

# Of graphs, sets, and molecules

**AI4S3M2**

2025-11-26

Daniel Probst, Wageningen University & Research

# It's always graphs

- SMILES (or SELFIES, DeepSMILES, ...) are a way to represent (and store) graphs as strings
- ... but then again, “Transformers are Graph Neural Networks” [1]
- The baseline for GNNs applied to molecular property prediction is quite strong (when also adding RDKit descriptors)

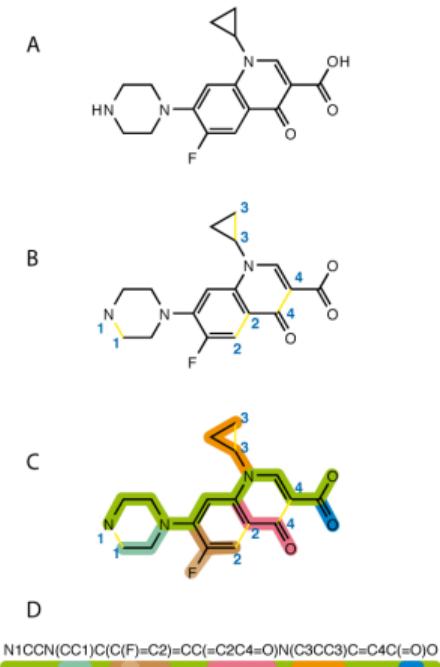


Figure 1: SMILES construction  
(Source: Wikipedia)

# What about sets?

- ... so we got (almost) rid of the bonds and thought of molecules as bags of atoms [2]
- This worked very well on the MoleculeNet benchmark... but everything does perform well on this

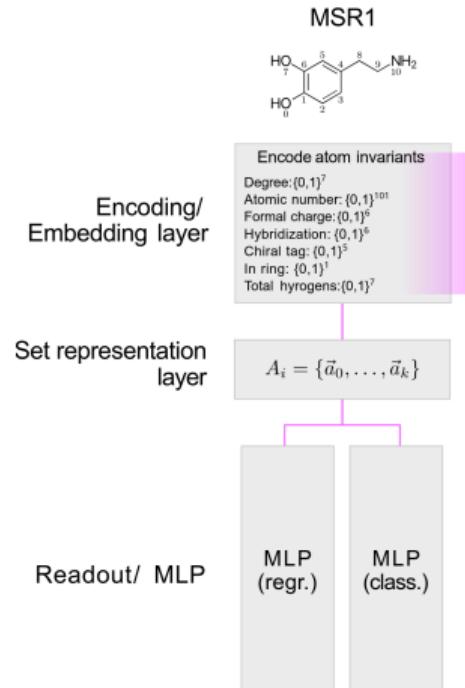


Figure 2: Molecular set representation architecture

# What about sets?

4 / 20

Data Set	Metric	MGCN	SchNet	GCN	GIN	D-MPNN	MSR1
HIV	AUROC	$73.8 \pm 1.6$	$70.2 \pm 3.4$	$71.6 \pm 4.0$	$75.3 \pm 1.9$	<u><math>75.0 \pm 2.1</math></u>	$72.3 \pm 2.2$
BACE	AUROC	$73.4 \pm 3.0$	$76.6 \pm 1.1$	$71.6 \pm 2.0$	$70.1 \pm 5.4$	<u><math>85.3 \pm 5.3</math></u>	$75.5 \pm 1.9$
BBBP	AUROC	<b><math>85.0 \pm 6.4</math></b>	$84.8 \pm 2.2$	$71.8 \pm 0.9$	$65.6 \pm 0.5$	$71.2 \pm 3.8$	<u><math>71.4 \pm 0.0</math></u>
Tox21	AUROC	$70.7 \pm 1.6$	<b><math>77.2 \pm 2.3</math></b>	$70.9 \pm 2.6$	$74.0 \pm 0.8$	$68.9 \pm 1.3$	<u><math>72.1 \pm 5.0</math></u>
SIDER	AUROC	$55.2 \pm 1.8$	$53.9 \pm 3.7$	$53.6 \pm 3.2$	$57.3 \pm 1.6$	<u><math>63.2 \pm 2.3</math></u>	$61.4 \pm 7.3$
ClinTox	AUROC	$63.4 \pm 4.2$	$71.5 \pm 3.7$	$62.5 \pm 2.8$	$58.0 \pm 4.4$	<u><math>90.5 \pm 5.3</math></u>	$86.6 \pm 1.2$
ESOL	RMSE	$1.27 \pm 0.15$	$1.05 \pm 0.06$	$1.43 \pm 0.05$	$1.45 \pm 0.02$	$0.98 \pm 0.26$	<u><math>0.59 \pm 0.03</math></u>
FreeSolv	RMSE	$3.35 \pm 0.01$	$3.22 \pm 0.76$	$2.87 \pm 0.14$	$2.76 \pm 0.18$	$2.18 \pm 0.91$	<u><math>1.94 \pm 0.24</math></u>
Lipo	RMSE	$1.11 \pm 0.04$	$0.91 \pm 0.10$	$0.85 \pm 0.08$	$0.85 \pm 0.07$	<u><math>0.65 \pm 0.05</math></u>	$0.85 \pm 0.03$
QM7	MAE	$77.6 \pm 4.7$	<b><math>74.2 \pm 6.0</math></b>	$122.9 \pm 2.2$	$124.8 \pm 0.7$	$105.8 \pm 13.2$	<u><math>85.9 \pm 13.1</math></u>
QM8	MAE	$0.022 \pm 0.002$	$0.020 \pm 0.002$	$0.037 \pm 0.001$	$0.037 \pm 0.001$	<b><u><math>0.014 \pm 0.002</math></u></b>	$0.023 \pm 0.010$

# What about sets of node embeddings?

- What about combining the set representation approach with GNNs?
- Embedding nodes, followed by a set representation layer instead of global pooling
- Turns out, this works quite well on more representative benchmarks

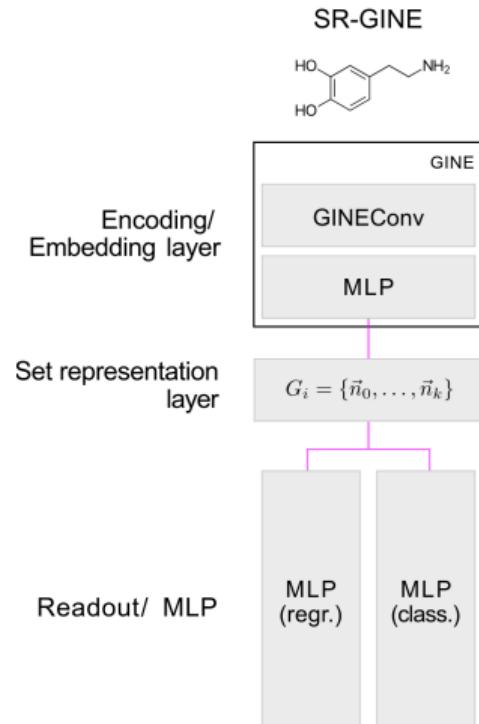


Figure 3: SR-GINE architecture

# What about sets of node embeddings?

6 / 20

Data from Biogen ADME in vitro assays. The metric used is the Pearson correlation. Similar results on OCELOT chromophore data containing chemically diverse  $\pi$ -conjugated molecules [3], [4].

Endpoint	D-MPNN	D-MPNN+	MSR1	GINE	SR-GINE
HLM	<u>0.65</u>	<b>0.68</b>	$0.58 \pm 0.00$	$0.59 \pm 0.01$	<b><math>0.68 \pm 0.01</math></b>
MDR1-MDCK ER	<u>0.72</u>	0.78	$0.68 \pm 0.02$	$0.67 \pm 0.00$	<b><math>0.80 \pm 0.01</math></b>
Solubility	<b>0.64</b>	0.59	$0.47 \pm 0.04$	$0.55 \pm 0.02$	<u><math>0.63 \pm 0.02</math></u>
RLM	<u>0.72</u>	<b>0.74</b>	$0.57 \pm 0.01$	$0.61 \pm 0.02$	<u><math>0.70 \pm 0.01</math></u>
hPPB	<u>0.74</u>	<b>0.77</b>	$0.68 \pm 0.01$	$0.68 \pm 0.02$	<b><math>0.77 \pm 0.01</math></b>
rPPB	<u>0.70</u>	0.70	$0.55 \pm 0.02$	$0.63 \pm 0.02$	<b><math>0.73 \pm 0.02</math></b>
Average	$0.695 \pm 0.037$	<u><math>0.710 \pm 0.064</math></u>	$0.588 \pm 0.08$	$0.622 \pm 0.05$	<b><math>0.718 \pm 0.06</math></b>

# Sets of molecules (reactions)

7 / 20

Yield prediction on Buchwald-Hartwig reactions from high-throughput experiments (HTE) and extracted from an electronic laboratory notebook (ELN) [5], [6]. The metric used is  $R^2$ .

Data Set	Subset/Split	Yield-BERT	DRFP	YieldGNN	MSR2-RXN
HTE	Rand 70/30	$0.95 \pm 0.005$	$0.95 \pm 0.005$	$0.96 \pm 0.005$	$0.94 \pm 0.005$
	Rand 50/50	$0.92 \pm 0.01$	$0.93 \pm 0.01$	-	$0.93 \pm 0.01$
	Rand 30/70	$0.88 \pm 0.01$	$0.89 \pm 0.01$	-	$0.90 \pm 0.01$
	Rand 20/80	$0.86 \pm 0.01$	$0.87 \pm 0.01$	-	$0.87 \pm 0.01$
	Rand 10/90	$0.79 \pm 0.02$	<b><math>0.81 \pm 0.01</math></b>	-	$0.80 \pm 0.02$
	Rand 5/95	$0.61 \pm 0.04$	$0.73 \pm 0.02$	-	$0.69 \pm 0.03$
	Rand 2.5/97.5	$0.45 \pm 0.05$	<b><math>0.62 \pm 0.04</math></b>	-	$0.57 \pm 0.05$
	Avg. 1-4	<b>0.73</b>	$0.71 \pm 0.16$	-	$0.72 \pm 0.15$
ELN	-	$-0.006 \pm 0.11$	<b><math>0.20 \pm 0.05</math></b>	$0.05 \pm 0.07$	$0.13 \pm 0.08$

# Speaking of sets and reactions ...

- We recently published an update on the differential reaction fingerprint (DRFP) where we set SOTA on the ELN data set
- See results on previous slide
- More on sets, molecules, and reactions coming soon ...

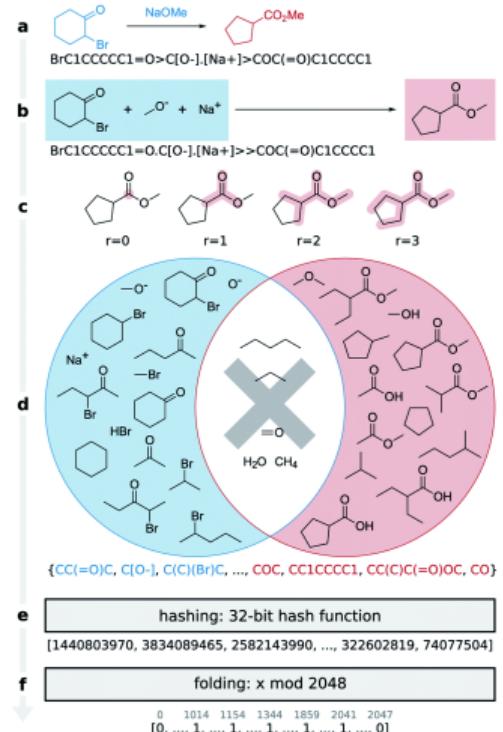


Figure 4: DRFP Scheme

# Speaking of sets and reactions ...

9 / 20

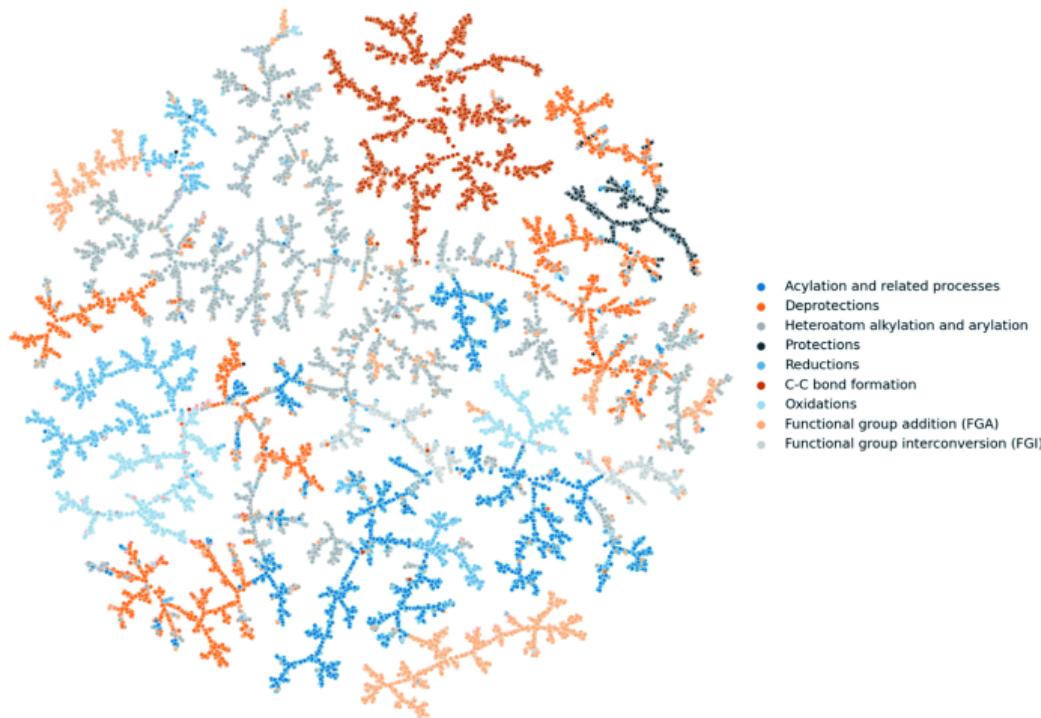


Figure 5: DRFP-encoded USPTO subset visualised using TMAP.

- Mass spectra are commonly binned to account for different numbers of peaks between samples
- Why not represent them as graphs and predict the quantitative estimate of drug-likeness (QED)?
- Intensities as node attributes and m/z differences as edge weights [7]

Model	Params	MAE (↓)	RMSE (↓)	Pearson's $r$ (↑)	$R^2$ (↑)
MLP	11.0 M	$0.145 \pm 0.008$	$0.200 \pm 0.008$	$0.736 \pm 0.011$	$0.437 \pm 0.043$
SetTransformer	0.4 M	$0.134 \pm 0.002$	$0.174 \pm 0.001$	$0.758 \pm 0.004$	$0.572 \pm 0.006$
GNN (GAT)	12.6 M	<b><math>0.110 \pm 0.006</math></b>	<b><math>0.144 \pm 0.006</math></b>	<b><math>0.843 \pm 0.015</math></b>	<b><math>0.709 \pm 0.025</math></b>

# Molecular dynamics as graphs

11 / 20

- Idea: Combine structure and dynamics in a heterogeneous graph [8]
- 8 ns MD per protein (or complex) from the MISATO data set [9]
- Significant performance improvements in binding site detection, binding affinity prediction, and atomic adaptability prediction

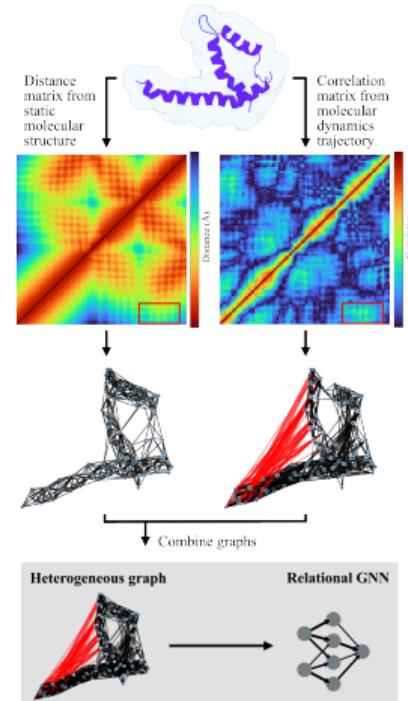


Figure 6: DRFP Scheme

# Molecular dynamics as graphs

- Hypothesis: In addition to the dynamics information, dynamics edges reduce over-squashing and improve long-range information propagation
- More recent research showed that random edges do not have a similar effect on performance

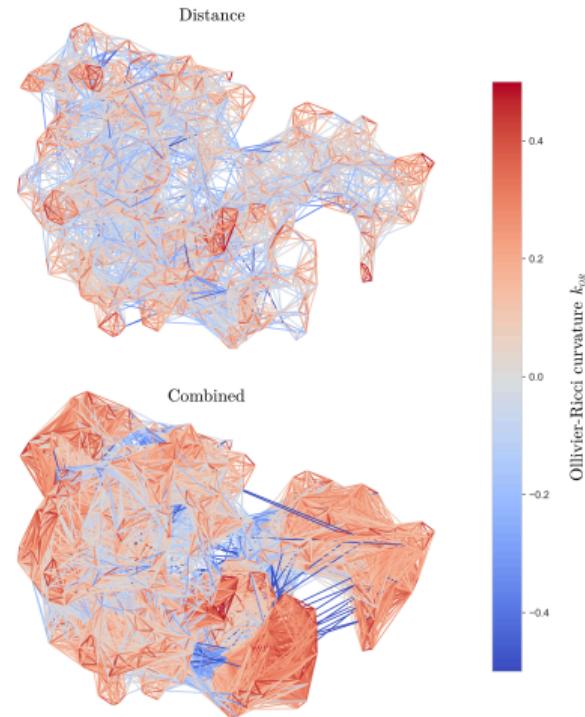


Figure 7: Ollivier-Ricci curvature

# Modelling protein conformations

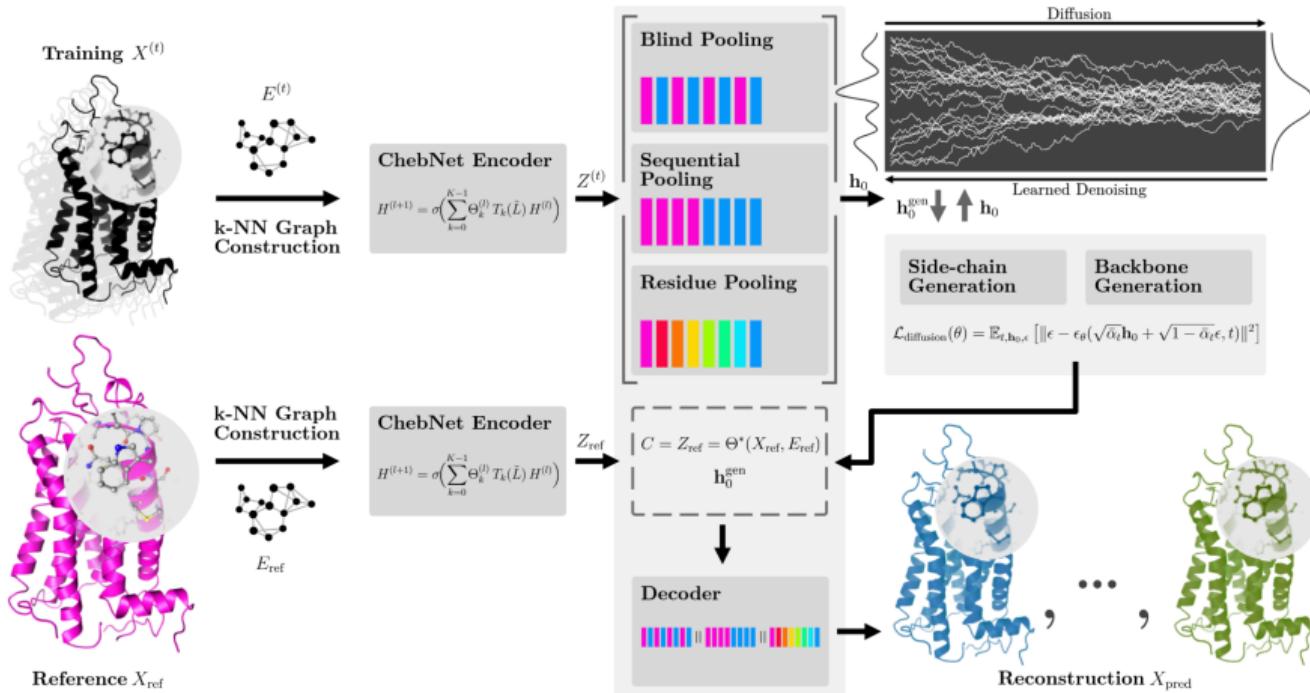


Figure 8: Schematic of the LD-FPG framework.

# Modelling protein conformations

14 / 20

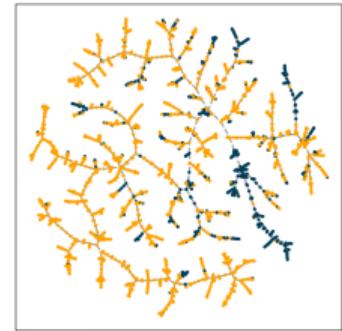
- Autoencoder architecture utilising **ChebNet** combined with residue-based pooling [10]
- As it turns out, ChebNet performs well compared to other GNN architectures on long-range tasks
- But you have probably heard enough about protein ensemble generation already ...

# Back to SMILES and ... GZIP

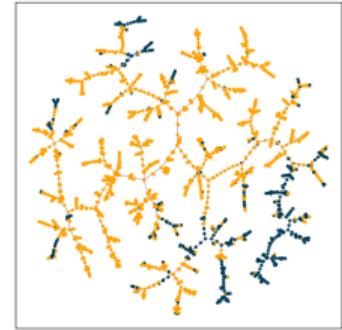
15 / 20

Model		RMSE	MAE	R
GraphDTA	GCN	$1.735 \pm 0.034$	$1.343 \pm 0.037$	$0.613 \pm 0.016$
	GAT	$1.765 \pm 0.026$	$1.354 \pm 0.033$	$0.601 \pm 0.016$
	GIN	$1.640 \pm 0.044$	$1.261 \pm 0.044$	$0.667 \pm 0.018$
	GAT-GCN	$1.562 \pm 0.022$	$1.191 \pm 0.016$	$0.697 \pm 0.008$
GNN-based	GNN-DTI	$1.492 \pm 0.025$	$1.192 \pm 0.032$	$0.736 \pm 0.021$
	D-MPNN	$1.493 \pm 0.016$	$1.188 \pm 0.009$	$0.729 \pm 0.006$
	MAT	$1.457 \pm 0.037$	$1.154 \pm 0.037$	$0.747 \pm 0.013$
	DimeNet	$1.453 \pm 0.027$	$1.138 \pm 0.026$	$0.752 \pm 0.010$
	CMPNN	<b><math>1.408 \pm 0.028</math></b>	<b><math>1.117 \pm 0.031</math></b>	<b><math>0.765 \pm 0.009</math></b>
Compression-based	MolZip	$1.508 \pm 0.000$	$1.190 \pm 0.000$	$0.720 \pm 0.000$
	MolZip Aug	<u><math>1.422 \pm 0.017</math></u>	<u><math>1.131 \pm 0.014</math></u>	<u><math>0.757 \pm 0.007</math></u>

- Based on the normalised compression distance
- $\text{NCD}(x, y) = \frac{0.5(C(xy)+C(yx)) - \min\{C(xx), C(yy)\}}{\max\{C(xx), C(yy)\}}$
- approximates normalised information distance, specifically the Kolmogorov complexity, which is uncomputable ( $\rightarrow$  halting problem) in
- $\text{NID}(x, y) = \frac{\max\{K(x|y), K(y|x)\}}{\max\{K(x), K(y)\}}$
- Based on this, we build kNN classifiers and regressors [11]



(a) ECFP-based TMAP



(b) MolZip-based TMAP

Figure 9: TMAPs of BBBP

- [daniel.probst@wur.nl](mailto:daniel.probst@wur.nl)
- [probstlab.science](http://probstlab.science)
- [github.com/daenuprobst/talks](https://github.com/daenuprobst/talks)

- [1] C. K. Joshi, “Transformers are Graph Neural Networks.” [Online]. Available: <https://arxiv.org/abs/2506.22084>
- [2] M. Boulougouri, P. Vandergheynst, and D. Probst, “Molecular set representation learning,” *Nature Machine Intelligence*, vol. 6, no. 7, pp. 754–763, July 2024, doi: [10.1038/s42256-024-00856-0](https://doi.org/10.1038/s42256-024-00856-0).
- [3] C. Fang *et al.*, “Prospective Validation of Machine Learning Algorithms for Absorption, Distribution, Metabolism, and Excretion Prediction: An Industrial Perspective,” *Journal of Chemical Information and Modeling*, vol. 63, no. 11, pp. 3263–3274, June 2023, doi: [10.1021/acs.jcim.3c00160](https://doi.org/10.1021/acs.jcim.3c00160).
- [4] V. Bhat, P. Sornberger, B. S. S. Pokuri, R. Duke, B. Ganapathysubramanian, and C. Risko, “Electronic, redox, and optical property prediction of organic *pi*-conjugated

## References (ii)

19 / 20

- molecules through a hierarchy of machine learning approaches”, *Chemical Science*, vol. 14, no. 1, pp. 203–213, Dec. 2022, doi: [10.1039/D2SC04676H](https://doi.org/10.1039/D2SC04676H).
- [5] D. T. Ahneman, J. G. Estrada, S. Lin, S. D. Dreher, and A. G. Doyle, “Predicting reaction performance in C–N cross-coupling using machine learning,” *Science*, vol. 360, no. 6385, pp. 186–190, Apr. 2018, doi: [10.1126/science.aar5169](https://doi.org/10.1126/science.aar5169).
- [6] M. Saebi *et al.*, “On the use of real-world datasets for reaction yield prediction,” *Chemical Science*, vol. 14, no. 19, pp. 4997–5005, 2023, doi: [10.1039/D2SC06041H](https://doi.org/10.1039/D2SC06041H).
- [7] N. de Jonge, J. J. J. van der Hooft, and D. Probst, “To Bin or not to Bin: Alternative Representations of Mass Spectra.” [Online]. Available: <https://arxiv.org/abs/2502.10851>
- [8] P. Guo, B. Correia, P. Vandergheynst, and D. Probst, “Boosting Protein Graph Representations through Static-Dynamic Fusion,” in *Proceedings of the 42nd International Conference on Machine Learning*, A. Singh, M. Fazel, D. Hsu, S. Lacoste-

## References (iii)

20 / 20

- Julien, F. Berkenkamp, T. Maharaj, K. Wagstaff, and J. Zhu, Eds., in Proceedings of Machine Learning Research, vol. 267. PMLR, 2025, pp. 20777–20792.
- [9] T. Siebenmorgen *et al.*, “MISATO: machine learning dataset of protein–ligand complexes for structure-based drug discovery,” *Nature Computational Science*, vol. 4, no. 5, pp. 367–378, May 2024, doi: [10.1038/s43588-024-00627-2](https://doi.org/10.1038/s43588-024-00627-2).
- [10] A. Sengar, A. Hariri, D. Probst, P. Barth, and P. Vandergheynst, “Generative Modeling of Full-Atom Protein Conformations using Latent Diffusion on Graph Embeddings.” [Online]. Available: <https://arxiv.org/abs/2506.17064>
- [11] J. Weinreich and D. Probst, “Learning on compressed molecular representations,” *Digital Discovery*, vol. 4, no. 1, pp. 84–92, 2025, doi: [10.1039/d4dd00162a](https://doi.org/10.1039/d4dd00162a).