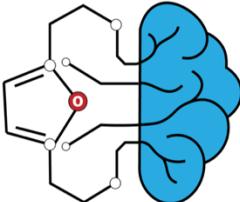
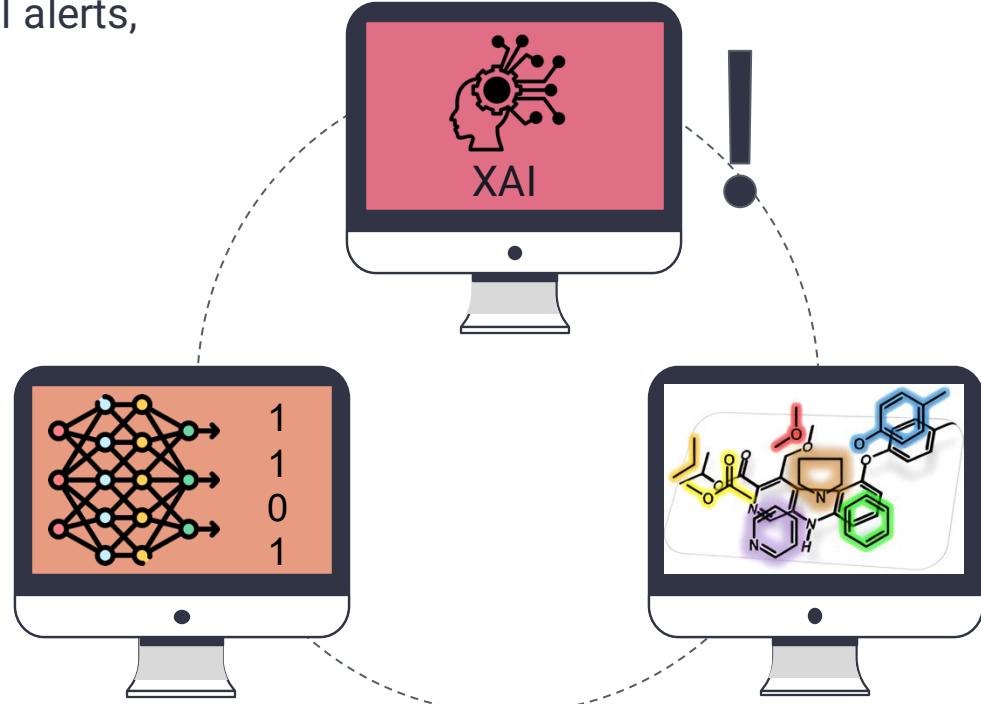


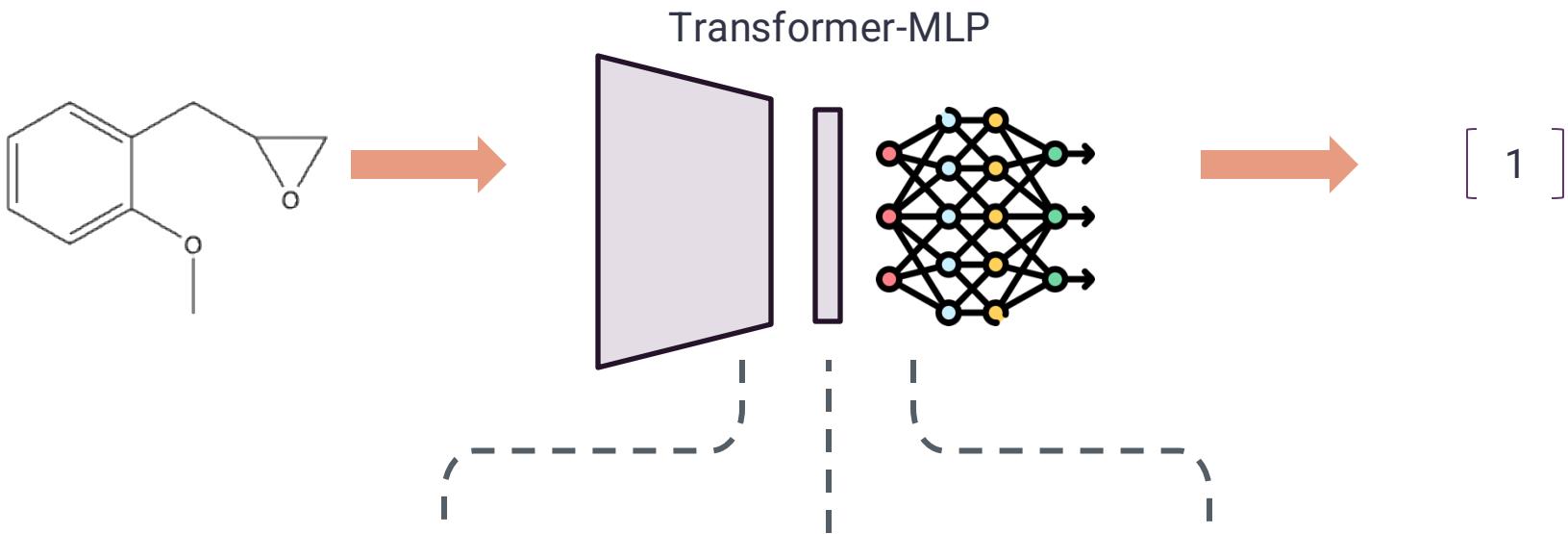
# Explainable AI for functional groups

DC2: Khasanova Dina

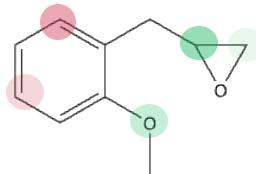


The reliability of XAI is important for eventually converting learned rules into structural alerts, CSRML, or SMARTS.

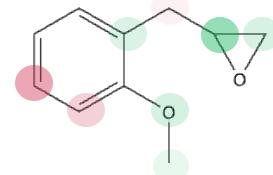




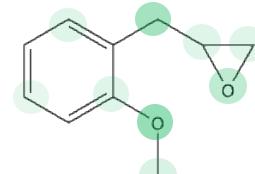
Integrated Gradients



SHAP

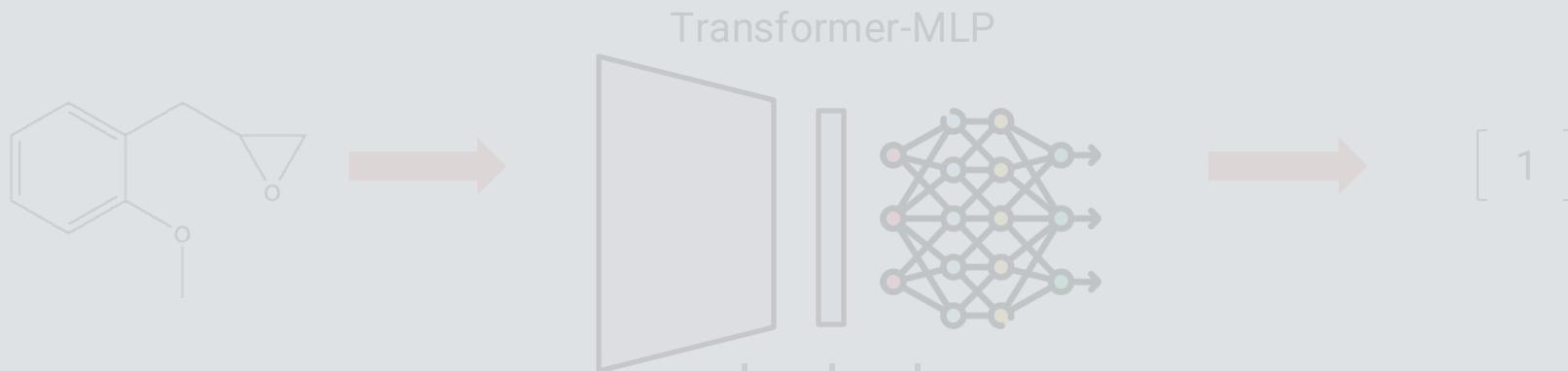


Attention Maps



Hartog, P.B.R., Krüger, F., Genheden, S. et al. Using test-time augmentation to investigate explainable AI: inconsistencies between method, model and human intuition.

J Cheminform 16, 39 (2024). <https://doi.org/10.1186/s13321-024-00824-1>



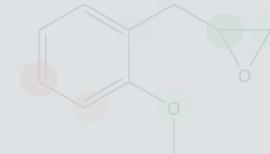
# Are XAI methods inconsistent?



Integrated Gradients



SHAP



Attention Maps



Hartog, P.B.R., Krüger, F., Genheden, S. et al. Using test-time augmentation to investigate explainable AI: inconsistencies between method, model and human intuition.

J Cheminform 16, 39 (2024). <https://doi.org/10.1186/s13321-024-00824-1>

**Table 7** AUROC, accuracy, F1, MCC precision and recall scores of MLP models transfer learned on Ames data

	Training	AUROC↑	Accuracy↑	F1↑	MCC↑	Precision↑	Recall↑
No training	Untrained	0.652	0.516	0.143	0.063	0.081	0.619
	Native	0.676	0.634	0.624	0.269	0.607	0.642
Variations	Random split	0.856	0.788	0.788	0.576	0.789	0.787
	Train set	0.873	0.792	0.792	0.584	0.792	0.792
Encoder only	CNN	0.709	0.650	0.657	0.300	0.671	0.644
	Enumerated	0.810	0.739	0.739	0.478	0.738	0.739
Encoder-decoder	C2C	0.734	0.666	0.666	0.332	0.665	0.666
	R2C	0.738	0.670	0.670	0.339	0.670	0.670
	E2C	0.731	0.665	0.665	0.331	0.665	0.665
	MC2C	0.754	0.682	0.682	0.364	0.683	0.682
	MR2C	0.694	0.653	0.652	0.305	0.651	0.653
	ME2C	0.804	0.719	0.719	0.438	0.719	0.719
	C2C	0.716	0.662	0.662	0.324	0.663	0.662
	R2C	0.698	0.634	0.634	0.269	0.634	0.634
	E2C	0.748	0.677	0.678	0.354	0.679	0.676
	MC2C	0.751	0.682	0.682	0.365	0.681	0.683
	MR2C	0.698	0.647	0.647	0.294	0.647	0.647
	ME2C	0.772	0.696	0.697	0.393	0.699	0.695

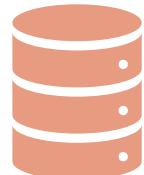
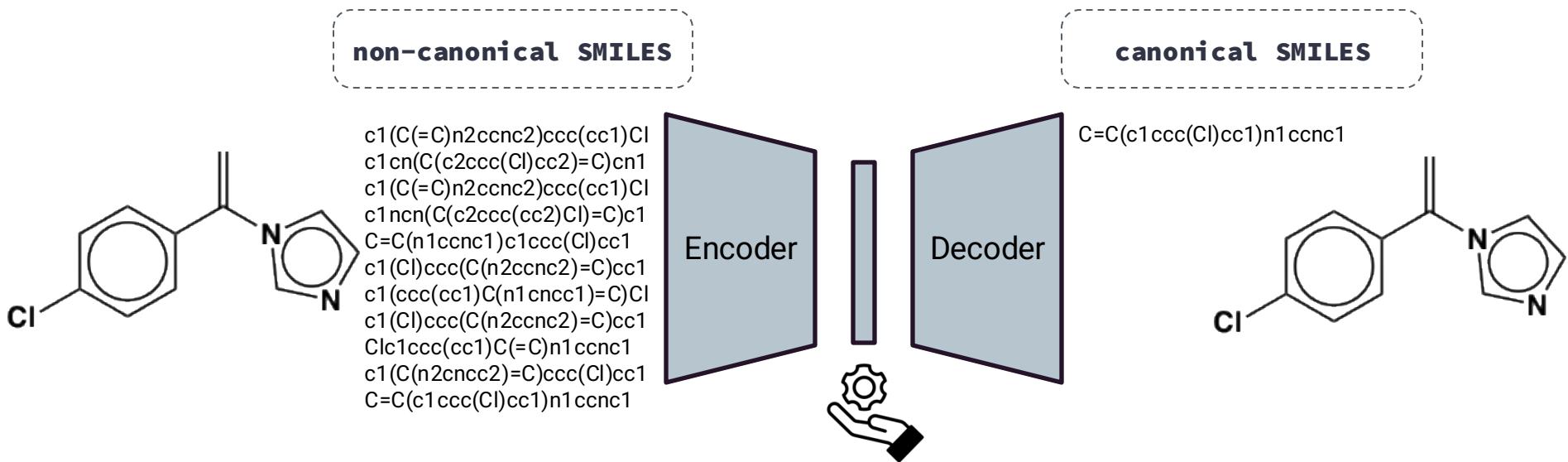
When the model's accuracy is low, each XAI method attempts to identify which features contributed to the incorrect predictions, leading to different features being highlighted and resulting in high cosine distances.

# Objective of the Project:

- To build a high-accuracy model for deterministic labels like **functional groups**, which are well-defined and universally agreed upon, to assess how accurate or inaccurate existing explainability methods are.

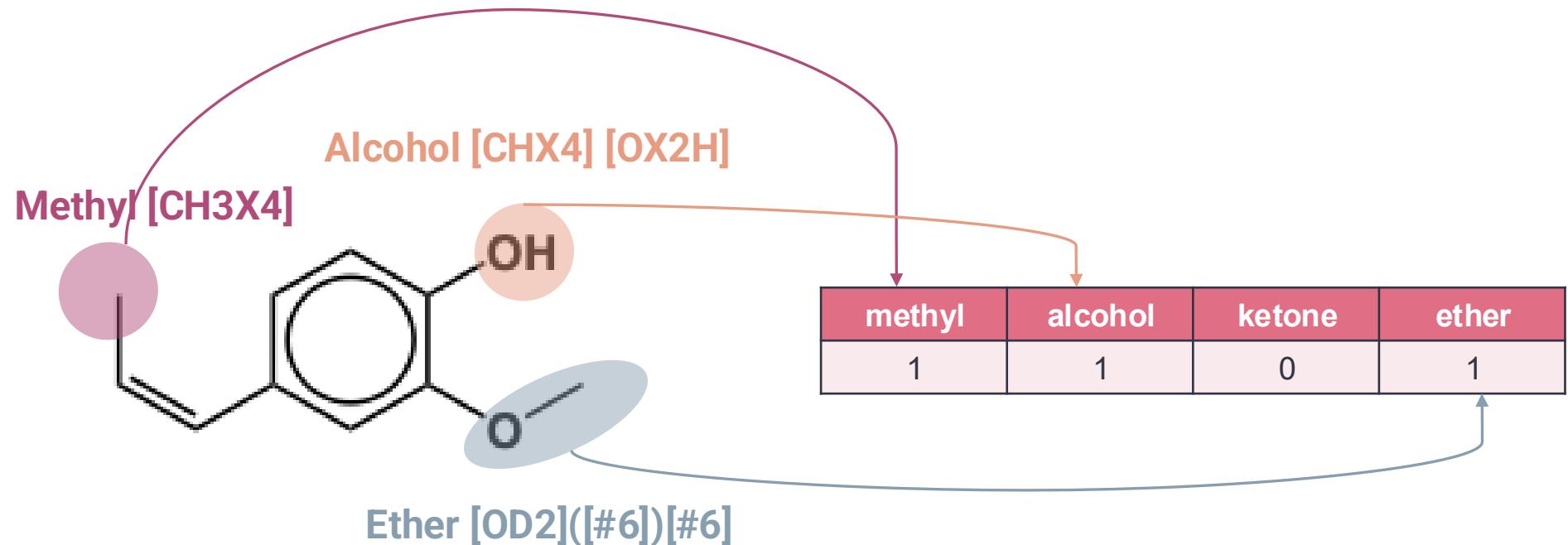


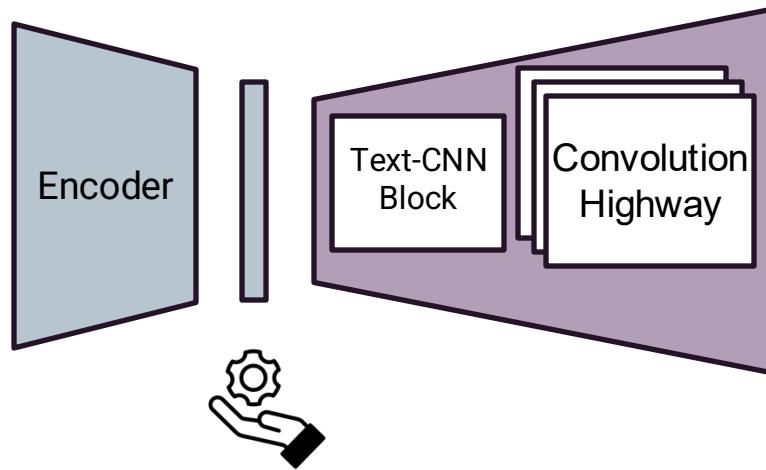
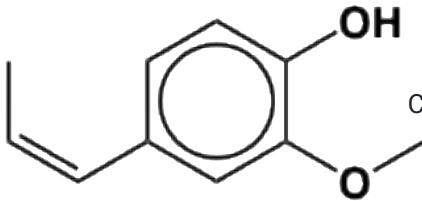
Objective



ChEMBL dataset  
(≈2M molecules)

Models





### Functional Groups

Alcohol	...	Methyl
1	...	0
...	...	...
1	...	1



TOX21 dataset  
(≈9K molecules)

Models

# XAI methods:

## Gradient-based

Integrated Gradients



SHAP GradientExplainer



## Backpropagation-Based

SHAP DeepExplainer



## Perturbation-based

SHAP KernelExplainer

LIME

Occlusion



## Attention-based

Grad-CAM

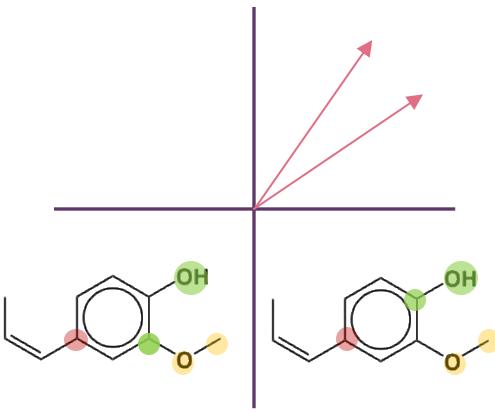


Attention Maps

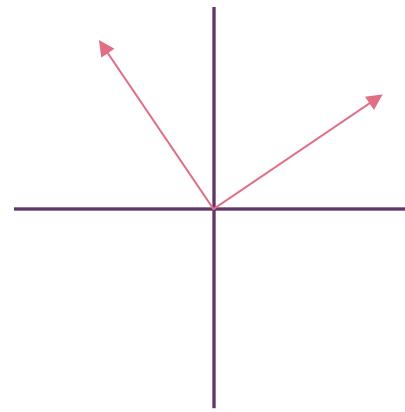
XAI

# Cosine similarity

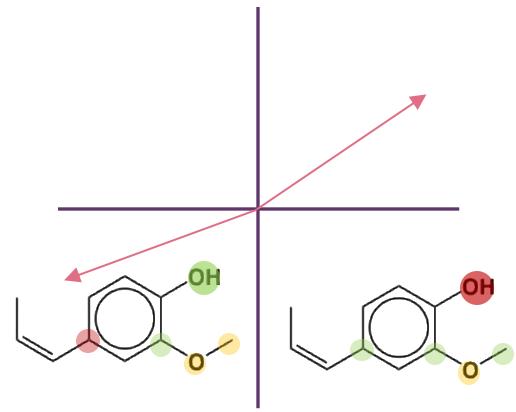
Similar



Unrelated



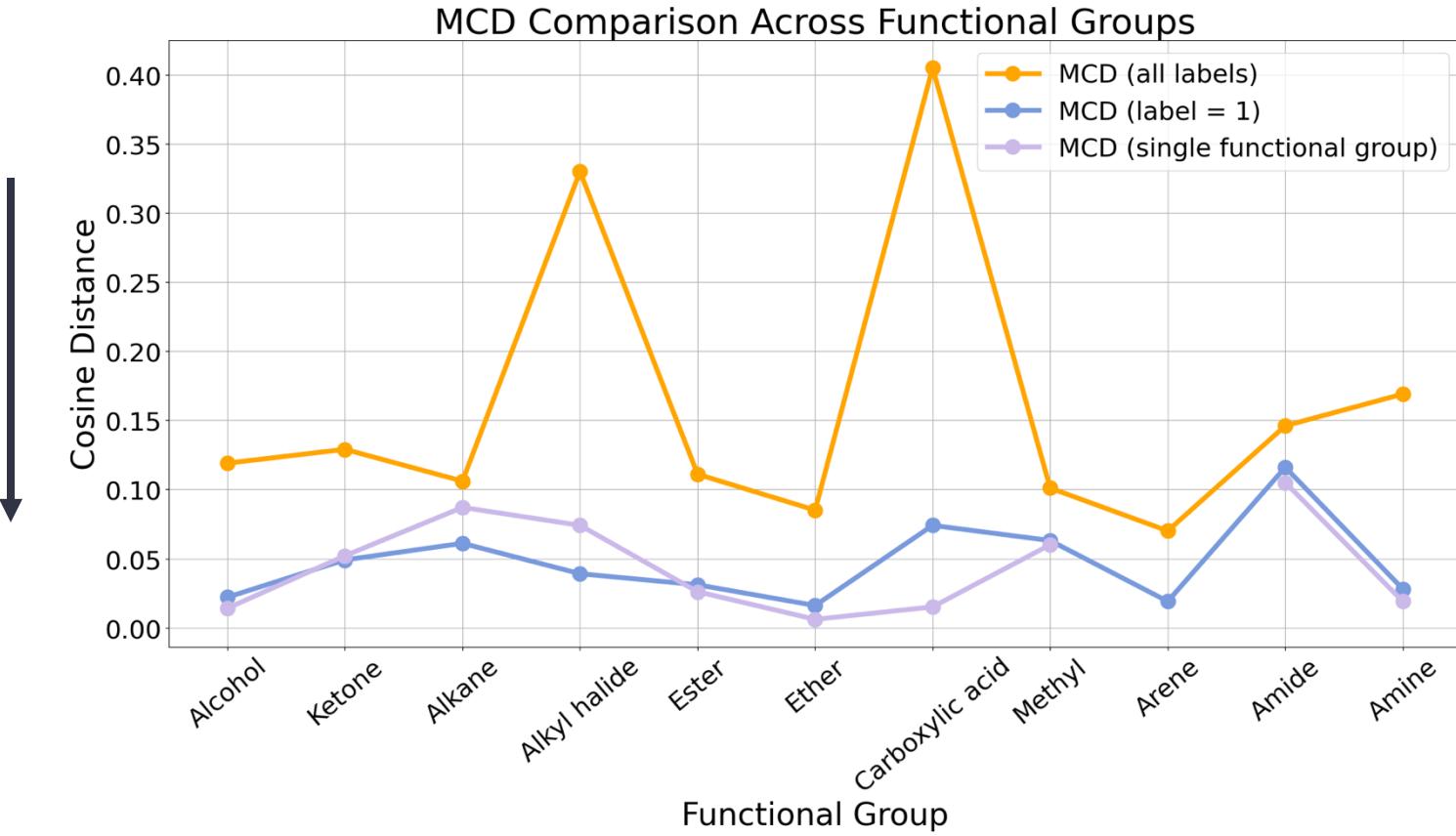
Opposite



$$\mathbf{A}^*\mathbf{B} = \|\mathbf{A}\| \|\mathbf{B}\| \cos\theta$$

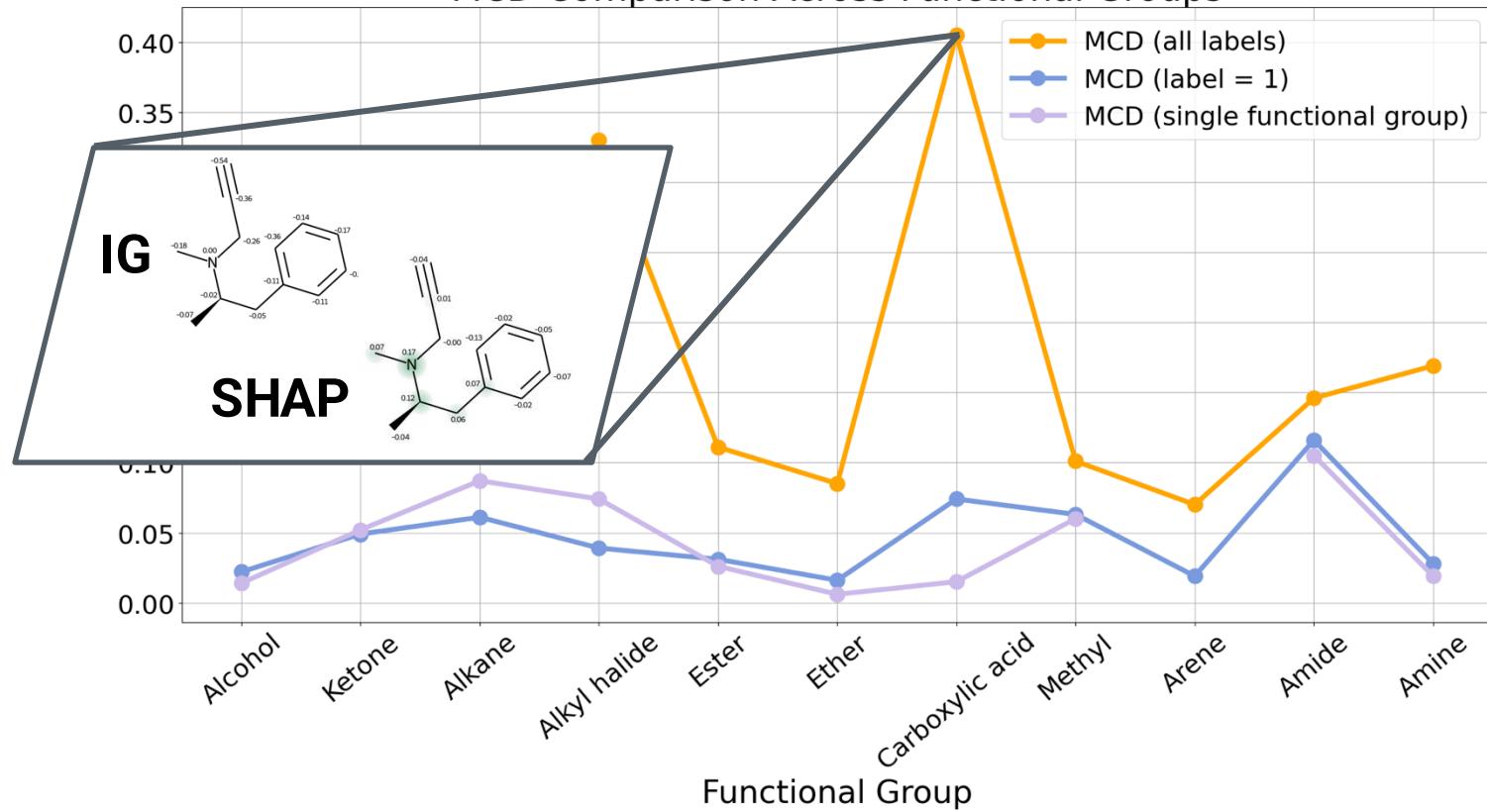
Cosine Distance = 1 - Cosine Similarity

XAI



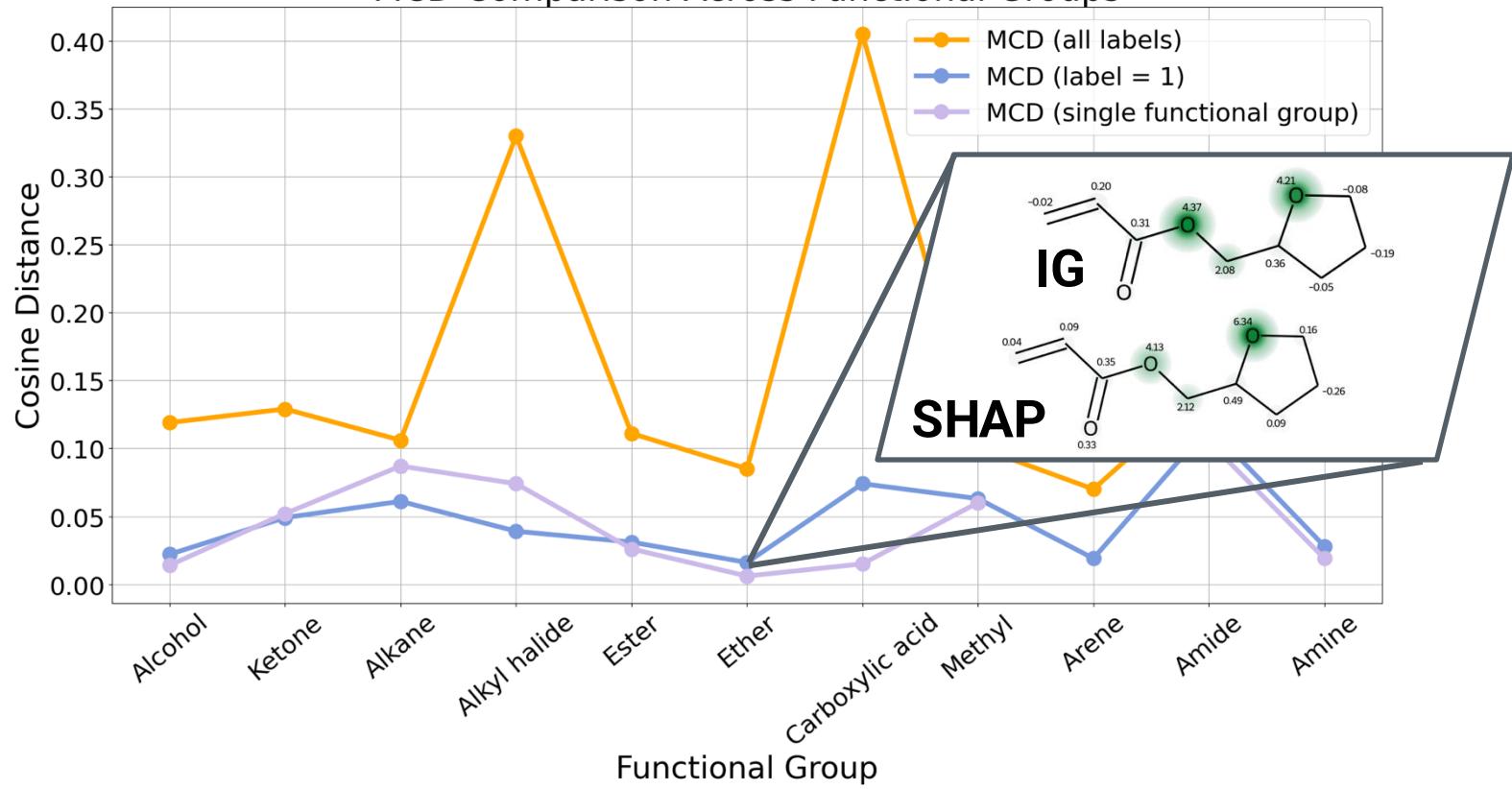
Results

## MCD Comparison Across Functional Groups



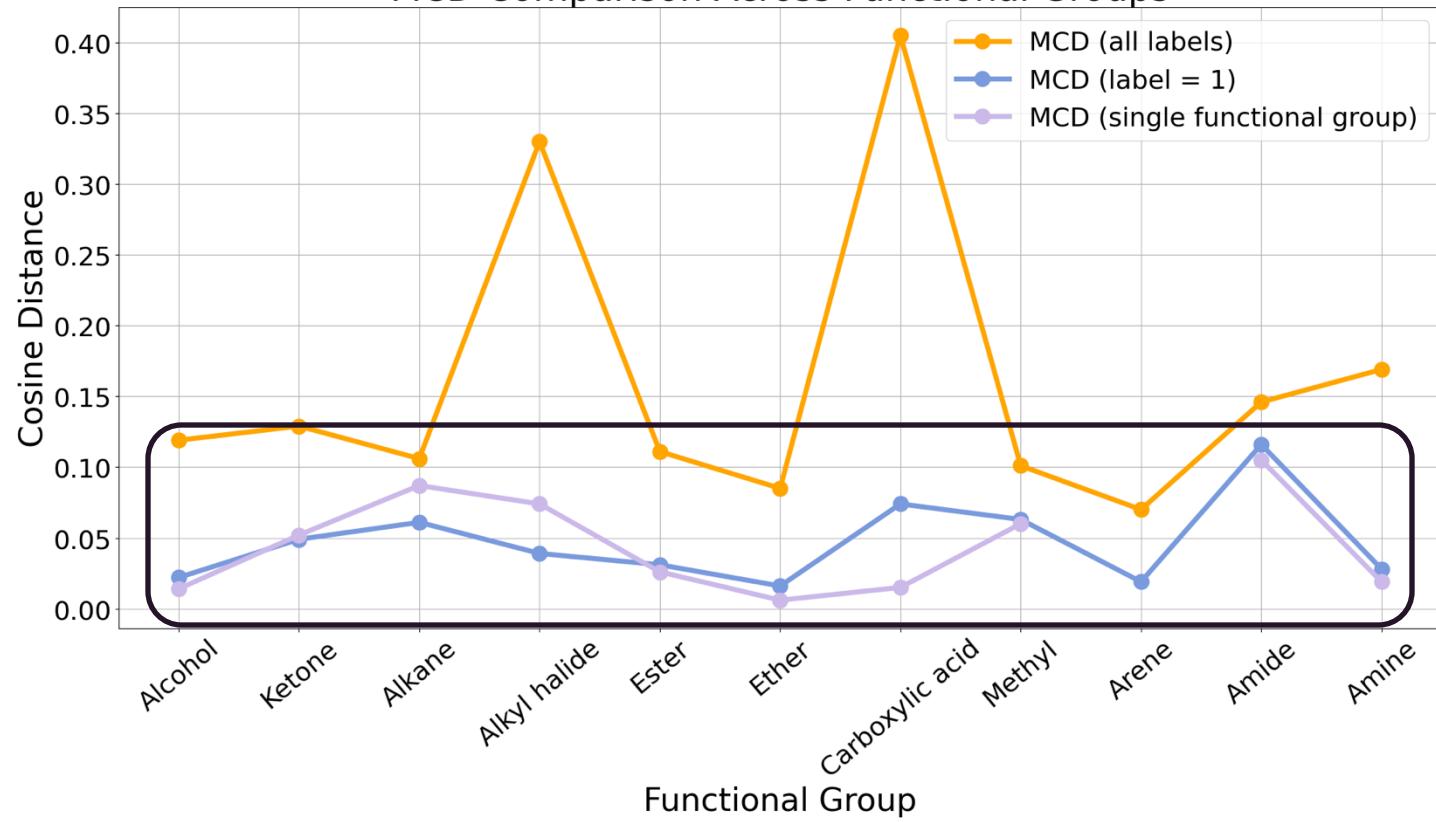
Results

## MCD Comparison Across Functional Groups

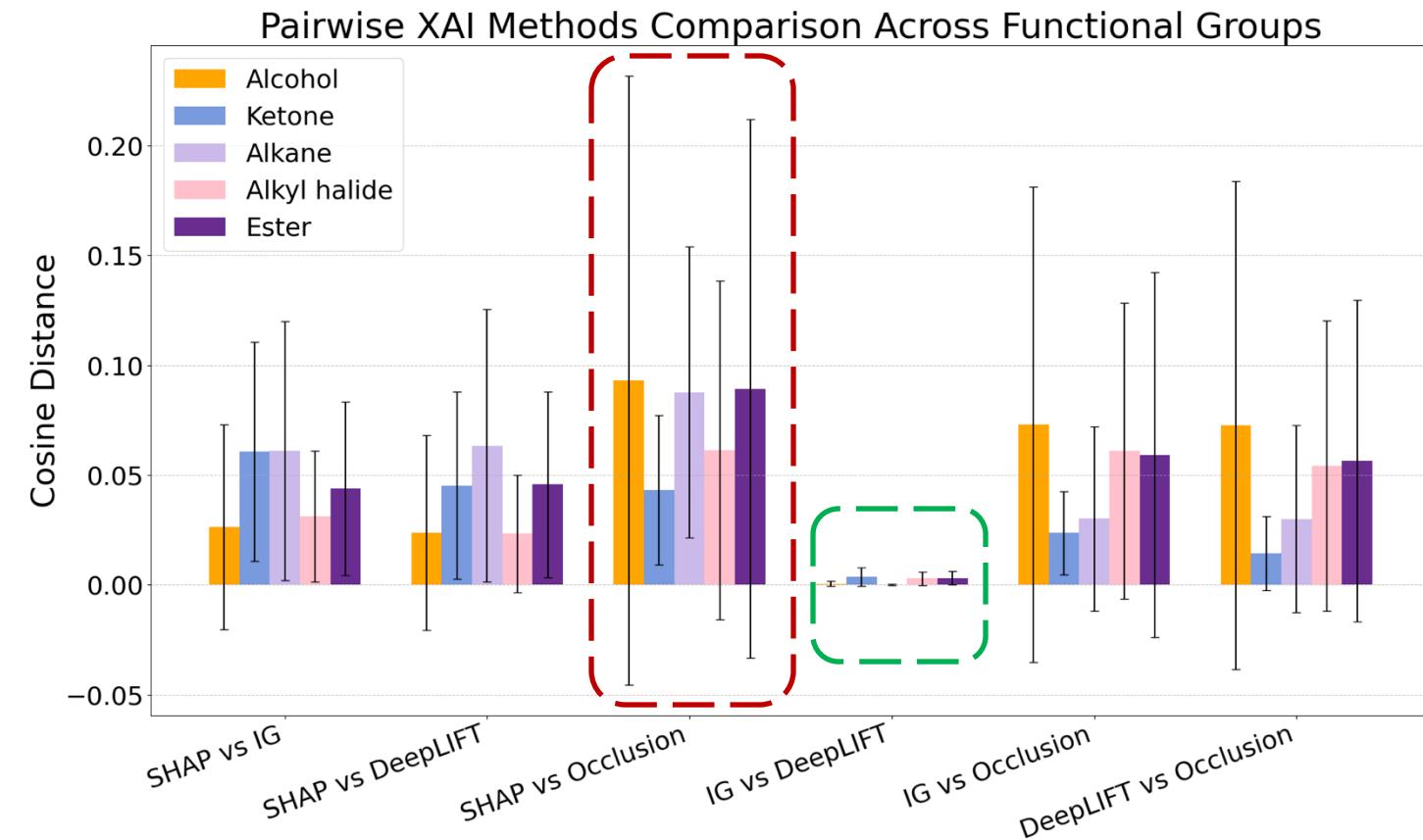


Results

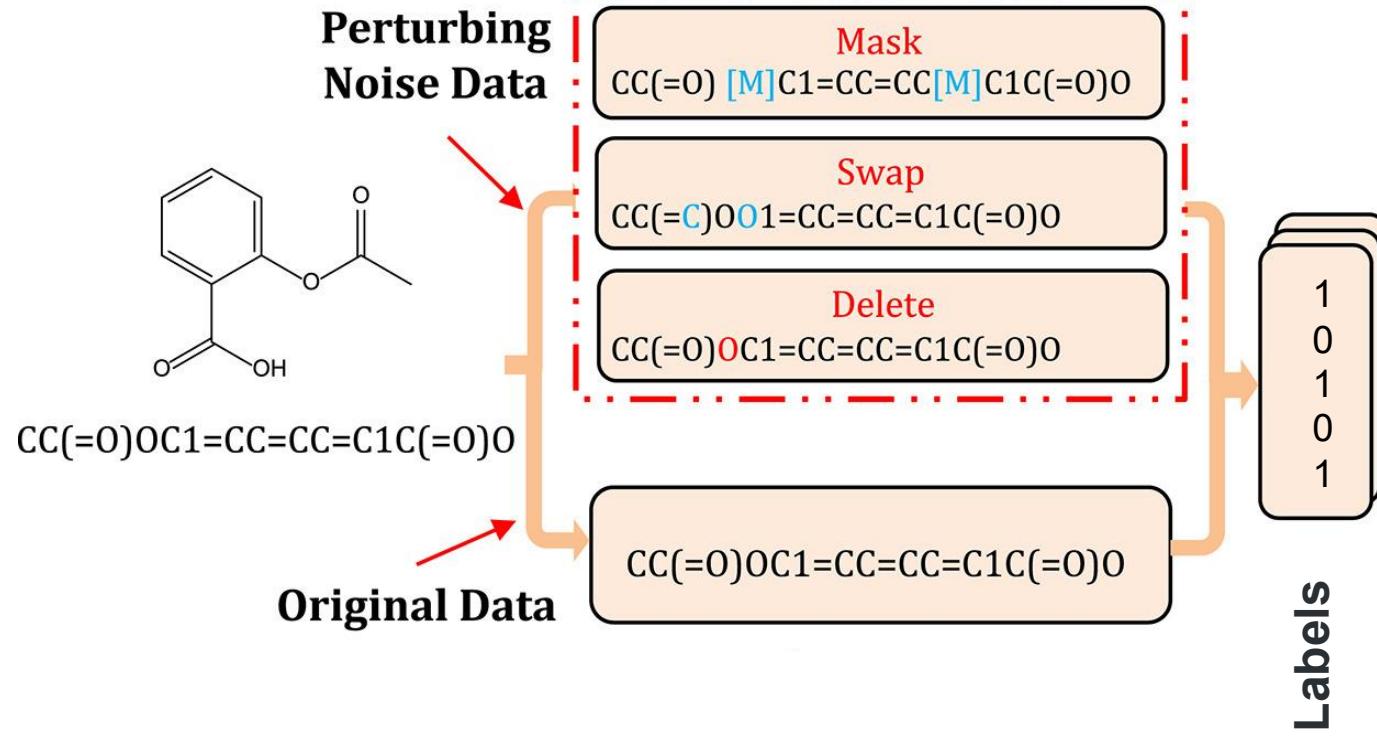
## MCD Comparison Across Functional Groups



Results

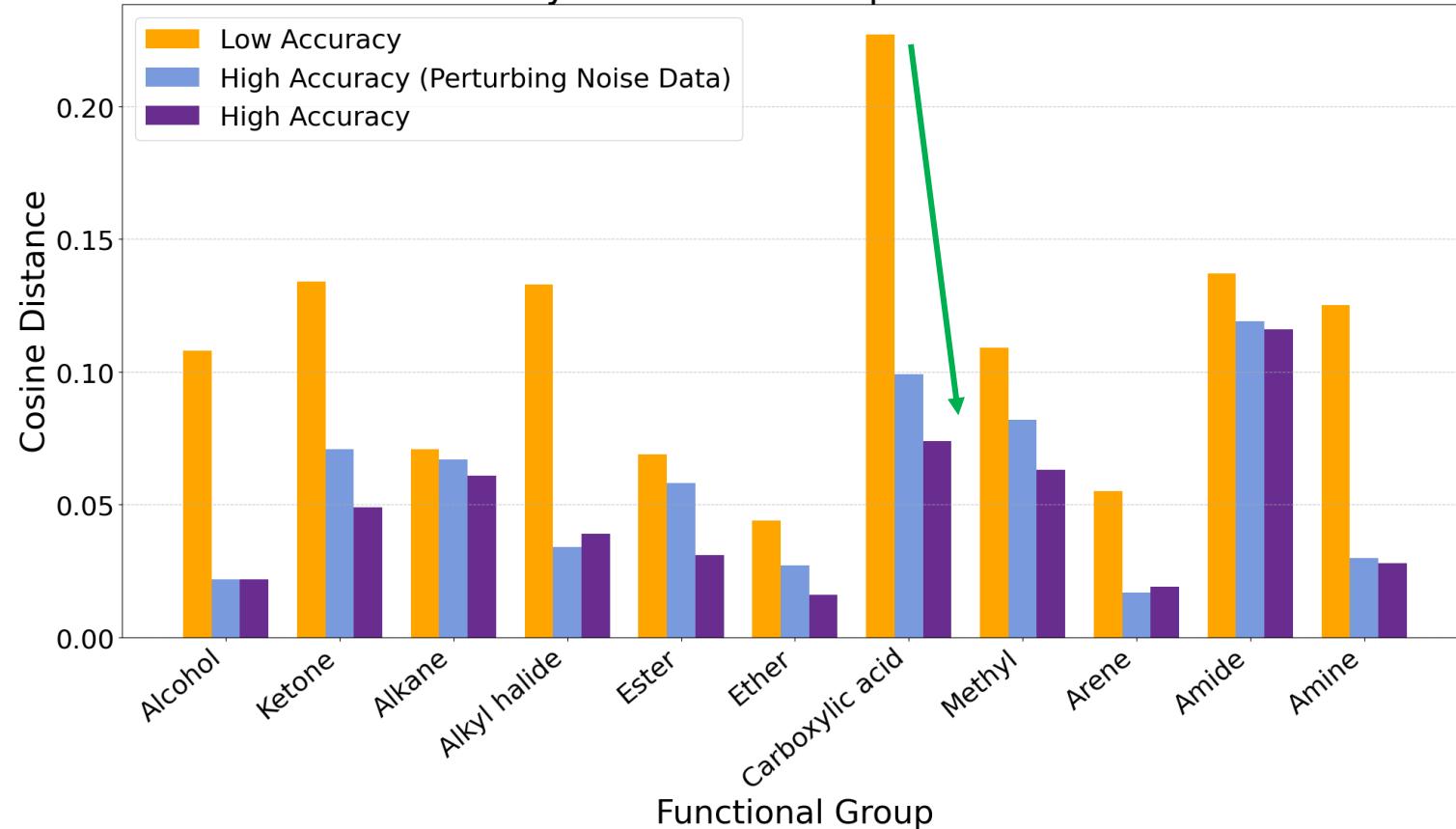


Results



Results

## MCD by Functional Group Across Models



Results

**Thank you  
for your attention!**