ring	1
Halogen (F)	6
CH ₃	0
OH	0
C	4

Problem:

We know what we have but there is none relations to other groups.

FEATURE ENGINEERING: STRUCTURE (2)

We look at every Atom as a center. It's connected to other atoms with a bond

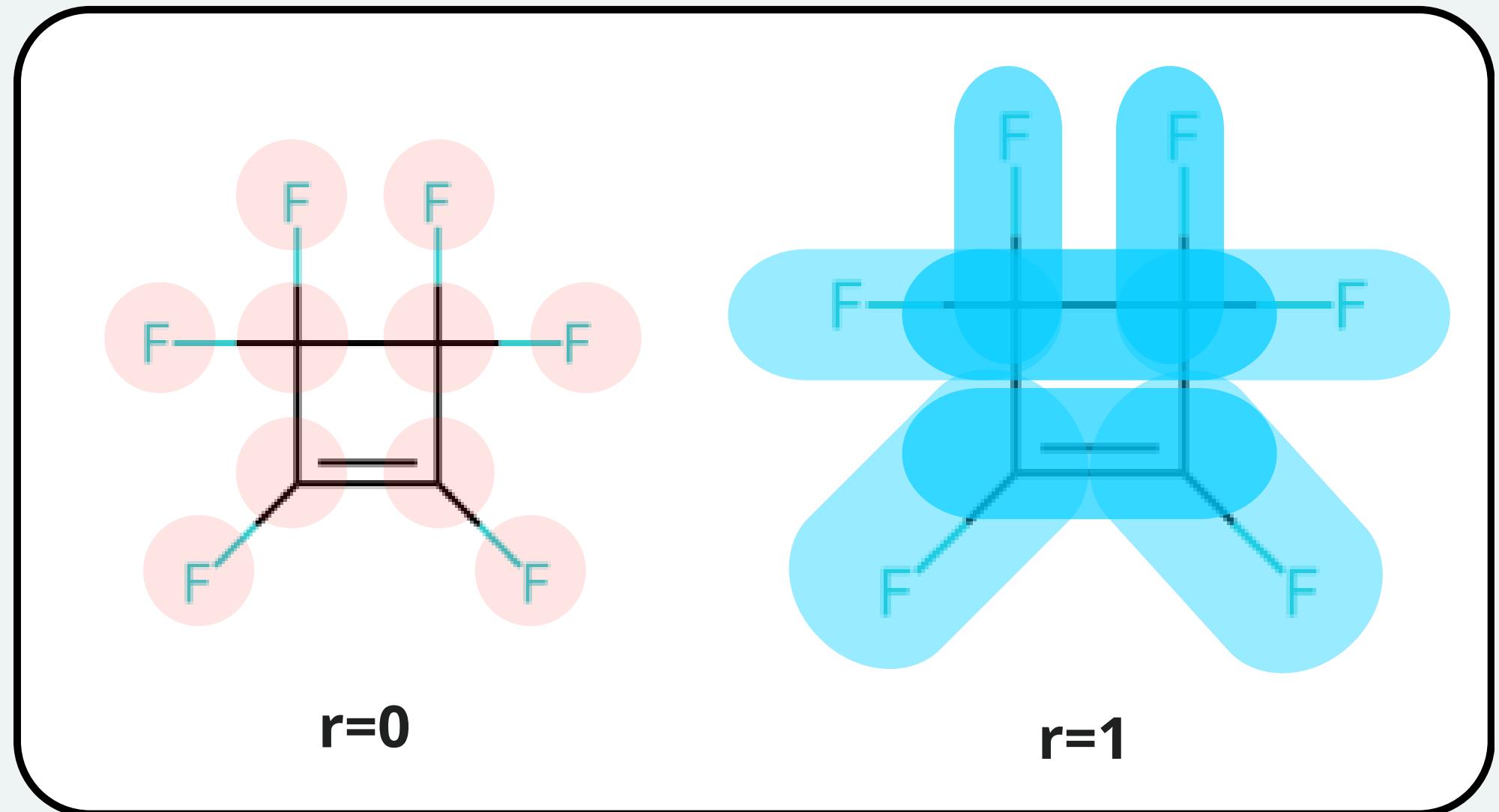
- A Feature which looks at the atom itself (**r=0**)
- A Feature which looks at the atom and the nearest neighbor (**r=1**)
- A Feature which looks at itself, it's neighbor and the direct neighbor of the neighbor (**r=2**)

Two Examples:

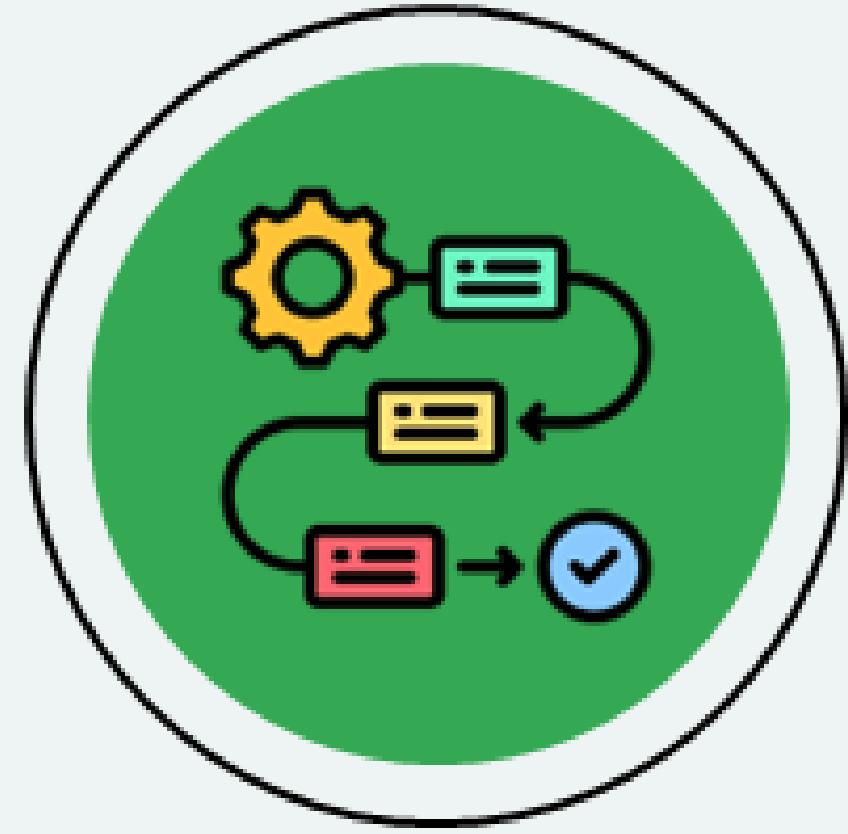
Feature “C” (r=0): 4

→ because there are 4 “C” which are sp₂ hybridized and part of a ring C (sp₂, ring): 4

For r=1 we look at the bonds: F-C (sp₂, ring): 6 and also C-2*F(ring): 2



→ We got 2432 new Features!



Methodology
Models
Results

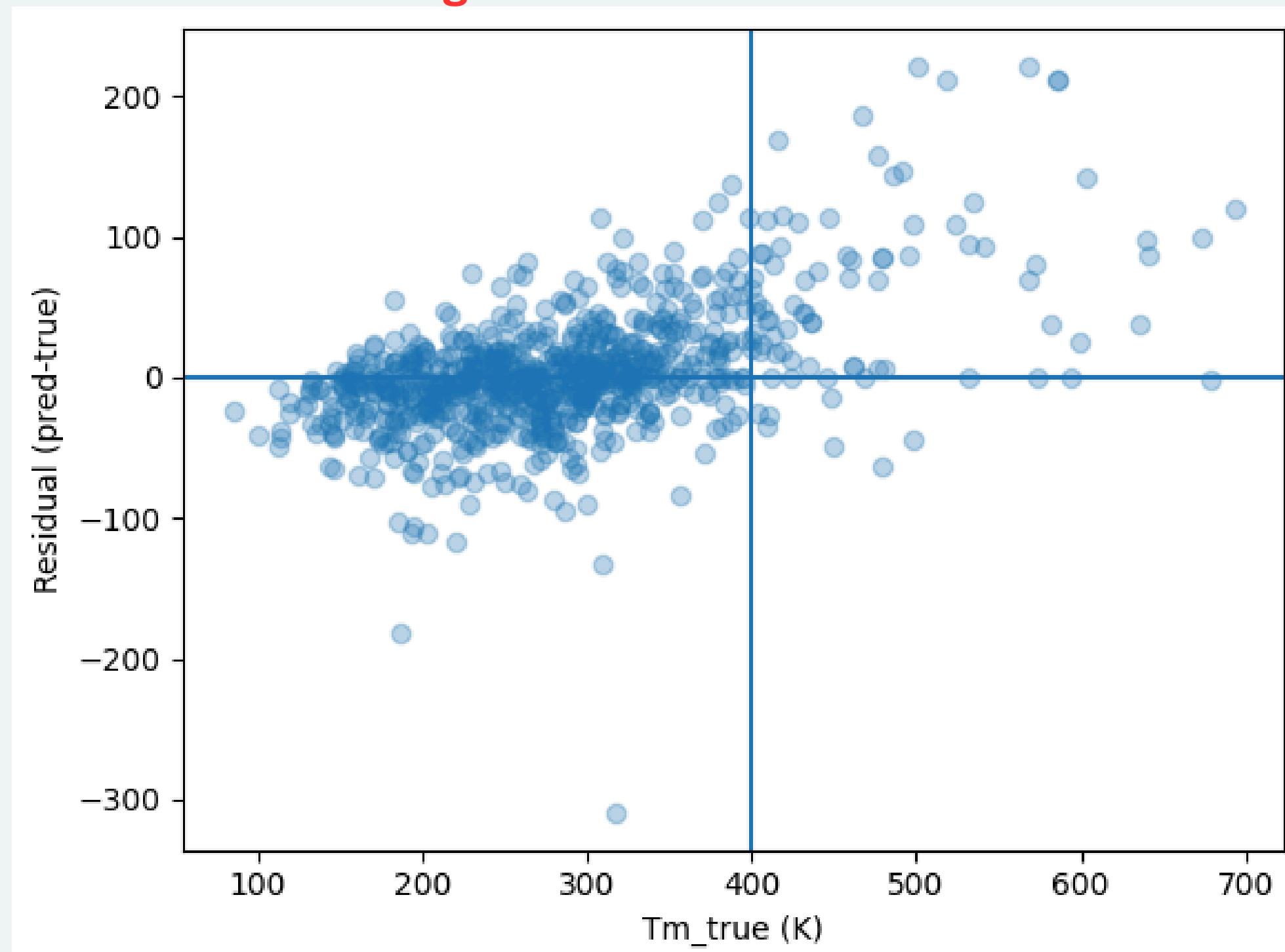
MODEL APPROACH



- **Gradient Boosted Decision Trees (XGBoost)**
- Suitable for high-dimensional, non-linear tabular data
- Input: SMILES-derived features (descriptors + fingerprints)
- Evaluation metric: Mean Absolute Error (MAE)

→ We recognized that our data has many outliers as soon as the Tm gets higher!

This gave us an MAE of about 25K

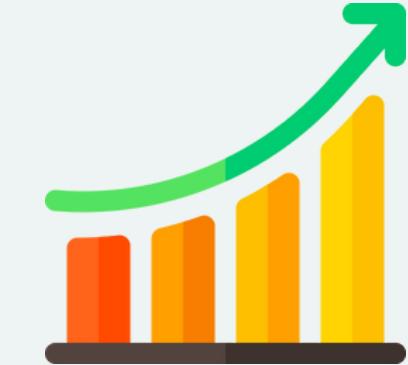


TWO-STAGE (RESIDUAL) MODEL (1)



→ Stage 1 – Base Model

- XGBoost regression on SMILES-derived features
- Sample weighting to emphasize high melting points ($T_m > 450$ K)
- Produces initial melting point prediction



→ Stage 2 – Residual Model

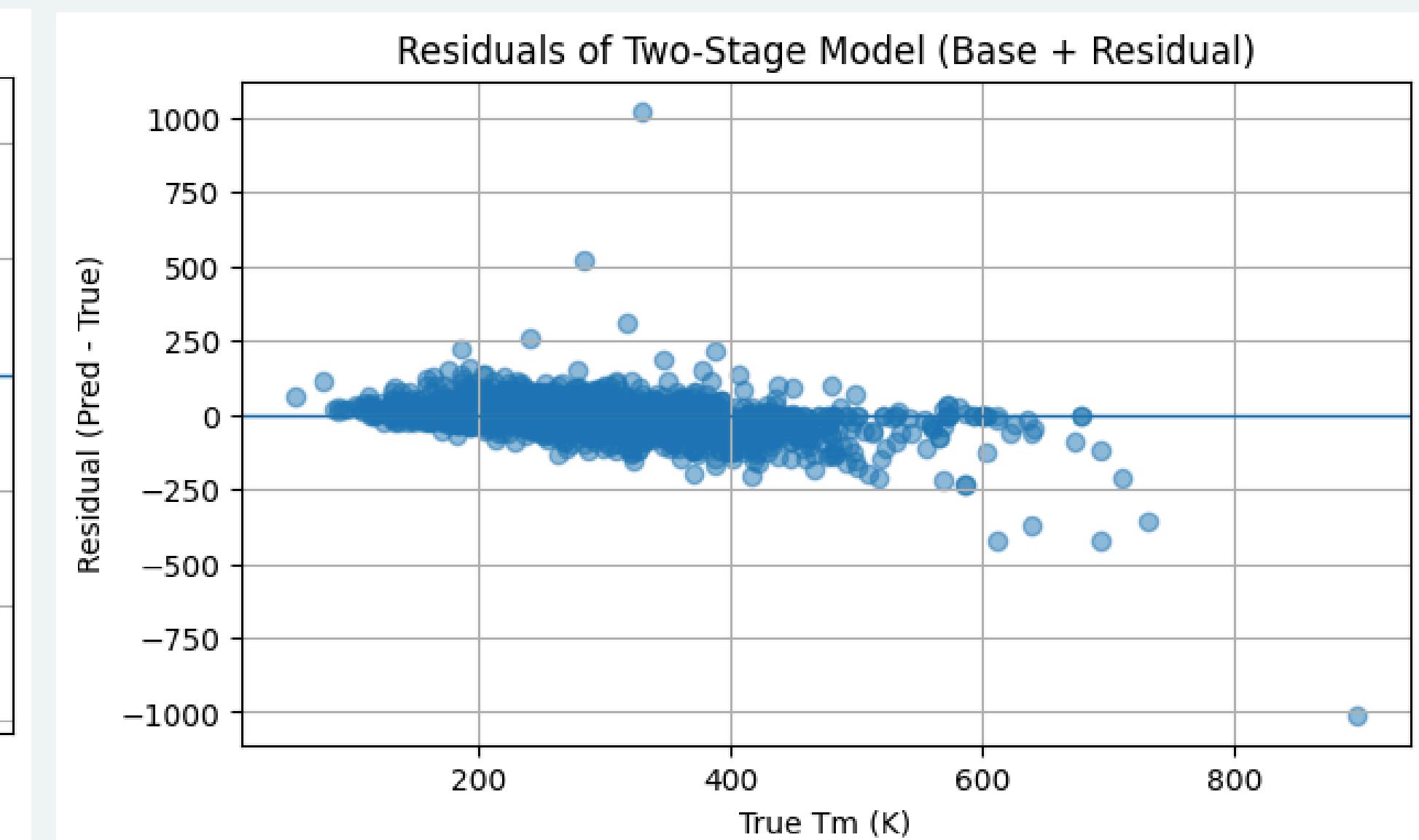
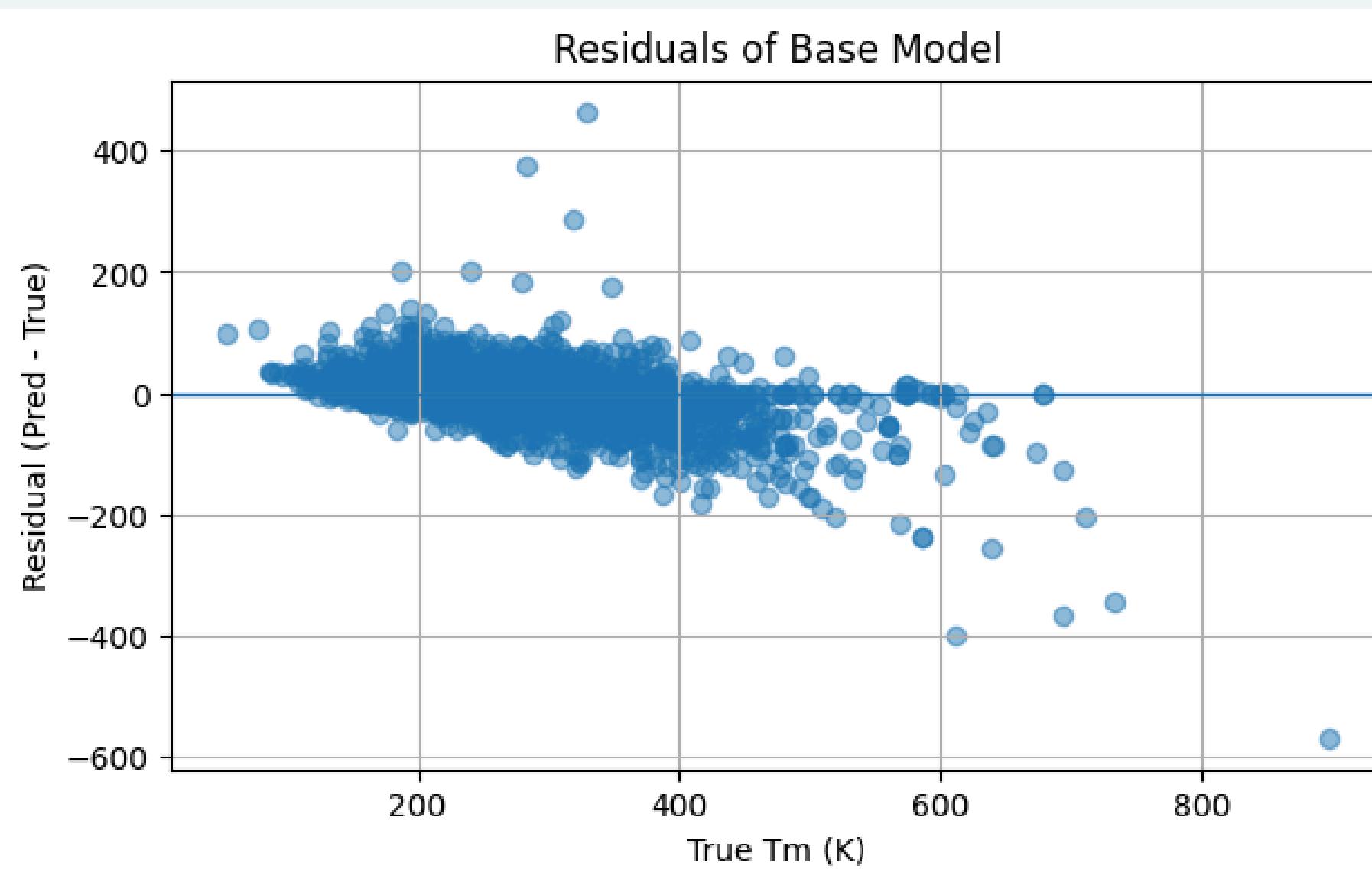
- Second XGBoost model trained on prediction residuals
- Learns systematic errors not captured by the base model

→ Final Prediction

Final T_m = Base prediction + Residual correction

With this approach we managed to achieve an MAE of 23.78K

TWO-STAGE (RESIDUAL) MODEL (2)



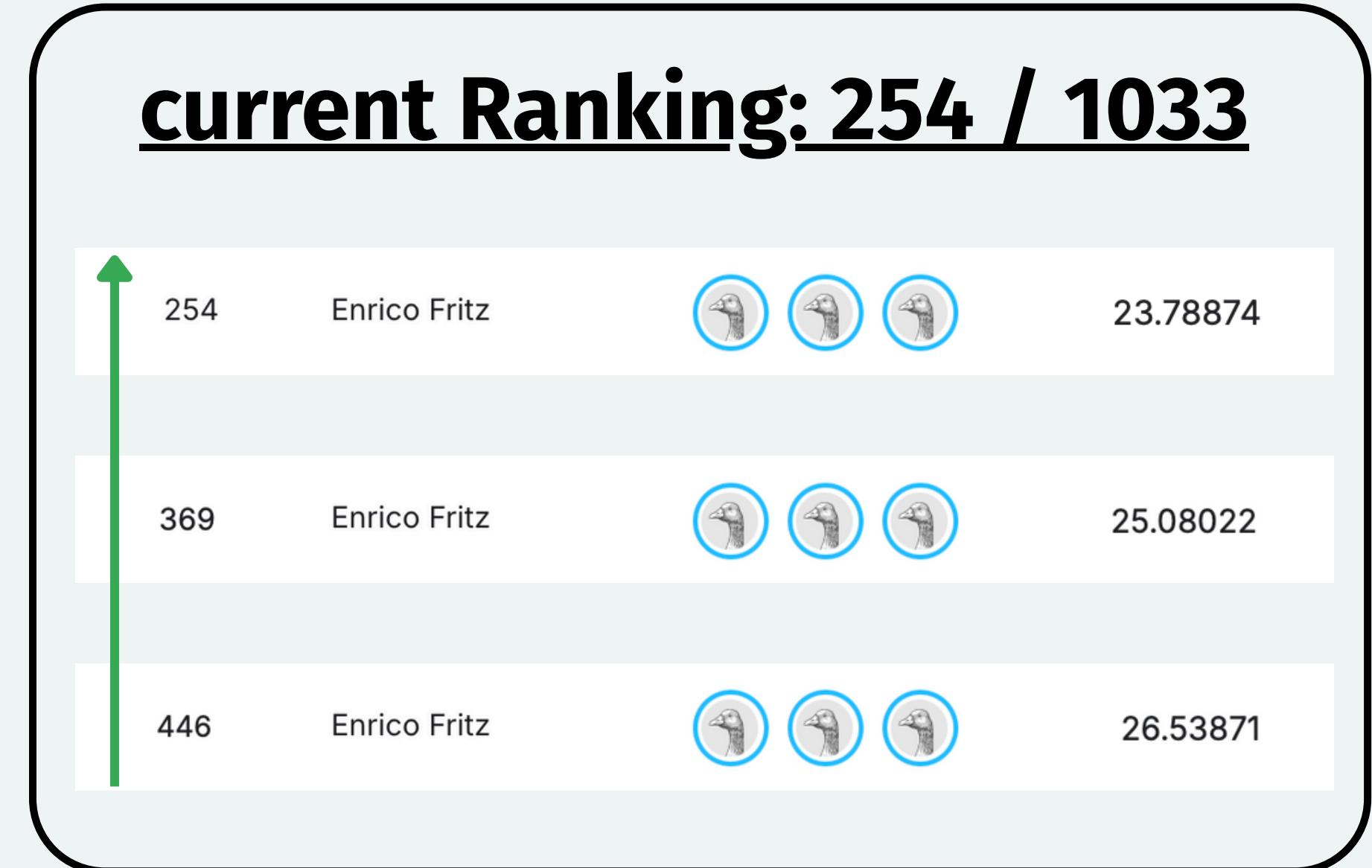


**Interpretation
Key Learnings
Future Work**

RESULTS OVERVIEW

Model	MAE (K)
Neural Network	~31
XGBoost Baseline	~25.4
Two-Stage XGBoost	~23.8

current Ranking: 254 / 1033



ANSWERING THE RESEARCH QUESTION

“

“To what extent can the melting point of organic molecules be predicted from molecular structure using machine learning methods?”

Answer:



- Melting points can be predicted to a large extent
- Reliable predictions for the majority of compounds
- Strong performance in the mid-range (200–400 K)
- Reduced accuracy for rare and extreme molecular structures



FUTURE OUTLOOK



Scale the model on an AWS server

Thank you!

Any questions?

