
Context-enriched molecule representations improve few-shot drug discovery

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A central task in computational drug discovery is to construct models from known
2 active molecules to find further promising molecules for subsequent screening.
3 However, typically only very few active molecules are known. Therefore, few-shot
4 learning methods have the potential to improve the effectiveness of this critical
5 phase of the drug discovery process. We introduce a new method for few-shot
6 drug discovery. Its main idea is to enrich a molecule representation by knowledge
7 about known context or reference molecules. Our novel concept for molecule
8 representation enrichment is to associate molecules from both the support set and
9 the query set with a large set of reference (context) molecules through a modern
10 Hopfield network. Intuitively, this enrichment step is analogous to a human expert
11 who would associate a given molecule with familiar molecules whose properties
12 are known. The enrichment step reinforces and amplifies the covariance structure
13 of the data, while simultaneously removing spurious correlations arising from the
14 decoration of molecules. Our approach is compared with other few-shot methods
15 for drug discovery on the FS-Mol benchmark dataset. On FS-Mol, our approach
16 outperforms all compared methods and therefore sets a new state-of-the art for
17 few-shot learning in drug discovery. An ablation study shows that the enrichment
18 step of our method is the key to improve the predictive quality. In a domain shift
19 experiment, we further demonstrate the robustness of our method.

20 1 Introduction

21 To improve human health, combat diseases, and tackle pandemics there is a steady need of discovering
22 new drugs in a fast and efficient way. However, the drug discovery process is time-consuming and
23 cost-intensive (Arrowsmith, 2011). Deep learning methods have recently been shown to reduce
24 time and costs of this process (Chen et al., 2018; Walters and Barzilay, 2021). They diminish the
25 required number of both wet-lab measurements and molecules that must be synthesized (Merk et al.,
26 2018; Schneider et al., 2020). However, as of now, deep learning approaches use only the molecular
27 information about the ligands after being trained on a large training set. At inference time, they yield
28 highly accurate property and activity prediction (Mayr et al., 2018; Yang et al., 2019), generative
29 (Segler et al., 2018a; Gómez-Bombarelli et al., 2018), or synthesis models (Segler et al., 2018b; Seidl
30 et al., 2022).

31 **Deep learning methods in drug discovery usually require large amounts of biological measure-**
32 **ments.** To train deep learning-based activity and property prediction models with high predictive per-
33 formance, hundreds or thousands of data points per task are required. For example, well-performing
34 predictive models for activity prediction tasks of ChEMBL have been trained with an average of 3,621
35 activity points per task, i.e., drug target, by (Mayr et al., 2018). The ExCAPE-DB dataset provides on
36 average 42,501 measurements per task (Sun et al., 2017; Sturm et al., 2020). (Wu et al., 2018) pub-

37 lished a large scale benchmark for molecular machine learning, including prediction models for the
38 SIDER dataset (Kuhn et al., 2016) with an average of 5,187 data points, Tox21 (Huang et al., 2016b;
39 Mayr et al., 2016) with on average 9,031, and ClinTox (Wu et al., 2018) with 1,491 measurements
40 per task. However, for typical drug design projects, the amount of available measurements is very
41 limited (Stanley et al., 2021; Waring et al., 2015; Hochreiter et al., 2018), since in-vitro experiments
42 are expensive and time-consuming. Therefore, methods that need only few measurements to build
43 precise prediction models are desirable. This problem — i.e., the challenge of learning from few
44 data points — is the focus of machine learning areas like meta-learning (Schmidhuber, 1987; Bengio
45 et al., 1991; Hochreiter et al., 2001) and few-shot learning (Miller et al., 2000; Bendre et al., 2020;
46 Wang et al., 2020).

47 **Few-shot learning tackles the low-data problem that is ubiquitous in drug discovery.** Few-shot
48 learning methods have been predominantly developed and tested on image datasets (Bendre et al.,
49 2020; Wang et al., 2020), and have recently been adapted to drug discovery problems (Chen et al.,
50 2022; Wang et al., 2021; Stanley et al., 2021; Altae-Tran et al., 2017). They are usually categorized
51 into three groups according to their main approach (Bendre et al., 2020; Wang et al., 2020; Adler et al.,
52 2020). a) Data-augmentation-based approaches augment the available samples and generate new, more
53 diverse data points (Chen et al., 2020; Zhao et al., 2019; Antoniou and Storkey, 2019). b) Embedding-
54 based and nearest neighbour approaches learn embedding space representations. Predictive models
55 can then be constructed from only few net data points by comparing these embeddings. For example
56 in Matching Networks (Vinyals et al., 2016) an attention mechanism that relies on embeddings is the
57 basis for the predictions. Prototypical Networks (Snell et al., 2017) create prototype representations
58 for each class using the above mentioned representations in the embedding space. c) Optimization-
59 based or fine-tuning methods utilize a meta-optimizer that focuses on efficiently navigating the
60 parameter space. For example, with MAML the meta-optimizer learns initial weights that can be
61 adapted to a novel task by few optimization steps (Finn et al., 2017).

62 Most of these approaches have already been applied to few-shot drug discovery (see Sec. 4). Surpris-
63 ingly, almost all these few-shot learning methods in drug discovery are worse than a naive baseline,
64 which does not even use the support set (see Section 5). We hypothesize that the under-performance
65 of these methods stems from disregarding the context — both in terms of similar molecules and
66 similar activities. Therefore, we propose a method that informs the representations of the query and
67 support set with a large number of context molecules covering the chemical space.

68 **Enriching molecule representations with context using associative memories.** In data-scarce
69 situations, humans extract co-occurrences and covariances by associating current perceptions with
70 memories (Bonner and Epstein, 2021; Potter, 2012). When we show a small set of active molecules to
71 a human expert in drug discovery, the expert associates them with known molecules to suggest further
72 active molecules (Gomez, 2018; He et al., 2021). In an analogous manner, our novel concept for
73 few-shot learning uses associative memories to extract co-occurrences and the covariance structure
74 of the original data and to amplify them in the representations (Fürst et al., 2021). We use Modern
75 Hopfield Networks (MHNs) as an associative memory, since they can store a large set of context
76 molecule representations (Ramsauer et al., 2021, Theorem 3). The representations that are retrieved
77 from the MHNs replace the original representations of the query and support set molecules. Those
78 retrieved representations have amplified co-occurrences and covariance structures, while peculiarities
79 and spurious co-occurrences of the query and support set molecules are averaged out.

80 In this work, our contributions are the following:

- 81 • We propose a new architecture **MHNfs** for few-shot learning in drug discovery.
- 82 • We achieve a new state-of-the-art on the benchmarking dataset FS-Mol.
- 83 • We introduce a novel concept to enrich the molecule representations with context by associ-
84 ating them with a large set of context molecules.
- 85 • We add a naive baseline to the FS-Mol benchmark that yields better results than almost all
86 other published few-shot learning methods.
- 87 • We provide results of an ablation study and a domain shift experiment to further demonstrate
88 the effectiveness of our new method.

89 2 Problem setting

90 Drug discovery projects revolve around models $g(\mathbf{m})$ that can predict a molecular property or activity
91 \hat{y} , given a representation \mathbf{m} of an input molecule from a chemical space \mathcal{M} . We consider machine
92 learning models $\hat{y} = g_{\mathbf{w}}(\mathbf{m})$ with parameters \mathbf{w} that have been selected using a training set. Typically,
93 deep learning based property prediction uses a molecule encoder $f^{\text{ME}} : \mathcal{M} \rightarrow \mathbb{R}^d$. The molecule
94 encoder can process different symbolic or low-level representations of molecules, such as molecular
95 descriptors (Bender et al., 2004; Untertiner et al., 2014; Mayr et al., 2016), SMILES (Weininger,
96 1988; Mayr et al., 2018; Winter et al., 2019; Segler et al., 2018a), or molecular graphs (Merkwirth
97 and Lengauer, 2005; Kearnes et al., 2016; Yang et al., 2019; Jiang et al., 2021) and can be pre-trained
98 on related property prediction tasks.

99 For few-shot learning, the goal is to select a high-quality predictive model based on a small set of
100 molecules $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ with associated measurements $\mathbf{y} = \{y_1, \dots, y_N\}$. The measurements are
101 usually assumed to be binary $y_n \in \{-1, 1\}$, corresponding to the molecule being inactive or active.
102 The set $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ is called the *support set* that contains samples from a prediction task and N is
103 the *support set size*. The goal is to construct a model that correctly predicts y for an \mathbf{x} that is not in
104 the support set — in other words, a model that generalizes well.

105 Standard supervised machine learning approaches typically just show limited predictive power at
106 this task (Stanley et al., 2021) since they tend to overfit on the support set due to a small number of
107 training samples. These approaches learn the parameters \mathbf{w} of the model $g_{\mathbf{w}}$ from the support set in
108 a supervised manner. However, they heavily overfit to the support set when N is small. Therefore,
109 few-shot learning methods are necessary to construct models from the support set that generalize
110 well to new data.

111 3 MHNfs: Hopfield-based molecular context enrichment for few-shot drug 112 discovery

113 We aim at increasing the generalization capabilities of few-shot learning methods in drug discovery
114 by enriching the molecule representations with molecular context. In comparison to the support set,
115 which encodes information about the task, the context set – i.e. a large set of molecules – includes
116 information about a large chemical space. The query and the support set molecules perform a retrieval
117 from the context set and thereby enrich their representations. We detail this in the following.

118 3.1 Model architecture

119 We propose an architecture which consists of three consecutive modules. The first module, a) the
120 *context module* f^{CM} , enriches molecule representations by retrieving from a large set of molecules.
121 The second module, b) the *cross-attention module* f^{CAM} (Hou et al., 2019; Chen et al., 2021), enables
122 the effective exchange of information between the query molecule and the support set molecules.
123 Finally the prediction for the query molecule is computed by using the usual c) *similarity module*
124 f^{SM} (Koch et al., 2015; Altae-Tran et al., 2017):

$$\begin{aligned} \text{context module:} \quad \mathbf{m}' &= f^{\text{CM}}(\mathbf{m}, \mathbf{C}) \\ \mathbf{X}' &= f^{\text{CM}}(\mathbf{X}, \mathbf{C}), \end{aligned} \quad (1)$$

$$\text{cross-attention module:} \quad [\mathbf{m}'', \mathbf{X}''] = f^{\text{CAM}}([\mathbf{m}', \mathbf{X}']), \quad (2)$$

$$\text{similarity module:} \quad \hat{y} = f^{\text{SM}}(\mathbf{m}'', \mathbf{X}'', \mathbf{y}), \quad (3)$$

125 where $\mathbf{m} \in \mathbb{R}^d$ is a molecule embedding from a trainable or fixed molecule encoder, and \mathbf{m}' and \mathbf{m}''
126 are enriched versions of it. Similarly, $\mathbf{X} \in \mathbb{R}^{d \times N}$ contains the stacked embeddings of the support set
127 molecules and \mathbf{X}' and \mathbf{X}'' are their enriched versions. $\mathbf{C} \in \mathbb{R}^{d \times M}$ is a large set of stacked molecule
128 embeddings, \mathbf{y} are the support set labels, and \hat{y} is the prediction for the query molecule. Square
129 brackets indicate concatenation, for example $[\mathbf{m}', \mathbf{X}']$ is a matrix with $N + 1$ columns. The modules
130 f^{CM} , f^{CAM} , and f^{SM} are detailed in the paragraphs below. An overview of our architecture is given
131 in Figure 1. The architecture also includes skip connections bypassing $f^{\text{CM}}(\cdot, \cdot)$ and $f^{\text{CAM}}(\cdot)$ and
132 layer normalization (Ba et al., 2016), which are not shown in Figure 1.

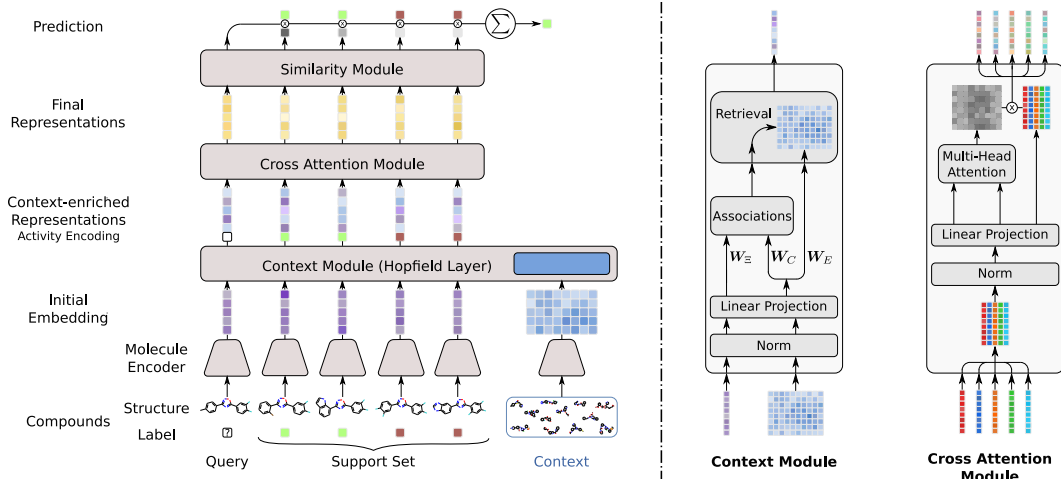


Figure 1: Schematic overview of our architecture. **Left:** All molecules are fed through a shared molecule encoder to obtain embeddings. Then, the context module (CM) enriches the representations by associating them with context molecules. The cross-attention module (CAM) enriches representations by mutually associating the query and support set molecules. Finally, the similarity module computes the prediction for the query molecule. **Right:** Detailed depiction of the operations in the CM and the CAM.

133 A shared molecule encoder f^{ME} creates embeddings for the query molecule $\mathbf{m} = f^{\text{ME}}(m)$, the
 134 support set molecules $\mathbf{x}_n = f^{\text{ME}}(x_n)$, and the context molecules $\mathbf{c}_m = f^{\text{ME}}(c_m)$. There are many
 135 possible choices for fixed or adaptive molecule encoders (see Section 2), of which we use descriptor-
 136 based fully-connected networks because of their computational efficiency and good accuracy (Dahl
 137 et al., 2014; Mayr et al., 2016, 2018). For notational clarity we denote the course of the representations
 138 through the architecture:

$$\begin{array}{ccccccc}
 m & \xrightarrow{f^{\text{ME}}} & \mathbf{m} & \xrightarrow{f^{\text{CM}}} & \mathbf{m}' & \xrightarrow{f^{\text{CAM}}} & \mathbf{m}'' \\
 \text{symbolic or} & & \text{molecule} & & \text{context} & & \text{similarity} \\
 \text{low-level repr.} & & \text{embedding} & & \text{repr.} & & \text{repr.}
 \end{array} \quad (4)$$

$$\begin{array}{ccccccc}
 x_n & \xrightarrow{f^{\text{ME}}} & \mathbf{x}_n & \xrightarrow{f^{\text{CM}}} & \mathbf{x}'_n & \xrightarrow{f^{\text{CAM}}} & \mathbf{x}''_n \\
 \text{symbolic or} & & \text{molecule} & & \text{context} & & \text{similarity} \\
 \text{low-level repr.} & & \text{embedding} & & \text{repr.} & & \text{repr.}
 \end{array} \quad (5)$$

139 3.2 Context module (CM)

140 The context module associates the query and support set molecules with a large set of context
 141 molecules, and represents them as weighted average of context molecule embeddings. The context
 142 module is realised by a continuous Modern Hopfield Network (MHN) (Ramsauer et al., 2021). An
 143 MHN is a content-addressable associative memory which can be built into deep learning architectures.
 144 There exists an analogy between the energy update of MHNs and the attention mechanism of
 145 Transformers (Vaswani et al., 2017; Ramsauer et al., 2021). MHNs are capable of storing and
 146 retrieving patterns from a memory $\mathbf{M} \in \mathbb{R}^{e \times M}$ given a state pattern $\boldsymbol{\xi} \in \mathbb{R}^e$ that represents the query.
 147 The retrieved pattern $\boldsymbol{\xi}^{\text{new}} \in \mathbb{R}^e$ is obtained by

$$\boldsymbol{\xi}^{\text{new}} = \mathbf{M} \mathbf{p} = \mathbf{M} \text{softmax}(\beta \mathbf{M}^T \boldsymbol{\xi}), \quad (6)$$

148 where \mathbf{p} is called the vector of associations and β is a scaling factor or inverse temperature. Modern
 149 Hopfield Networks have been successfully applied to chemistry and computational immunology
 150 (Seidl et al., 2022; Widrich et al., 2020).

151 We use this mechanism in the form of a *Hopfield layer*, which first maps raw patterns to an associative
 152 space using linear transformations, and uses multiple simultaneous queries $\boldsymbol{\Xi} \in \mathbb{R}^{d \times N}$:

$$\text{Hopfield}(\boldsymbol{\Xi}, \mathbf{C}) := (\mathbf{W}_E \mathbf{C}) \text{softmax}(\beta (\mathbf{W}_C \mathbf{C})^T (\mathbf{W}_\Xi \boldsymbol{\Xi})), \quad (7)$$

153 where $\mathbf{W}_E \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_C, \mathbf{W}_\Xi \in \mathbb{R}^{e \times d}$ are trainable parameters of the Hopfield layer, softmax
 154 is applied column-wise, and β is a hyperparameter. Note that in principle the Ξ and C could have a
 155 different second dimension as long as the linear transformations map to the same dimension e . Note
 156 that all embeddings that enter this module are first layer normalized (Ba et al., 2016). Several of
 157 these Hopfield layers can run in parallel and we refer to them as "heads" in analogy to Transformers
 158 (Vaswani et al., 2017).

159 The context module of our new architecture uses a Hopfield layer, where the query patterns are the
 160 embeddings of the query molecule \mathbf{m} and the support set molecules \mathbf{X} . The memory is composed of
 161 embeddings of a large set of M molecules from a chemical space, for example reference molecules,
 162 here called context molecules C . Then the original embeddings \mathbf{m} and \mathbf{X} are replaced by the
 163 retrieved embeddings, which are weighted averages of context molecule embeddings:

$$\mathbf{m}' = \text{Hopfield}(\mathbf{m}, C) \quad \text{and} \quad \mathbf{X}' = \text{Hopfield}(\mathbf{X}, C). \quad (8)$$

164 This retrieval step reinforces the covariance structure of the retrieved representations (see Ap-
 165 pendix A.7). Note that the embeddings of the query and the support set molecules have not yet
 166 influenced each other. These updated representations \mathbf{m}' , \mathbf{X}' are passed to the cross-attention module.

167 3.3 Cross-attention module (CAM)

168 For embedding-based few-shot learning methods in the field of drug discovery, Altae-Tran et al. (2017)
 169 showed that the representations of the molecules can be enriched, if the architecture allows information
 170 exchange between query and support set molecules. Altae-Tran et al. (2017) uses an attention-
 171 enhanced LSTM variant which updates the query and the support set molecule representations in an
 172 iterative fashion, being aware of each other. We further develop this idea and combine it with the idea
 173 of using a transformer encoder layer (Vaswani et al., 2017) as a cross-attention module (Hou et al.,
 174 2019; Chen et al., 2021).

175 The cross-attention module updates the query molecule representation \mathbf{m}' and the support set
 176 molecule representations \mathbf{X}' by mutually exchanging information, using the usual Transformer
 177 mechanism:

$$[\mathbf{m}'', \mathbf{X}''] = \text{Hopfield}([\mathbf{m}', \mathbf{X}'], [\mathbf{m}', \mathbf{X}']), \quad (9)$$

178 where $[\mathbf{m}', \mathbf{X}'] \in \mathbb{R}^{d \times (N+1)}$ is the concatenation of the representations of the query molecule \mathbf{m}'
 179 with the support set molecules \mathbf{X}' and we exploited that the Transformer is a special case of the
 180 Hopfield layer. Again, normalization is applied (Ba et al., 2016) and multiple Hopfield layers, i.e.,
 181 heads, can run in parallel, be stacked, and equipped with skip-connections. The representations \mathbf{m}''
 182 and \mathbf{X}'' are passed to the similarity module.

183 3.4 Similarity module (SM)

184 In this module, pairwise similarity values $k(\mathbf{m}'', \mathbf{x}_n'')$ are computed between the representation of
 185 a query molecule \mathbf{m}'' and each molecule \mathbf{x}_n'' in the support set as done recently (Koch et al., 2015;
 186 Altae-Tran et al., 2017). Based on these similarity values, the activity for the query molecule is
 187 predicted, building a weighted mean over the support set labels:

$$\hat{y} = \sigma \left(\tau^{-1} \frac{1}{N} \sum_{n=1}^N y_n' k(\mathbf{m}'', \mathbf{x}_n'') \right), \quad (10)$$

188 where our architecture employs dot product similarity of normalized representations $k(\mathbf{m}'', \mathbf{x}_n'') =$
 189 $\mathbf{m}''^T \mathbf{x}_n''$. $\sigma(\cdot)$ is the sigmoid function and τ is a hyperparameter. Note that we use a balancing
 190 strategy for the labels $y_n' = \begin{cases} N/(2\sqrt{N_A}) & \text{if } y_n = 1 \\ -N/(2\sqrt{N_I}) & \text{else} \end{cases}$, where N_A is the number of actives and N_I
 191 is the number of inactives of the support set.

192 3.5 Architecture, hyperparameter selection, and training details

193 **Hyperparameters.** The main hyperparameters of our architecture are the number of heads, the
 194 embedding dimension, the dimension of the association space of the CAM and CM, the learning

195 rate schedule, the scaling parameter β , and the molecule encoder. The following hyperparameters
196 were selected by manual hyperparameter selection on the validation tasks. The molecule encoder
197 consists of a single layer with output size $d = 1024$ and SELU activation (Klambauer et al., 2017).
198 The CM consists of one Hopfield layer with 8 heads. The dimension e of the association space is set
199 to 512 and $\beta = 1/\sqrt{e}$. Since we use skip connections between all modules the output dimension of
200 the CM and CAM matches the input dimension. The CAM comprises one layer with 8 heads and an
201 association-space dimension of 1088. For the input to the CAM, an activity encoding was added to
202 the support set molecule representations to provide label information. The SM uses $\tau = 22.6$. For the
203 context set, we randomly sample 5% from a large set of molecules – i.e., the molecules in the FS-Mol
204 training split – for each batch. For inference, we used a fixed set of 5% of training set molecules as
205 the context set for each seed. We hypothesize that these choices about the context could be further
206 improved (Section 6). We provide considered and selected hyperparameters in Appendix A.1.6.

207 **Loss function, regularization and optimization.** We use the Adam optimizer (Kingma and Ba,
208 2014) to minimize the cross-entropy loss between the predicted and known activity labels. We use
209 a learning rate scheduler which includes a warm up phase, followed by a section with a constant
210 learning rate, which is 0.0001, and a third phase in which the learning rate steadily decreases. As a
211 regularization strategy, for the CM and the CAM a dropout rate of 0.5 is used. The molecule encoder
212 has a dropout with rate 0.1 for the input and 0.5 elsewhere (see also Appendix A.1.6).

213 **Compute time and resources.** Training a single MHNfs model on the benchmarking dataset FS-
214 Mol takes roughly 90 hours of wall-clock time on an A100 GPU. In total, roughly 15,000 GPU hours
215 were consumed for this work.

216 4 Related work

217 Several approaches to few-shot learning in drug discovery have been suggested (Altae-Tran et al.,
218 2017; Nguyen et al., 2020; Guo et al., 2021; Wang et al., 2021). (Nguyen et al., 2020) evaluated
219 the applicability of MAML and its variants to graph neural networks (GNNs) and (Guo et al., 2021)
220 also combine GNNs and meta-learning. (Altae-Tran et al., 2017) suggested an approach called
221 Iterative Refinement Long Short-Term Memory, in which query and support set embeddings can
222 share information and update their embeddings. Property-aware relation networks (PAR) (Wang
223 et al., 2021) use an attention mechanism to enrich representations from cluster centers and then learn
224 a relation graph between molecules. (Chen et al., 2022) propose to adaptively learn kernels and apply
225 their method to few-shot drug discovery with predictive performance for larger support set sizes.
226 Recently, (Stanley et al., 2021) generated a benchmark dataset for few-shot learning methods in drug
227 discovery and provided some baseline results.

228 Many successful deep neural network architectures use external memories, such as the neural Turing
229 machine (Graves et al., 2014), memory networks (Weston et al., 2014), end-to-end memory networks
230 (Sukhbaatar et al., 2015). Recently, the connection between continuous modern Hopfield networks
231 (Ramsauer et al., 2021), which are content-addressable associative memories, and Transformer
232 architectures (Vaswani et al., 2017) has been established. We refer to (Le, 2021) for an extensive
233 overview of memory-based architectures. Architectures with external memories have also been used
234 for meta-learning (Vinyals et al., 2016; Santoro et al., 2016) and few-shot learning (Munkhdalai and
235 Yu, 2017; Ramalho and Garnelo, 2018; Ma et al., 2021).

236 5 Experiments

237 5.1 Benchmarking on FS-Mol

238 **Experimental setup.** Recently, the dataset FS-Mol (Stanley et al., 2021) was proposed to benchmark
239 few-shot learning methods in drug discovery. It was extracted from ChEMBL27 and comprises in
240 total 489,133 measurements, 233,786 compounds and 5,120 tasks. Per task, the mean number of
241 data points is 94. The dataset is well balanced as the mean ratio of active and inactive molecules is
242 close to 1. The FS-Mol benchmark dataset defines 4,938 training, 40 validation and 157 test tasks,
243 guaranteeing disjoint task sets. (Stanley et al., 2021) precomputed extended connectivity fingerprints
244 (ECFP) (Rogers and Hahn, 2010) and key molecular physical descriptors, which were defined by
245 RDKit (Landrum et al., 2006). While methods would be allowed to use other representations of
246 the input molecules, such as the molecular graph, we used a concatenation of these ECFPs and

Table 1: Results on FS-MOL [Δ AUC-PR]. The best method is marked bold. Error bars represent standard errors across tasks according to Stanley et al. (2021). The metrics are also averaged across five training re-runs and ten draws of support sets. In brackets the number of tasks per category is reported.

Method	All [157]	Kin. [125]	Hydrol. [20]	Oxid.[7]
GNN-ST ^a (Stanley et al., 2021)	.029 \pm .004	.027 \pm .004	.040 \pm .018	.020 \pm .016
MAT ^a (Maziarka et al., 2020)	.052 \pm .005	.043 \pm .005	.095 \pm .019	.062 \pm .024
Random Forest ^a (Breiman, 2001)	.092 \pm .007	.081 \pm .009	.158 \pm .028	.080 \pm .029
GNN-MT ^a (Stanley et al., 2021)	.093 \pm .006	.093 \pm .006	.108 \pm .025	.053 \pm .018
Similarity Search	.118 \pm .008	.109 \pm .008	.166 \pm .029	.097 \pm .033
GNN-MAML ^a (Guo et al., 2021)	.159 \pm .009	.177 \pm .009	.105 \pm .024	.054 \pm .028
PAR(Wang et al., 2021)	.164 \pm .008	.182 \pm .009	.109 \pm .020	.039 \pm .008
Frequent hitters	.182 \pm .010	.207 \pm .009	.098 \pm .009	.041 \pm .005
ProtoNet ^a (Snell et al., 2017)	.207 \pm .008	.215 \pm .009	.209 \pm .030	.095 \pm .029
Siamese Networks (Koch et al., 2015)	.223 \pm .010	.241 \pm .010	.178 \pm .026	.082 \pm .025
IterRefLSTM (Altae-Tran et al., 2017)	.234 \pm .010	.251 \pm .010	.199 \pm .026	.098 \pm .027
ADKF-IFT ^b (Chen et al., 2022)	.234 \pm .009	.248 \pm .020	.217 \pm .017	.106 \pm .008
MHNfs (ours)	.241 \pm .009	.259 \pm .010	.199 \pm .027	.096 \pm .019

^a metrics from Stanley et al. (2021). ^b results from Chen et al. (2022).

247 RDKit-based descriptors. For the main benchmark, the support set size was fixed to 16, using a
 248 stratified random split. We use all these settings of FS-Mol and therefore ensure a fair method
 249 comparison.

250 **Methods compared.** Baselines for few-shot learning and our proposed method **MHNfs** were com-
 251 pared against each other. The **Frequent Hitters** model is a naive baseline that ignores the provided
 252 support set and therefore has to learn to predict the average activity of a molecule. This method can
 253 potentially discriminate so-called frequent-hitter molecules (Stork et al., 2019) against molecules
 254 that are inactive across many tasks. We also added **Similarity Search** (Cereto-Massagué et al., 2015)
 255 as a baseline. Similarity search is a standard cheminformatics technique, used in situations with
 256 single or few known actives. In the simplest case, the search finds similar molecules by computing
 257 a fingerprint or descriptor-representation of the molecules and using a similarity measure $k(\cdot, \cdot)$ —
 258 such as Tanimoto Similarity (Tanimoto, 1960). Thus, Similarity Search, as used in cheminformatics,
 259 can be formally written as $\hat{y} = 1/N \sum_{n=1}^N y_n k(\mathbf{m}, \mathbf{x}_n)$; where $\mathbf{x}_1, \dots, \mathbf{x}_n$ come from a
 260 fixed molecule encoder, such as chemical fingerprint or descriptor calculation. A natural extension
 261 of Similarity Search with fixed chemical descriptors is **Neural Similarity Search** or **Siamese**
 262 **networks** (Koch et al., 2015), which extend the classic similarity search by learning a molecule
 263 encoder: $\hat{y} = \sigma\left(\tau^{-1} \frac{1}{N} \sum_{n=1}^N y'_n f_w^{\text{ME}}(\mathbf{m})^T f_w^{\text{ME}}(\mathbf{x}_n)\right)$. Furthermore, we re-implemented the
 264 **IterRefLSTM** (Altae-Tran et al., 2017) in Pytorch. The **IterRefLSTM** model consists of three
 265 modules. First, a molecule encoder maps the query and support set molecules to its representations
 266 \mathbf{m} and \mathbf{X} . Second, an attention-enhanced LSTM variant, the actual **IterRefLSTM**, iteratively
 267 updates the query and support set molecules, enabling information sharing between the molecules:
 268 $[\mathbf{m}', \mathbf{X}'] = \text{IterRefLSTM}_L([\mathbf{m}, \mathbf{X}])$, where the hyperparameter L controls the number of iteration
 269 steps of the **IterRefLSTM**. Third, a similarity module computes attention weights based on the rep-
 270 resentations: $\mathbf{a} = \text{softmax}(k(\mathbf{m}', \mathbf{X}'))$. These representations are then used for the final prediction:
 271 $\hat{y} = \sum_{i=1}^N a_i y_i$. For further details, see Appendix A.1.5. The **Random Forest** baseline uses the
 272 chemical descriptors and is trained in standard supervised manner on the support set molecules for
 273 each task. The method **GNN-ST** is a graph neural network (Stanley et al., 2021; Gilmer et al., 2017)
 274 that is trained from scratch for each task. The **GNN-MT** uses a two step strategy: First, the model is
 275 pretrained on a large dataset on related tasks; second, an output layer is constructed to the few-shot
 276 task via linear probing (Stanley et al., 2021; Alain and Bengio, 2016). The **Molecule Attention**
 277 **Transformer (MAT)** is pre-trained in a self-supervised fashion and fine-tuning is performed for the
 278 few-shot task (Maziarka et al., 2020). **GNN-MAML** is based on MAML (Finn et al., 2017), and uses
 279 a model-agnostic meta-learning strategy to find a general core model from which one can easily adapt
 280 to single tasks. **ProtoNet** (Snell et al., 2017) includes a molecule encoder, which maps query and
 281 support set molecules to representations in an embedding space. In this embedding space, prototypical

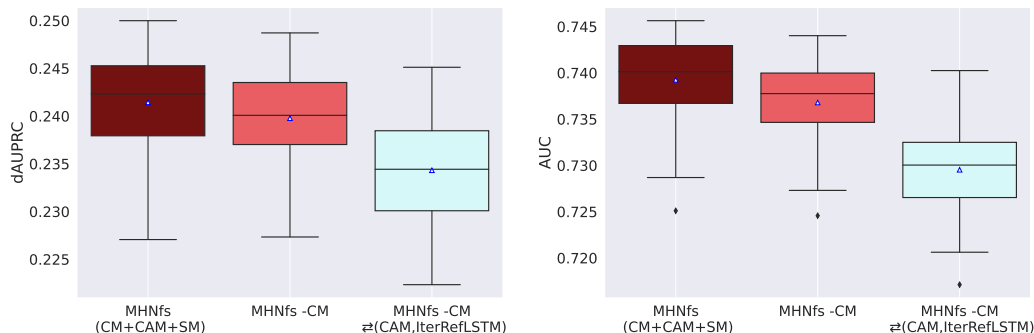


Figure 2: Results of the ablation study. The boxes show the median, mean and the variability of the average predictive performance of the methods across training re-runs and draws of support sets. The performance significantly drops when the context module is removed (light red bars), and when additionally the cross-attention module is replaced with the **IterRefLSTM** module (light blue bars). This indicates that our two newly introduced modules, CM and CAM, play a crucial role in MHNfs.

282 representations of each class are built by taking the mean across all related support set molecules for
 283 each class (details in Appendix A.1.4). For all methods the most important hyperparameters were
 284 adjusted on the validation tasks of FS-Mol. The **PAR** model (Wang et al., 2021) includes a GNN
 285 which creates initial molecule embeddings. These molecule embeddings are then enriched by an
 286 attention mechanism. Finally, another GNN learns relations between support and query set molecules.
 287 The **PAR** model has shown good results for datasets which just include very few tasks such as Tox21
 288 (Wang et al., 2021). Chen et al. (2022) suggest a framework for learning deep kernels by interpolating
 289 between meta-learning and conventional deep kernels, which results in the **ADKF-IFT** model. The
 290 model has exhibited especially high performance for large support set sizes.

291 **Training and evaluation.** For the model implementations, we used PyTorch (Paszke et al., 2019,
 292 BSD license). We used PyTorch Lightning (Falcon et al., 2019, Apache 2.0 license) as a framework
 293 for training and test logic, hydra for config file handling (Yadan, 2019, Apache 2.0 license) and
 294 Weights & Biases (Biewald, 2020, MIT license) as an experiment tracking tool. We performed five
 295 training reruns with different seeds for all methods, except Classic Similarity Search as there is
 296 no variability across seeds. Each model was evaluated ten times by drawing support sets with ten
 297 different seeds.

298 **Results.** The results in terms of area under precision-recall curve (AUC-PR) are presented in Table 1,
 299 where the difference to a random classifier is reported (Δ AUC-PR). The standard error is reported
 300 across tasks. Surprisingly, the naive baseline **Frequent Hitters**, that neglects the support set, has out-
 301 performed most of the few-shot learning methods, except for the embedding-based methods **Siamese**
 302 **Networks**, **ProtoNet**, **IterRefLSTM**, and **MHNfs**. **IterRefLSTM**, which has not been included
 303 in the FS-Mol benchmark study, reaches the second best performance. **MHNfs** has outperformed
 304 all other methods with respect to Δ AUC-PR across all tasks, including the **IterRefLSTM** model
 305 (p -value $1.72e-7$, paired Wilcoxon test), the **ADKF-IFT** model (p -value $<1.0e-8$, Wilcoxon test), and
 306 the **PAR** model (p -value $<1.0e-8$, paired Wilcoxon test).

307 5.2 Ablation study

308 **MHNfs** has two new main components compared to the previous state-of-the-art method **Iter-**
 309 **RefLSTM**: i) the context module, and ii) the cross-attention module which replaces the LSTM-like
 310 module. To assess the effects of these components, we performed an ablation study. Therefore,
 311 we compared **MHNfs** to a method that does not have the context module ("MHNfs -CM") and to
 312 a method that does not have the context module and uses an LSTM-like module instead of the
 313 CAM ("MHNfs -CM \rightleftharpoons (CAM,IterRefLSTM)"). For the ablation study, we used all 5 training reruns
 314 and evaluated each model 10 times on the test set with different support sets. The results of this
 315 ablation steps are presented in Figure 2. Both removing the CM and exchanging the CAM with
 316 the **IterRefLSTM** module were detrimental for the performance of the method (p -value 0.002 and
 317 $1.72e-7$, respectively; paired Wilcoxon test). The difference was even more pronounced under
 318 domain shift (see Appendix A.3.3). Appendix A.3.2 contains a second ablation study that examines

319 the overall effects of the context, the cross-attention, the similarity module, and the molecule encoder
320 of **MHNfs**.

321 5.3 Domain shift experiment

322 We performed an experiment in which we evaluate models, that were pretrained on FS-Mol, on the
323 Tox21 (Mayr et al., 2016) dataset. There is a strong domain shift from the drug-like molecules of
324 FS-Mol to the environmental chemicals, pesticides, and food additives of Tox21, such this dataset
325 poses a challenging setting for few-shot learning methods. The experiment is described in detail
326 in Appendix A.2. Our **MHNfs** approach has reached an AUC of $.679 \pm .018$ and has significantly
327 outperformed the **IterRefLSTM**-based model ($p_{\Delta\text{AUC-PR}}$ -value $3.4e-5$, paired Wilcoxon test) and
328 the Classic Similarity Search ($p_{\Delta\text{AUC-PR}}$ -value $2.4e-9$ paired Wilcoxon test) and therefore showed
329 robust performance on the toxicity domain, see Table A6.

330 6 Conclusion and discussion

331 We have introduced a new architecture for few-shot learning in drug discovery that is based on
332 the novel concept to enrich molecule representations with context. In a benchmarking experiment,
333 the architecture was assessed for its ability to learn accurate predictive models from small sets of
334 labelled molecules and in this setting it outperformed all other methods. In a domain shift study, the
335 robustness and transferability of the learned models has been assessed and again **MHNfs** exhibited
336 the best performance. The resulting predictive models often reach an AUC larger than .70, which
337 means that enrichment of active molecules is expected (Simm et al., 2018) when the models are used
338 for virtual screening. It has not escaped our notice that the specific context module we have proposed
339 could immediately be used for few-shot learning tasks in computer vision, but might be hampered
340 by computational constraints. **Limitations.** While the implementation of our method is currently
341 limited to small, organic drug-like molecules as inputs, our conceptual approach can also be used
342 for macro-molecules such as RNA, DNA or proteins. The output domain of our method comprises
343 biological effects, such that the prediction must be understood in that domain. Our method demands
344 higher computational costs and memory footprint as other embedding-based methods because of
345 the calculations necessary for the context module. While we hypothesize that our approach could
346 also be successful for similar data in the materials science domain, this has not been assessed. Our
347 study is also constrained by a limited amount of hyperparameter search for all methods. Deep
348 learning methods usually have a large number of hyperparameters, such as hidden dimensions,
349 number of layers, learning rates, of which we were only able to explore the most important ones. The
350 composition and choice of the context set is also under-explored and might be improved by selecting
351 reference molecules with an appropriate strategy. **Broader impact.** *Impact on machine learning and*
352 *related scientific fields.* We envision that with (a) the increasing availability of drug discovery and
353 material science datasets, (b) further improved biotechnologies, and (c) accounting for characteristics
354 of individuals, the drug and materials discovery process will be made more efficient. For machine
355 learning and artificial intelligence, the novel way in which representations are enriched with context
356 might strengthen the general research stream to include more context into deep learning systems. Our
357 approach also shows that such a system is more robust against domain shifts, which could be a step
358 towards Broad AI (Chollet, 2019; Hochreiter, 2022). *Impact on society.* If the approach proves useful,
359 it could lead to a faster and more cost-efficient drug discovery process. Especially the COVID-19
360 pandemic has shown that it is crucial for humanity to speed up the drug discovery process to few years
361 or even months. We hope that this work contributes to this effort and eventually leads to safer drugs
362 developed faster. *Consequences of failures of the method.* As common with methods in machine
363 learning, potential danger lies in the possibility that users rely too much on our new approach and use
364 it without reflecting on the outcomes. Failures of the proposed method would lead to unsuccessful
365 wet lab validation and negative wet lab tests. Since the proposed algorithm does not directly suggest
366 treatment or therapy, human beings are not directly at risk of being treated with a harmful therapy.
367 Wet lab and in-vitro testing would indicate wrong decisions by the system. *Leveraging of biases in*
368 *the data and potential discrimination.* As for almost all machine learning methods, confounding
369 factors, lab or batch effects, could be used for classification. This might lead to biases in predictions
370 or uneven predictive performance across different drug targets or bioassays.

371 **References**

- 372 Adler, T., Brandstetter, J., Widrich, M., Mayr, A., Kreil, D., Kopp, M., Klambauer, G., and
373 Hochreiter, S. (2020). Cross-domain few-shot learning by representation fusion. *arXiv preprint*
374 *arXiv:2010.06498*.
- 375 Alain, G. and Bengio, Y. (2016). Understanding intermediate layers using linear classifier probes.
376 *arXiv preprint arXiv:1610.01644*.
- 377 Alperstein, Z., Cherkasov, A., and Rolfe, J. T. (2019). All smiles variational autoencoder. *arXiv*
378 *preprint arXiv:1905.13343*.
- 379 Altae-Tran, H., Ramsundar, B., Pappu, A. S., and Pande, V. (2017). Low data drug discovery with
380 one-shot learning. *ACS central science*, 3(4):283–293.
- 381 Antoniou, A. and Storkey, A. (2019). Assume, augment and learn: Unsupervised few-shot meta-
382 learning via random labels and data augmentation. *arXiv preprint arXiv:1902.09884*.
- 383 Arrowsmith, J. (2011). Phase ii failures: 2008-2010. *Nature reviews drug discovery*, 10(5).
- 384 Axelrod, S. and Gomez-Bombarelli, R. (2022). Geom, energy-annotated molecular conformations
385 for property prediction and molecular generation. *Scientific Data*, 9(1):1–14.
- 386 Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer normalization. *arXiv preprint*
387 *arXiv:1607.06450*.
- 388 Bender, A., Mussa, H. Y., Glen, R. C., and Reiling, S. (2004). Similarity searching of chemical
389 databases using atom environment descriptors (molprint 2d): evaluation of performance. *Journal*
390 *of chemical information and computer sciences*, 44(5):1708–1718.
- 391 Bendre, N., Marín, H. T., and Najafirad, P. (2020). Learning from few samples: A survey. *arXiv*
392 *preprint arXiv:2007.15484*.
- 393 Bengio, Y., Bengio, S., and Cloutier, J. (1991). Learning a synaptic learning rule. In *Seattle*
394 *international joint conference on neural networks*.
- 395 Biewald, L. (2020). Experiment tracking with weights and biases. Software available from
396 wandb.com.
- 397 Bonner, M. F. and Epstein, R. A. (2021). Object representations in the human brain reflect the
398 co-occurrence statistics of vision and language. *Nature Communications*, 12(4081).
- 399 Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- 400 Cereto-Massagué, A., Ojeda, M. J., Valls, C., Mulero, M., Garcia-Vallvé, S., and Pujadas, G. (2015).
401 Molecular fingerprint similarity search in virtual screening. *Methods*, 71:58–63.
- 402 Chen, H., Engkvist, O., Wang, Y., Olivecrona, M., and Blaschke, T. (2018). The rise of deep learning
403 in drug discovery. *Drug discovery today*, 23(6):1241–1250.
- 404 Chen, H., Li, H., Li, Y., and Chen, C. (2021). Sparse spatial transformers for few-shot learning. *arXiv*
405 *preprint arXiv:2109.12932*.
- 406 Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020). A simple framework for contrastive
407 learning of visual representations. In *International conference on machine learning*, pages 1597–
408 1607. PMLR.
- 409 Chen, W., Tripp, A., and Hernández-Lobato, J. M. (2022). Meta-learning feature representations for
410 adaptive gaussian processes via implicit differentiation. *arXiv preprint arXiv:2205.02708*.
- 411 Chollet, F. (2019). On the measure of intelligence. *arXiv preprint arXiv:1911.01547*.
- 412 Dahl, G. E., Jaitly, N., and Salakhutdinov, R. (2014). Multi-task neural networks for qsar predictions.
413 *arXiv preprint arXiv:1406.1231*.

- 414 Duvenaud, D., Maclaurin, D., Aguilera-Iparraguirre, J., Gómez-Bombarelli, R., Hirzel, T., Aspuru-
415 Guzik, A., and Adams, R. P. (2015). Convolutional networks on graphs for learning molecular
416 fingerprints. *arXiv preprint arXiv:1509.09292*.
- 417 Eckert, H. and Bajorath, J. (2007). Molecular similarity analysis in virtual screening: foundations,
418 limitations and novel approaches. *Drug discovery today*, 12(5-6):225–233.
- 419 Falcon, W. et al. (2019). Pytorch lightning. *GitHub*. Note: <https://github.com/PyTorchLightning/pytorch-lightning>, 3:6.
420
- 421 Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep
422 networks. In *International conference on machine learning*, pages 1126–1135. PMLR.
- 423 Fürst, A., Rumetshofer, E., Tran, V., Ramsauer, H., Tang, F., Lehner, J., Kreil, D., Kopp, M.,
424 Klambauer, G., Bitto-Nemling, A., et al. (2021). Cloob: Modern hopfield networks with infoloob
425 outperform clip. *arXiv preprint arXiv:2110.11316*.
- 426 Geppert, H., Horváth, T., Gärtner, T., Wrobel, S., and Bajorath, J. (2008). Support-vector-machine-
427 based ranking significantly improves the effectiveness of similarity searching using 2d fingerprints
428 and multiple reference compounds. *Journal of chemical information and modeling*, 48(4):742–746.
- 429 Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural message
430 passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272.
431 PMLR.
- 432 Gomez, L. (2018). Decision making in medicinal chemistry: The power of our intuition. *ACS*
433 *Medicinal Chemistry Letters*, 9(10):956–958.
- 434 Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling,
435 B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A.
436 (2018). Automatic chemical design using a data-driven continuous representation of molecules.
437 *ACS central science*, 4(2):268–276.
- 438 Graves, A., Wayne, G., and Danihelka, I. (2014). Neural turing machines. *arXiv preprint*
439 *arXiv:1410.5401*.
- 440 Guo, Z., Zhang, C., Yu, W., Herr, J., Wiest, O., Jiang, M., and Chawla, N. V. (2021). Few-shot graph
441 learning for molecular property prediction. In *Proceedings of the web conference 2021*, pages
442 2559–2567.
- 443 He, J., You, H., Sandström, E., Nittinger, E., Bjerrum, E. J., Tyrchan, C., Czechtizky, W., and
444 Engkvist, O. (2021). Molecular optimization by capturing chemist’s intuition using deep neural
445 networks. *Journal of cheminformatics*, 13(1):1–17.
- 446 Hertz, T., Hillel, A. B., and Weinshall, D. (2006). Learning a kernel function for classification with
447 small training samples. In *Proceedings of the 23rd international conference on machine learning*,
448 pages 401–408.
- 449 Hochreiter, S. (2022). Toward a broad ai. *Communications of the ACM*, 65(4):56–57.
- 450 Hochreiter, S., Klambauer, G., and Rarey, M. (2018). Machine learning in drug discovery. *Journal of*
451 *Chemical Information and Modeling*, 58(9):1723–1724.
- 452 Hochreiter, S., Younger, A. S., and Conwell, P. R. (2001). Learning to learn using gradient descent.
453 In *International conference on artificial neural networks*, pages 87–94. Springer.
- 454 Hou, R., Chang, H., Ma, B., Shan, S., and Chen, X. (2019). Cross attention network for few-shot
455 classification. *Advances in neural information processing systems* 32.
- 456 Huang, R., Xia, M., Nguyen, D.-T., Zhao, T., Sakamuru, S., Zhao, J., Shahane, S. A., Rossoshek,
457 A., and Simeonov, A. (2016a). Tox21challenge to build predictive models of nuclear receptor and
458 stress response pathways as mediated by exposure to environmental chemicals and drugs. *Frontiers*
459 *in Environmental Science*, 3:85.

460 Huang, R., Xia, M., Sakamuru, S., Zhao, J., Shahane, S. A., Attene-Ramos, M., Zhao, T., Austin,
461 C. P., and Simeonov, A. (2016b). Modelling the tox21 10 k chemical profiles for in vivo toxicity
462 prediction and mechanism characterization. *Nature communications*, 7(1):1–10.

463 Jiang, D., Wu, Z., Hsieh, C.-Y., Chen, G., Liao, B., Wang, Z., Shen, C., Cao, D., Wu, J., and Hou, T.
464 (2021). Could graph neural networks learn better molecular representation for drug discovery?
465 a comparison study of descriptor-based and graph-based models. *Journal of cheminformatics*,
466 13(1):1–23.

467 Kearnes, S., McCloskey, K., Berndl, M., Pande, V., and Riley, P. (2016). Molecular graph convolu-
468 tions: moving beyond fingerprints. *Journal of computer-aided molecular design*, 30(8):595–608.

469 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint*
470 *arXiv:1412.6980*.

471 Klambauer, G., Unterthiner, T., Mayr, A., and Hochreiter, S. (2017). Self-normalizing neural networks.
472 In *Advances in neural information processing systems 30*, pages 972–981.

473 Koch, G., Zemel, R., Salakhutdinov, R., et al. (2015). Siamese neural networks for one-shot image
474 recognition. In *ICML deep learning workshop*, volume 2. Lille.

475 Kuhn, M., Letunic, I., Jensen, L. J., and Bork, P. (2016). The sider database of drugs and side effects.
476 *Nucleic acids research*, 44(D1):D1075–D1079.

477 Landrum, G. et al. (2006). Rdkit: Open-source cheminformatics.

478 Le, H. (2021). Memory and attention in deep learning. *arXiv preprint arXiv:2107.01390*.

479 Li, J., Cai, D., and He, X. (2017). Learning graph-level representation for drug discovery. *arXiv*
480 *preprint arXiv:1709.03741*.

481 Li, P., Li, Y., Hsieh, C.-Y., Zhang, S., Liu, X., Liu, H., Song, S., and Yao, X. (2021). Trimnet: learning
482 molecular representation from triplet messages for biomedicine. *Briefings in Bioinformatics*,
483 22(4):bbaa266.

484 Ma, Y., Liu, W., Bai, S., Zhang, Q., Liu, A., Chen, W., and Liu, X. (2021). Few-shot visual
485 learning with contextual memory and fine-grained calibration. In *Proceedings of the Twenty-Ninth*
486 *International Conference on International Joint Conferences on Artificial Intelligence*, pages
487 811–817.

488 Mayr, A., Klambauer, G., Unterthiner, T., and Hochreiter, S. (2016). Deeptox: toxicity prediction
489 using deep learning. *Frontiers in environmental science*, 3:80.

490 Mayr, A., Klambauer, G., Unterthiner, T., Steijaert, M., Wegner, J. K., Ceulemans, H., Clevert, D.-A.,
491 and Hochreiter, S. (2018). Large-scale comparison of machine learning methods for drug target
492 prediction on chembl. *Chemical science*, 9(24):5441–5451.

493 Maziarka, Ł., Danel, T., Mucha, S., Rataj, K., Tabor, J., and Jastrzębski, S. (2020). Molecule attention
494 transformer. *arXiv preprint arXiv:2002.08264*.

495 Merk, D., Friedrich, L., Grisoni, F., and Schneider, G. (2018). De novo design of bioactive small
496 molecules by artificial intelligence. *Molecular informatics*, 37(1-2):1700153.

497 Merkwirth, C. and Lengauer, T. (2005). Automatic generation of complementary descriptors with
498 molecular graph networks. *Journal of chemical information and modeling*, 45(5):1159–1168.

499 Miller, E. G., Matsakis, N. E., and Viola, P. A. (2000). Learning from one example through shared
500 densities on transforms. In *Proceedings ieee conference on computer vision and pattern recognition.*
501 *cvpr 2000 (cat. no. PR00662)*, volume 1, pages 464–471.

502 Munkhdalai, T. and Yu, H. (2017). Meta networks. In *International Conference on Machine Learning*,
503 pages 2554–2563. PMLR.

504 Nguyen, C. Q., Kretsoulas, C., and Branson, K. M. (2020). Meta-learning gnn initializations for
505 low-resource molecular property prediction. *arXiv preprint arXiv:2003.05996*.

- 506 Oord, A. v. d., Li, Y., and Vinyals, O. (2018). Representation learning with contrastive predictive
507 coding. *arXiv preprint arXiv:1807.03748*.
- 508 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga,
509 L., and Lerer, A. (2019). Automatic differentiation in pytorch. In *Conference on neural information
510 processing systems*.
- 511 Potter, M. (2012). Conceptual short term memory in perception and thought. *Frontiers in Psychology*,
512 3:113.
- 513 Ramalho, T. and Garnelo, M. (2018). Adaptive posterior learning: few-shot learning with a surprise-
514 based memory module. In *International Conference on Learning Representations*.
- 515 Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Adler, T.,
516 Kreil, D., Kopp, M. K., Klambauer, G., Brandstetter, J., and Hochreiter, S. (2021). Hopfield
517 networks is all you need. In *International conference on learning representations*.
- 518 Riniker, S. and Landrum, G. A. (2013). Open-source platform to benchmark fingerprints for ligand-
519 based virtual screening. *Journal of cheminformatics*, 5(1):1–17.
- 520 Rogers, D. and Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of chemical information
521 and modeling*, 50(5):742–754.
- 522 Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. (2016). Meta-learning with
523 memory-augmented neural networks. In *International conference on machine learning*, pages
524 1842–1850. PMLR.
- 525 Schmidhuber, J. (1987). Evolutionary principles in self-referential learning.
- 526 Schneider, P., Walters, W. P., Plowright, A. T., Sieroka, N., Listgarten, J., Goodnow, R. A., Fisher,
527 J., Jansen, J. M., Duca, J. S., Rush, T. S., et al. (2020). Rethinking drug design in the artificial
528 intelligence era. *Nature reviews drug discovery*, 19(5):353–364.
- 529 Segler, M. H., Kogej, T., Tyrchan, C., and Waller, M. P. (2018a). Generating focused molecule
530 libraries for drug discovery with recurrent neural networks. *ACS central science*, 4(1):120–131.
- 531 Segler, M. H., Preuss, M., and Waller, M. P. (2018b). Planning chemical syntheses with deep neural
532 networks and symbolic ai. *Nature*, 555(7698):604–610.
- 533 Seidl, P., Renz, P., Dyubankova, N., Neves, P., Verhoeven, J., Wegner, J. K., Segler, M., Hochreiter,
534 S., and Klambauer, G. (2022). Improving few- and zero-shot reaction template prediction using
535 modern hopfield networks. *Journal of chemical information and modeling*, 62(9):2111–2120.
- 536 Sheridan, R. P. and Kearsley, S. K. (2002). Why do we need so many chemical similarity search
537 methods? *Drug discovery today*, 7(17):903–911.
- 538 Simm, J., Klambauer, G., Arany, A., Steijaert, M., Wegner, J. K., Gustin, E., Chupakhin, V., Chong,
539 Y. T., Vialard, J., Buijnsters, P., et al. (2018). Repurposed high-throughput image assays enables
540 biological activity prediction for drug discovery. *Cell Chemical Biology*, page 108399.
- 541 Snell, J., Swersky, K., and Zemel, R. S. (2017). Prototypical networks for few-shot learning. *arXiv
542 preprint arXiv:1703.05175*.
- 543 Stanley, M., Bronskill, J. F., Maziarz, K., Misztela, H., Lanini, J., Segler, M., Schneider, N., and
544 Brockschmidt, M. (2021). Fs-mol: A few-shot learning dataset of molecules. In *Conference on
545 neural information processing systems workshop*.
- 546 Stork, C., Chen, Y., Sicho, M., and Kirchmair, J. (2019). Hit dexter 2.0: machine-learning models for
547 the prediction of frequent hitters. *Journal of chemical information and modeling*, 59(3):1030–1043.
- 548 Sturm, N., Mayr, A., Le Van, T., Chupakhin, V., Ceulemans, H., Wegner, J., Golib-Dzib, J.-F.,
549 Jeliaskova, N., Vandriessche, Y., Böhm, S., et al. (2020). Industry-scale application and evaluation
550 of deep learning for drug target prediction. *Journal of Cheminformatics*, 12(1):1–13.

- 551 Sukhbaatar, S., Weston, J., Fergus, R., et al. (2015). End-to-end memory networks. *Advances in*
552 *neural information processing systems*, 28.
- 553 Sun, J., Jeliaskova, N., Chupakhin, V., Golib-Dzib, J.-F., Engkvist, O., Carlsson, L., Wegner, J.,
554 Ceulemans, H., Georgiev, I., Jeliaskov, V., et al. (2017). Excape-db: an integrated large scale
555 dataset facilitating big data analysis in chemogenomics. *Journal of cheminformatics*, 9(1):1–9.
- 556 Tanimoto, T. (1960). Ibm type 704 medical diagnosis program. *IRE transactions on medical*
557 *electronics*, (4):280–283.
- 558 Torres, L., Monteiro, N., Oliveira, J., Arrais, J., and Ribeiro, B. (2020). Exploring a siamese neural
559 network architecture for one-shot drug discovery. In *2020 IEEE 20th International Conference on*
560 *Bioinformatics and Bioengineering (BIBE)*, pages 168–175.
- 561 Unterthiner, T., Mayr, A., Klambauer, G., Steijaert, M., Wegner, J. K., Ceulemans, H., and Hochreiter,
562 S. (2014). Deep learning as an opportunity in virtual screening. In *Advances in neural information*
563 *processing systems workshop*.
- 564 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and
565 Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing*
566 *systems*, pages 5998–6008.
- 567 Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al. (2016). Matching networks for one shot
568 learning. *Advances in neural information processing systems*, 29:3630–3638.
- 569 Walters, W. P. and Barzilay, R. (2021). Critical assessment of ai in drug discovery. *Expert opinion on*
570 *drug discovery*, pages 1–11.
- 571 Wang, X., Huan, J., Smalter, A., and Lushington, G. H. (2010). Application of kernel functions for
572 accurate similarity search in large chemical databases. In *BMC bioinformatics*, volume 11, pages
573 1–14. BioMed Central.
- 574 Wang, Y., Abuduweili, A., Yao, Q., and Dou, D. (2021). Property-aware relation networks for
575 few-shot molecular property prediction. *Advances in Neural Information Processing Systems*,
576 34:17441–17454.
- 577 Wang, Y., Yao, Q., Kwok, J. T., and Ni, L. M. (2020). Generalizing from a few examples: A survey
578 on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34.
- 579 Waring, M. J., Arrowsmith, J., Leach, A. R., Leeson, P. D., Mandrell, S., Owen, R. M., Pairaudeau, G.,
580 Pennie, W. D., Pickett, S. D., Wang, J., et al. (2015). An analysis of the attrition of drug candidates
581 from four major pharmaceutical companies. *Nature reviews drug discovery*, 14(7):475–486.
- 582 Weininger, D. (1988). Smiles, a chemical language and information system. 1. introduction to
583 methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–
584 36.
- 585 Weston, J., Chopra, S., and Bordes, A. (2014). Memory networks. *arXiv preprint arXiv:1410.3916*.
- 586 Widrich, M., Schäfl, B., Pavlović, M., Ramsauer, H., Gruber, L., Holzleitner, M., Brandstetter, J.,
587 Sandve, G. K., Greiff, V., Hochreiter, S., et al. (2020). Modern hopfield networks and attention for
588 immune repertoire classification. In *Advances in neural information processing systems* 33.
- 589 Willett, P. (2014). The calculation of molecular structural similarity: principles and practice. *Molecu-*
590 *lar informatics*, 33(6-7):403–413.
- 591 Winter, R., Montanari, F., Noé, F., and Clevert, D.-A. (2019). Learning continuous and data-driven
592 molecular descriptors by translating equivalent chemical representations. *Chemical science*,
593 10(6):1692–1701.
- 594 Wu, Z., Ramsundar, B., Feinberg, E. N., Gomes, J., Geniesse, C., Pappu, A. S., Leswing, K., and
595 Pande, V. (2018). Moleculenet: a benchmark for molecular machine learning. *Chemical science*,
596 9(2):513–530.

- 597 Yadan, O. (2019). Hydra - a framework for elegantly configuring complex applications. Github.
598 Visited 2022-04-25.
- 599 Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A., Hopper, T., Kelley,
600 B., Mathea, M., et al. (2019). Analyzing learned molecular representations for property prediction.
601 *Journal of chemical information and modeling*, 59(8):3370–3388.
- 602 Ye, M. and Guo, Y. (2018). Deep triplet ranking networks for one-shot recognition. *arXiv preprint*
603 *arXiv:1804.07275*.
- 604 Zaslavskiy, M., Jégou, S., Tramel, E. W., and Wainrib, G. (2019). Toxicblend: Virtual screening of
605 toxic compounds with ensemble predictors. *Computational Toxicology*, 10:81–88.
- 606 Zhao, A., Balakrishnan, G., Durand, F., Guttag, J. V., and Dalca, A. V. (2019). Data augmentation
607 using learned transformations for one-shot medical image segmentation. In *Proceedings of the*
608 *ieee conference on computer vision and pattern recognition*, pages 8543–8553.

609 **Checklist**

- 610 1. For all authors...
- 611 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
612 contributions and scope? [Yes]
- 613 (b) Did you describe the limitations of your work? [Yes] See Section 6.
- 614 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
615 Section 6.
- 616 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
617 them? [Yes]
- 618 2. If you are including theoretical results...
- 619 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 620 (b) Did you include complete proofs of all theoretical results? [N/A]
- 621 3. If you ran experiments...
- 622 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
623 mental results (either in the supplemental material or as a URL)? [Yes]
- 624 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
625 were chosen)? [Yes] We refer e.g. to Section 5 in which we provide this information.
- 626 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
627 ments multiple times)? [Yes] We report error bars for all performance metrics. For all
628 experiments the variability across re-runs, different support sets and prediction tasks
629 was assessed.
- 630 (d) Did you include the total amount of compute and the type of resources used (e.g., type
631 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 3.5.
- 632 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 633 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 634 (b) Did you mention the license of the assets? [Yes]
- 635 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 636 (d) Did you discuss whether and how consent was obtained from people whose data you’re
637 using/curating? [N/A]
- 638 (e) Did you discuss whether the data you are using/curating contains personally identifiable
639 information or offensive content? [Yes] See Section 6.
- 640 5. If you used crowdsourcing or conducted research with human subjects...
- 641 (a) Did you include the full text of instructions given to participants and screenshots, if
642 applicable? [N/A]
- 643 (b) Did you describe any potential participant risks, with links to Institutional Review
644 Board (IRB) approvals, if applicable? [N/A]
- 645 (c) Did you include the estimated hourly wage paid to participants and the total amount
646 spent on participant compensation? [N/A]

647 **Contents of the appendix**

648	A Appendix	17
649	A.1 Details on methods	17
650	A.1.1 Frequent hitters: details and hyperparameters	17
651	A.1.2 Classic similarity search: details and hyperparameters	18
652	A.1.3 Neural Similarity Search or Siamese networks: details and hyperparameters	19
653	A.1.4 ProtoNet: details and hyperparameters	19
654	A.1.5 IterRefLSTM: details and hyperparameters	20
655	A.1.6 MHNfs: details and hyperparameters	21
656	A.1.7 PAR: details and hyperparameters	22
657	A.2 Domain shift experiment	24
658	A.3 Details on the ablation study	25
659	A.3.1 Ablation study A: comparison against IterRefLSTM	25
660	A.3.2 Ablation study B: all design elements	25
661	A.3.3 Ablation study C: Under domain shift on Tox21	26
662	A.4 Generalization to different support set sizes	26
663	A.5 Generalization to different context sets	26
664	A.6 Details and insights on the context module	27
665	A.7 Reinforcing the covariance structure in the data using modern Hopfield networks .	27

666 **A Appendix**

667 **A.1 Details on methods**

668 Few-shot learning methods in drug discovery can be described as models with adaptive parameters w
669 that use a support set $\mathbf{Z} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}^1$ as additional input to predict a label \hat{y} for a
670 molecule m

$$\hat{y} = g_w(m, \mathbf{Z}). \quad (\text{A1})$$

671 Optimization-based methods, such as MAML (Finn et al., 2017), use the support set to update the
672 parameters w

$$\hat{y} = g_{a(w; \mathbf{Z})}(m), \quad (\text{A2})$$

673 where $a(\cdot)$ is a function that adapts w of g based on \mathbf{Z} for example via gradient-descent.

674 Embedding-based methods use a different approach and learn representations of the support set
675 molecules $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, sometimes written as stacked embeddings $\mathbf{X} \in \mathbb{R}^{d \times N}$, and the query
676 molecule m , and some function that associates these two types of information with each other. We
677 describe the embedding-based methods Similarity Search in Section A.1.2, Neural Similarity Search
678 in Section A.1.3, ProtoNet in Section A.1.4, IterRefLSTM in Section A.1.5, PAR in Section A.1.7,
679 and MHNfs in the main paper and details in Section A.1.6. The "frequent hitters" baseline is described
680 in Section A.1.1.

681 **A.1.1 Frequent hitters: details and hyperparameters**

682 The "frequent hitters" model g^{FH} is a baseline that we implemented and included in the method
683 comparison. This method uses the usual training scheme of sampling a query molecule m with a
684 label y , having access to a support set \mathbf{Z} . In contrast to the usual models of the type $g_w(m, \mathbf{Z})$, the
685 frequent hitters model g^{FH} neglects the support set:

$$\hat{y} = g_w^{\text{FH}}(m). \quad (\text{A3})$$

686 Thus, during training for the same molecule m , the model might have to predict both $y = 1$ and
687 $y = -1$, since the molecule can be active in one task and inactive in another task. Therefore, the

¹We use \mathbf{Z} to denote the support set of already embedded molecules to keep the notation uncluttered. More correctly, the methods have access to the raw support set $Z = \{(x_1, y_1), \dots, (x_N, y_N)\}$, where x_n is a symbolic, such as the molecular graph, or low-level representation of the molecule.

Table A1: Hyperparameter space considered for the Frequent hitters model. The hyperparameters of the best configuration are marked bold.

Hyperparameter	Explored values
Number of hidden layers	1, 2 , 4
Number of units per hidden layer	1024, 2048 , 4096
Output dimension	512 , 1024
Activation function	ReLU
Learning rate	0.0001 , 0.001
Optimizer	Adam, AdamW
Weight decay	0, 0.01
Batch size	32, 128, 512, 2048, 4096
Input Dropout	0, 0.1
Dropout	0.1, 0.2, 0.3, 0.4 , 0.5
Layer-normalization	False, True
• Affine	False , True
Similarity function	dot product

688 model tends to predict average activity of a molecule to minimize the cross-entropy loss. We chose
 689 an additive combination of the Morgan fingerprints, RDKit fingerprints, and MACCS keys for the
 690 input representation to the MLP.

691 **Hyperparameter search.** We performed manual hyperparameter search on the validation set and
 692 report the explored hyperparameter space (Table A1). We use early-stopping based on validation
 693 average-precision, a patience of 3 epochs and train for a maximum of 20 epochs with a linear warm-up
 694 learning-rate schedule for the first 3 epochs.

695 A.1.2 Classic similarity search: details and hyperparameters

696 Similarity Search (Cereto-Massagué et al., 2015) is a classic chemoinformatics technique used in
 697 situations in which a single or few actives are known. In the simplest case, molecules that are similar to
 698 a given active molecule are searched by computing a fingerprint or descriptor-representation $f^{\text{desc}}(\mathbf{m})$
 699 of the molecules and using a similarity measure $k(\cdot, \cdot)$, such as Tanimoto Similarity (Tanimoto, 1960).
 700 Thus, the Similarity Search as used in chemoinformatics can be formally written as:

$$\hat{y} = 1/N \sum_{n=1}^N y_n k(f^{\text{desc}}(\mathbf{m}), f^{\text{desc}}(x_n)), \quad (\text{A4})$$

701 where the function f^{desc} maps the molecule to its chemical descriptors or fingerprints and takes
 702 the role of both the molecule encoder and the support set encoder. The association function f^{assoc}
 703 consists of a) the similarity measure $k(\cdot, \cdot)$ and then b) mean pooling across molecules weighted by
 704 their similarity and activity.

705 Notably, there are many variants of Similarity Search (Cereto-Massagué et al., 2015; Wang et al.,
 706 2010; Eckert and Bajorath, 2007; Geppert et al., 2008; Willett, 2014; Sheridan and Kearsley, 2002;
 707 Riniker and Landrum, 2013) of which some correspond to recent few-shot learning methods with a
 708 fixed molecule encoder. For example, (Geppert et al., 2008) suggest to use centroid molecules, i.e.,
 709 prototypes or averages of active molecules. This is equivalent to the idea of Prototypical Networks
 710 (Snell et al., 2017). Riniker and Landrum (2013) are aware of different fusion strategies for sets of
 711 active or inactive molecules, which corresponds to different pooling strategies of the support set.
 712 Overall, the variants of the classic Similarity Search are highly similar to embedding-based few-shot
 713 learning methods except that they have a fixed instead of a learned molecule encoder.

714 **Hyperparameter search.** For the Similarity Search, there were two decisions to make which was
 715 firstly the similarity metric and secondly the question whether we should use a balancing strategy like
 716 shown in Section 3.4. We decided for the dot-product as a similarity metric and using the balancing
 717 strategy. These decisions were made by evaluating the models on the validation set.

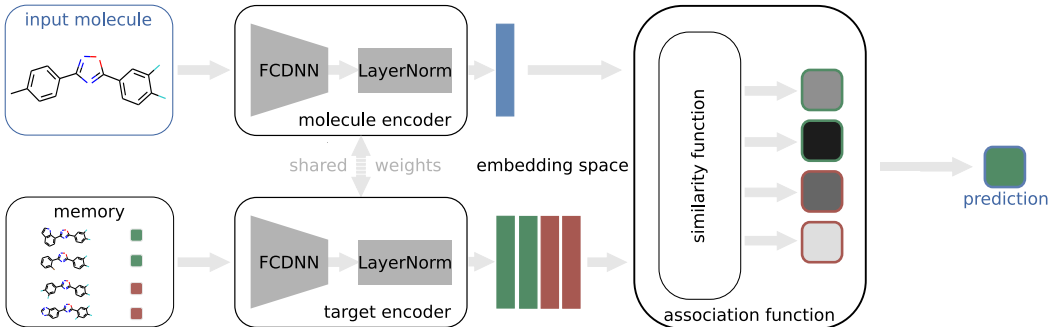


Figure A1: Schematic overview of the implemented Neural Similarity Search variant

718 A.1.3 Neural Similarity Search or Siamese networks: details and hyperparameters

719 A lot of related work already was done (Koch et al., 2015; Hertz et al., 2006; Ye and Guo, 2018;
 720 Torres et al., 2020). We adapted these ideas, such that a fully-connected deep neural network followed
 721 by a Layer Normalization (Ba et al., 2016) operation, f_w^{ME} with adaptive parameters w , is used in a
 722 Siamese fashion to compute the embeddings for the input molecule and the support set molecules.
 723 Within the association function block, pairwise similarity values for the input molecule and each
 724 support set molecule are computed, associating both embeddings via the dot product. Based on these
 725 similarity values, the activity for the input molecule is predicted, building the weighted mean over
 726 the support set molecule labels:

$$\hat{y} = \sigma \left(\tau^{-1} \frac{1}{N} \sum_{n=1}^N y'_n f^{\text{ME}}(m)^T f^{\text{ME}}(x_n) \right), \quad (\text{A5})$$

727 where $\sigma(\cdot)$ is the sigmoid function and τ is a hyperparameter in the range of \sqrt{d} . Note that this
 728 method uses a balancing strategy for the labels $y'_n = \begin{cases} N/(2\sqrt{N_A}) & \text{if } y_n = 1 \\ -N/(2\sqrt{N_I}) & \text{else} \end{cases}$, where N_A is the
 729 number of actives and N_I is the number of inactives of the support set. Figure A1 provides an
 730 schematic overview of the Neural Similarity Search variant.

731 We trained the networks using the Adam optimizer (Kingma and Ba, 2014) to minimize binary
 732 cross-entropy loss.

733 **Hyperparameter search.** We performed manual hyperparameter search on the validation set. We
 734 report the explored hyperparameter space (Table A2). Bold values indicate the selected hyperparam-
 735 eters for the final model.

736 A.1.4 ProtoNet: details and hyperparameters

737 Prototypical Networks (ProtoNet) (Snell et al., 2017), learn a prototype r for each class. Concretely,
 738 the support set Z is class-wise separated into $Z^+ := \{(x, y) \in Z \mid y = 1\}$ and $Z^- := \{(x, y) \in Z \mid$
 739 $y = -1\}$. For the subsets Z^+ and Z^- prototypical representations r^+ and r^- can be computed by

$$r^+ = \frac{1}{|Z^+|} \cdot \sum_{(x,y) \in Z^+} f^{\text{ME}}(x) \quad (\text{A6})$$

740 and

$$r^- = \frac{1}{|Z^-|} \cdot \sum_{(x,y) \in Z^-} f^{\text{ME}}(x). \quad (\text{A7})$$

741 The prototypical representations r^+ , $r^- \in \mathbb{R}^d$ and the query molecule embedding $m \in \mathbb{R}^d$ are then
 742 used to make the final prediction:

$$\hat{y} = \frac{\exp(-d(m, r^+))}{\exp(-d(m, r^+)) + \exp(-d(m, r^-))}, \quad (\text{A8})$$

Table A2: Hyperparameter space considered for the Neural Search model selection. The hyperparameters of the best configuration are marked bold.

Hyperparameter	Explored values
Number of hidden layers	1, 2 , 4
Number of units per hidden layer	1024 , 4096
Output dimension	512 , 1024
Activation function	ReLU, SELU
Learning rate	0.0001, 0.001 , 0.01
Optimizer	Adam
Weight decay	0 , $1 \cdot 10^{-4}$
Batch size	4096
Input Dropout	0.1
Dropout	0.5
Layer-normalization	False, True
• Affine	False
Similarity function	cosine similarity, dot product , MinMax similarity

743 where d is a distance metric.

744 **Hyperparameter search.** Hyperparameter search has been done in Stanley et al. (2021), to which
 745 we refer here. ECFP fingerprints and descriptors created by a GNN, which operates on the molecular
 746 graph, are fed into a fully connected neural network, which maps the input into an embedding space
 747 with the dimension of 512. (Stanley et al., 2021) use the Mahalanobis distance to measure the
 748 similarity between a query molecule and the prototypical representations in the embedding space.
 749 The learning rate is 0.001 and the batch size is 256. The implementation can be found here https://github.com/microsoft/FS-Mol/blob/main/fs_mol/protonet_train.py and important
 750 hyperparameters are chosen here https://github.com/microsoft/FS-Mol/blob/main/fs_mol/utis/protonet_utis.py.

753 **Connection to Siamese networks and contrastive learning with InfoNCE.** If instead of the neg-
 754 ative distance $-d(\cdot, \cdot)$ the dot product similarity measure with appropriate scaling is used, ProtoNet
 755 for two classes becomes equivalent to Siamese Networks. Note that in our study, another differ-
 756 ence is that ProtoNet uses a GNN as encoder, whereas Siamese Networks use a descriptor-based
 757 fully-connected network as encoder. In case of dot product as similarity measure, the objective also
 758 becomes equivalent to contrastive learning with the InfoNCE objective (Oord et al., 2018).

759 A.1.5 IterRefLSTM: details and hyperparameters

760 (Altae-Tran et al., 2017) modified the idea of Matching Networks (Vinyals et al., 2016) by replacing
 761 the LSTM with their Iterative Refinement Long Short-Term Memory (IterRefLSTM). The use of the
 762 IterRefLSTM empowers the architecture to update not only the embeddings for the input molecule
 763 but also adjust the representations of the support set molecules.

764 For IterRefLSTM, $\mathbf{m} = f_{\theta_1}^{\text{ME}}(m)$ and $\mathbf{x}_n = f_{\theta_2}^{\text{ME}}(x_n)$ are two potentially different molecule en-
 765 coders for input molecule m and the support set molecules x_1, \dots, x_N . The next step in IterRefLSTM
 766 is:

$$[\mathbf{m}', \mathbf{X}'] = \text{IterRefLSTM}_L([\mathbf{m}, \mathbf{X}]).$$

767 Here, \mathbf{m}' and \mathbf{X}' contain the updated representations for the query molecule and the support
 768 set molecules. The IterRefLSTM denotes the function which updates these representations. The
 769 main property of the IterRefLSTM module is that it is permutation-equivariant, thus a permu-
 770 tation $\pi(\cdot)$ of the input elements results in the permutation of output elements: $\pi([\mathbf{m}', \mathbf{X}']) =$
 771 $\text{IterRefLSTM}_L(\pi([\mathbf{m}, \mathbf{X}]))$. The full architecture is invariant to permutations of the support set
 772 elements. For details, we refer to (Altae-Tran et al., 2017). The hyperparameter $L \in \mathbb{N}$ controls the
 773 number of iteration steps of the IterRefLSTM.

Table A3: Hyperparameter space considered for the IterRef model selection. The hyperparameters of the best configuration are marked bold.

Hyperparameter	Explored values
Molecule encoder	
• Number of hidden layers	0 , 1, 2, 4
• Number of units per hidden layer	1024 , 4096
• Output dimension	512 , 1024
• Activation function	ReLU, SELU
• Input dropout	0.1
• Dropout	0.5
IterRef embedding layer	
• L	1, 3
Similarity module:	
• Metric	cosine similarity, dot product , MinMax similarity
• Similarity space dimension	512, 1024
Layer-normalization	False, True
• Affine	False , True
Training	
• Learning rate	0.0001, 0.001 , 0.01
• Optimizer	Adam , AdamW
• Weight decay	0 , 0.0001
• Batch size	2048 , 4096

774 As similarity module, the IterRefLSTM uses the following:

$$\mathbf{a} = \text{softmax}(\mathbf{k}(\mathbf{m}', \mathbf{X}'))$$

$$\hat{y} = \sum_{n=1}^N a_n y_n,$$

775 where \hat{y} is the prediction for the query molecule. For the computation of the attention values \mathbf{a} , the
776 softmax function is used. \mathbf{k} is a similarity metric, such as the cosine similarity.

777 **Hyperparameter search.** All hyperparameters were selected based on manual tuning on the
778 validation set. We report the explored hyperparameter space in Table A3. Bold values indicate the
779 selected hyperparameters for the final model.

780 A.1.6 MHNfs: details and hyperparameters

781 The MHNfs consists of a molecule encoder, the context module, the cross-attention-module, and the
782 similarity module. The molecule encoder is a fully-connected Neural Network, consisting of one
783 layer with 1024 units. For the context module, a Hopfield layer with 8 heads is used and also the cross-
784 attention module include 8 heads. We chose a concatenation of ECFPs and RDKit-based descriptors
785 as the inputs for the MHNfs model. Notably, the RDKit-based descriptors were pre-processed in a
786 way that instead of raw values quantils, which were computed by comparing a raw value with the
787 distribution of all FS-Mol training molecules, were used. All descriptors were normalized based on
788 the FS-Mol training data.

789 **Hyperparameter search.** All hyperparameters were selected based on manual tuning on the
790 validation set. We report the explored hyperparameter space in Table A4. Bold values indicate the
791 selected hyperparameters for the final model. Early stopping points for the different re-runs are
792 chosen based on the $\Delta\text{AUC-PR}$ metric on the validation set. For the five re-runs the early-stopping
793 points, that were automatically chosen by their validation metrics, were the checkpoints at epoch 94,
794 192, 253, 253 and 309.

795 **Model training.** Figure A2 shows the learning curve of an exemplary training run of a MHNfs
796 model on FS-Mol. The left plot shows the loss on the training set and the right plot shows the

Table A4: Hyperparameter space considered for the MHNfs model selection. The hyperparameters of the best configuration are marked bold.

Hyperparameter	Explored values
Molecule encoder	
• Number of hidden layers	0 , 1, 2, 4
• Number of units per hidden layer	1024 , 4096
• Output dimension	512 , 1024
• Activation function	ReLU, SELU
• Input dropout	0.1
• Dropout	0.5
Context module (hopfield layer)	
• Heads	8 , 16
• Association space dimension	512 [512;2048]
• τ	22.6 [15;40]
• Dropout	0.1, 0.5
Cross-attention module (transformer mechanism)	
• Heads	1, 8 , 10, 16, 32, 64
• Number units in the hidden feedforward layer	567 [512; 4096]
• Association space dimension	1088 [512;2048]
• Dropout	0.1, 0.5 , 0.6, 0.7
• Number of layers:	1 , 2, 3
Similarity module:	
• Metric	cosine similarity, dot product , MinMax similarity
• Similarity space dimension	512, 1024
Layer-normalization	False, True
• Affine	False , True
Training	
• Learning rate	0.0001 , 0.001, 0.01
• Optimizer	Adam , AdamW
• Weight decay	0 , 0.0001
• Batch size	4096
• Warm-up phase (epochs)	5
• Constant learning rate phase (epochs)	25, 35
• Decay rate	0.994
• Max. number of epochs	350

797 validation loss. The dashed line indicates the checkpoint of the model which was saved and then used
798 for inference on the test set, whereas the stopping point was evaluated maximizing the Δ AUC-PR
799 metric on the validation set.

800 **Performance improvements in comparison to a naive baseline.** Figure A3 shows a task-wise
801 performance comparison between MHNfs and the frequent hitter model. Each point indicates a task
802 in the test set and is colored according to their super-class membership. In 132 cases the MHNfs
803 outperforms the frequent hitter model. In 25 cases the frequent hitter model yields better performance.

804 A.1.7 PAR: details and hyperparameters

805 The PAR model (Wang et al., 2021) includes a pre-trained GNN encoder, which creates initial
806 embeddings for the query and support set molecules. These embeddings are fed into an attention
807 mechanism module which also uses activity information of the support set molecules to create
808 enriched representations. Another GNN learns relations between query and support set molecules.

809 **Hyperparameter search.** For details we refer to (Wang et al., 2021) and <https://github.com/tata1661/PAR-NeurIPS21/blob/main/parser.py>. All hyperparameters were selected based
810 on manual tuning on the validation set. The hyperparameter choice for Tox21 (Wang et al., 2021)
811 was used as a starting point. We report the explored hyperparameter space in Table A5. Bold values
812

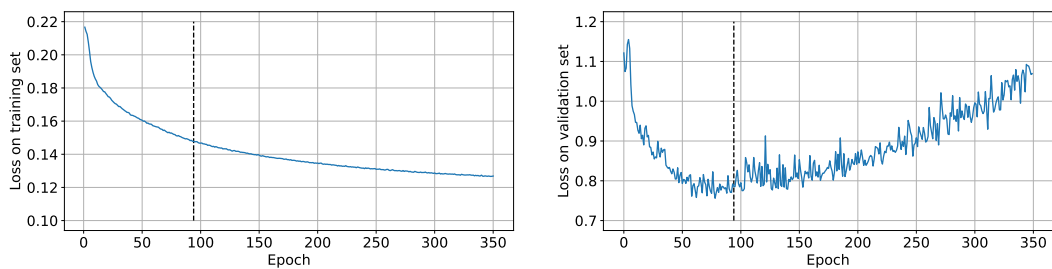


Figure A2: Exemplary MHNfs learning curve on FS-Mol. On the x-axis the number of epochs is displayed and on the y-axis the training loss (left) and the validation loss (right). The dashed line indicates the determined early-stopping point which is determined based on $\Delta\text{AUC-PR}$ on the validation set.

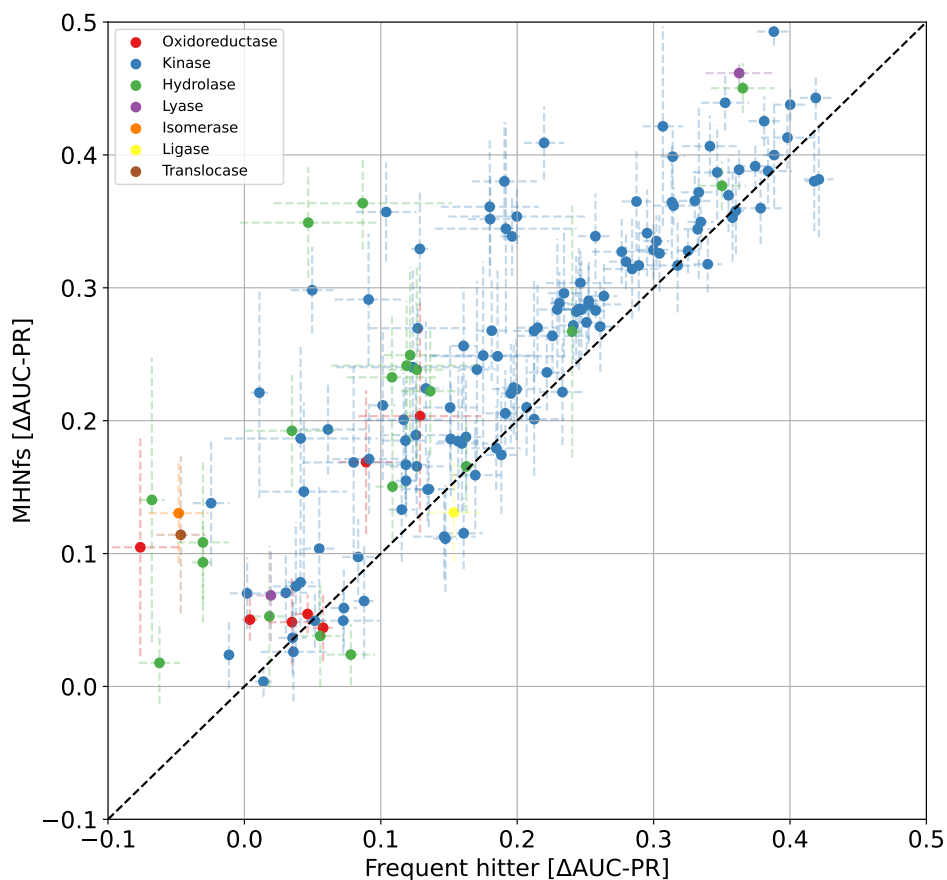


Figure A3: Performance comparison of MHNfs with the frequent hitter model. Each point refers to a task in the test set. Dashed lines indicate variability across training re-runs and different test support sets. The most points are located above the dashed line, which indicates that MHNfs performs better than the FH baseline at this task.

Table A5: Hyperparameter space considered for the PAR model selection. The hyperparameters of the best configuration are marked bold.

Hyperparameter	Explored values
Training	
• Meta learning rate	$1.0 \cdot 10^{-05}$, $1.0 \cdot 10^{-04}$, $1.0 \cdot 10^{-03}$, $1.0 \cdot 10^{-02}$
• Inner learning rate	0.01, 0.1
• Update step	1 , 2
• Update step test	1 , 2
• Weight decay	$5.0 \cdot 10^{-05}$, $1.0 \cdot 10^{-03}$
• Epochs	200000
• Eval. steps	2000
Encoder	
• Use pre-trained GNN	yes , no
Attention-based module	
• Map dimension	128, 512
• Map layer	2 , 3
• Pre fc layer	0 , 2
• Map dropout	0.1 , 0.5
• Context layer	2 , 3, 4
Relation graph	
• Hidden dimension	8, 128, 512
• Number of layers	2, 4
• Number of layers for relation edge update	2, 3
• Batch norm	yes, no
• Relation dropout 1	0 , 0.25, 0.5
• Relation dropout 2	0.2 , 0.25, 0.5

813 indicate the selected hyperparameters for the final model. Notably, we just report hyperparameter
814 choices which were different from standard choices. We used a training script provided by (Wang
815 et al., 2021), which can be found here <https://github.com/tata1661/PAR-NeurIPS21>.

816 A.2 Domain shift experiment

817 **Experimental setup.** For the domain shift experiment, we used the Tox21 dataset. This dataset
818 consists of 12,707 chemical compounds, for which measurements for up to 12 different toxic effects
819 are reported (Mayr et al., 2016; Huang et al., 2016a). It was published with a fixed training, validation
820 and test split. State-of-the-art supervised learning methods that have access to the full training set
821 reach AUC performance values between 0.845 and 0.871 (Klambauer et al., 2017; Duvenaud et al.,
822 2015; Li et al., 2017, 2021; Zaslavskiy et al., 2019; Alperstein et al., 2019). For our evaluation, we
823 re-cast Tox21 as a few-shot learning setting and draw small support sets from the 12 tasks. The
824 compared methods were pre-trained on FS-Mol and obtain small support sets from Tox21. Based
825 on the support sets, the methods had to predict the activities of the Tox21 test set. Note that there is
826 a strong domain shift from drug-like molecules of FS-Mol to environmental chemicals, pesticides,
827 food additives of Tox21. The domain shift also concerns the outputs where a shift from kinases,
828 hydrolases, and oxidoreductases of FS-Mol to nuclear receptors and stress responses of Tox21 is
829 present.

830 **Methods compared.** We compared the new method **MHNfs**, the runner-up method **IterRefLSTM**,
831 and **Similarity Search** — since it has been widely used for such purposes for decades (Cereto-
832 Massagué et al., 2015).

833 **Training and evaluation.** We followed the procedure of Stanley et al. (2021) for data-cleaning,
834 preprocessing and extraction of the fingerprints and descriptors used in FS-Mol. After running the
835 cleanup step, 8,423 molecules remained for the domain shift experiments. From the training set, 8
836 active and 8 inactive molecules per task were randomly selected to build the support set. The test set
837 molecules were used as query molecules. The validation set molecules were not used at all. During
838 test-time, a support set was drawn ten times for each task. Then, the performance of the models were

Table A6: Results of the domain shift experiment on Tox21 [AUC, Δ AUC-PR]. The best method is marked bold. Error bars represent standard deviation across training re-runs and draws of support sets

Method	AUC	Δ AUC-PR
Similarity Search (baseline)	.629 \pm .015	.061 \pm .008
IterRefLSTM (Altae-Tran et al., 2017)	.664 \pm .018	.067 \pm .008
MHNfs (ours)	.679 \pm .018	.073 \pm .008

Table A7: Results of the ablation study on FS-Mol [AUC, Δ AUC-PR]. The error bars represent standard deviation across training re-runs and draws of support sets. The p -values indicate whether the difference between two models in consecutive rows is significant.

Method	AUC	Δ AUC-PR	$p_{\text{AUC}}^{\text{a}}$	$p_{\Delta\text{AUC-PR}}^{\text{a}}$
MHNfs (CM+CAM+SM)	.739 \pm .005	.241 \pm .006		
MHNfs -CM	.737 \pm .004	.240 \pm .005	0.030	0.002
MHNfs -CM -CAM	.719 \pm .006	.223 \pm .006	< 1.0e-8	< 1.0e-8
Similarity Search	.604 \pm .003	.113 \pm .004	< 1.0e-8	< 1.0e-8
IterRefLSTM (Altae-Tran et al., 2017) ^b	.730 \pm .005	.234 \pm .005	< 1.0e-8	8.73e-7

^a paired Wilcoxon rank sum test ^b IterRefLSTM is compared to MHNfs -CM

839 evaluated for these support sets, using the area under precision-recall curve (AUC-PR), analogously to
 840 the FS-Mol benchmarking experiment reported as the difference to a random classifier (Δ AUC-PR),
 841 and the area under receiver operating characteristic curve (AUC) metrics. The performance values
 842 report the mean over all combinations regarding the training reruns and the support set sampling
 843 iterations. Error bars indicate the standard deviation.

844 **Results.** The Hopfield-based context retrieval method has significantly outperformed the
 845 IterRefLSTM-based model ($p_{\Delta\text{AUC-PR}}$ -value $3.4\text{e-}5$, p_{AUC} -value $2.5\text{e-}6$, paired Wilcoxon test)
 846 and the Classic Similarity Search ($p_{\Delta\text{AUC-PR}}$ -value $2.4\text{e-}9$, p_{AUC} -value $7.6\text{e-}10$, paired Wilcoxon
 847 test) and therefore showed robust performance on the toxicity domain, see Table A6. Notably, all
 848 models were trained on the FS-Mol dataset and then applied to the Tox21 dataset without adjusting
 849 any weight parameter.

850 A.3 Details on the ablation study

851 The MHNfs has two new main elements compared to the previous state-of-the art method Iter-
 852 RefLSTM, which are the context module and the cross-attention-module. In this ablation study
 853 we aim to investigate i) the importance of all design elements, which are the context module, the
 854 cross-attention module, and the similarity module, and ii) the superiority of the cross-attention module
 855 compared to the IterRefLSTM module.

856 A.3.1 Ablation study A: comparison against IterRefLSTM

857 For a fair comparison between the cross-attention module and the IterRefLSTM we used a pruned
 858 MHN version ("MHNfs -CM") which has no context module and compared it with the IterRefLSTM
 859 model. The evaluation includes five training re-runs each and ten different support set samplings.
 860 The results, reported as the mean across training re-runs and support sets, can be seen in Table A7.
 861 We performed a paired Wilcoxon rank sum test for both the AUC and the Δ AUC-PR metric. Both
 862 p -values indicate high significance.

863 A.3.2 Ablation study B: all design elements

864 We evaluate the performance of all main elements within the MHNfs, which are the context module,
 865 the cross-attention module, the similarity module and the molecule encoder. For this analysis,
 866 we start with the complete MHNfs which includes all modules and report AUC and Δ AUC-PR
 867 performance values. Then, we iteratively omit the individual modules, measuring whether there is a

Table A8: Results of the ablation study on Tox21 [AUC, Δ AUC-PR]. The error bars represent standard deviation across training re-runs and draws of support sets. The p -values indicate whether a model is significantly different to the MHNfs in terms of the AUC and Δ AUC-PR metric.

Method	AUC	Δ AUC-PR	$p_{\text{AUC}}^{\text{a}}$	$p_{\Delta\text{AUC-PR}}^{\text{a}}$
MHNfs (CM+CAM+SM)	.679 \pm .018	.073 \pm .008		
MHNfs -CM	.662 \pm .028	.069 \pm .012	6.28e-8	0.002
MHNfs -CM -CAM	.640 \pm .018	.057 \pm .009	<1.0e-8	<1.0e-8
Similarity Search	.629 \pm .015	.061 \pm .008	<1.0e-8	<1.0e-8
IterRefLSTM	.664 \pm .018	.067 \pm .008	2.53e-6	3.38e-5

^a paired Wilcoxon rank sum test

868 significant performance difference with and without the module. Table A7 shows the results, where
 869 performance values for the full MHNfs, a MHNfs model without the context module ("MHNfs -CM")
 870 and a MHNfs module without the context and the cross-attention module ("MHNfs -CM -CAM") is
 871 included. Notably, the model without the context module and without the cross-attention module
 872 just consists of a learned molecule encoder and the similarity module. We evaluated the impact of
 873 the learned molecule encoder by replacing it with a fixed encoder, which maps a molecule to its
 874 descriptors. The model with the fixed encoder is a classic chemoinformatics method which is called
 875 Similarity Search (Cereto-Massagué et al., 2015).

876 For the evaluation, we performed five training re-runs for every model and sampled ten different
 877 support sets for every task. Table A7 shows the results in terms of AUC and Δ AUC-PR. We performed
 878 paired Wilcoxon rank sum tests on both metrics, comparing two methods in consecutive rows in the
 879 table. The table shows that every module has a significant impact as omitting a module results in
 880 a significant performance drop. The comparison between the MHNfs version without the context
 881 module and without the cross-attention module with the Similarity Search showed a significant
 882 superiority of the learned molecule encoder in comparison to the fixed encoder.

883 A.3.3 Ablation study C: Under domain shift on Tox21

884 Referring to Section A.3.2, the context module and the cross-attention module showed their impor-
 885 tance for the global architecture. This importance gets even more pronounced for the domain shift
 886 experiment on Tox21 as one can see in Table A8.

887 Again, five training re-runs and ten support set draws are used for evaluation. Including the context
 888 module makes a clear and significant difference for both metrics AUC and Δ AUC-PR.

889 A.4 Generalization to different support set sizes

890 In this section, we test the ability of MHNfs to generalize to different support set sizes. During
 891 training in the FS-Mol benchmarking setting, the MHNfs model has access to support sets of size
 892 16. However, at inference, the support set size might be different. Figure A4 provides performance
 893 estimates of the support-set-size-16 MHNfs models on other support set sizes. Note that the estimates
 894 could be seen as approximate lower bounds of the predictive performance on settings with different
 895 support set sizes (y-axis labels). For a model used in production or in a real-world drug discovery
 896 setting, MHNfs should be trained with varying support set sizes that resemble the distribution of real
 897 drug discovery projects.

898 A.5 Generalization to different context sets

899 In this section, we test the ability of MHNfs to generalize to different context sets. While the FS-Mol
 900 training split is used as a context during training, we assessed whether our model is robust to different
 901 context sets for inference. To this end we preprocessed the GEOM dataset (Axelrod and Gomez-
 902 Bombarelli, 2022) from which we used 100,000 molecules that passed all pre-processing checks.
 903 From this set, we sample 10,000 molecules as context set for MHNfs. Because GEOM contains
 904 drug-like molecules, similar to FS-Mol the predictive performance remains stable (see Table A9).

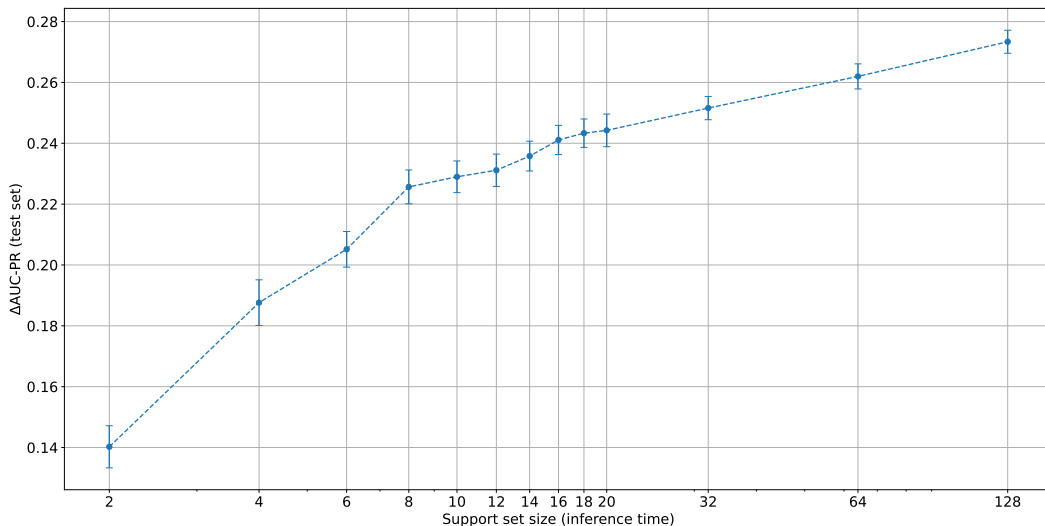


Figure A4: Performance of MHNfs for different support set sizes during inference time. The MHNfs models are trained with support sets of the size 16.

Table A9: MHNfs performance for different context sets [Δ AUC-PR]. The error bars represent standard deviation across training re-runs and draws of support sets.

Dataset used as a context	Δ AUC-PR
FS-Mol (Stanley et al., 2021)	.2414 \pm .006
GEOM (Axelrod and Gomez-Bombarelli, 2022)	.2415 \pm .005

905 A.6 Details and insights on the context module

906 The context module replaces the initial representations of query and support set molecules by a
 907 retrieval from the context set. The context set is a large set of molecules and covers a large chemical
 908 space. The context module learns how to replace the initial molecule embeddings such that the
 909 context-enriched representations are put in relation to this large chemical space and still contains
 910 all necessary information for the similarity-based prediction part. Figure A5 shows the effect of the
 911 context module for the MHNfs model. Extreme initial embeddings, such as the purple embedding
 912 on the right, are pulled more into the known chemical space, represented by the context molecules.
 913 Notably, the replacement described above is a soft replacement, because also the initial embeddings
 914 contribute to the context-enriched representations due to skip-connections.

915 A.7 Reinforcing the covariance structure in the data using modern Hopfield networks

916 We follow the argumentation of (Fürst et al., 2021, Theorem A3) that retrieval from an associative
 917 memory of a MHN reinforces the covariance structure.

918 Let us assume that we have one molecule embedding from the query set $\mathbf{m} \in \mathbb{R}^d$ and one molecule
 919 embedding from the support set $\mathbf{x} \in \mathbb{R}^d$ and both have been enriched with the context module with
 920 memory $\mathbf{C} \in \mathbb{R}^{d \times M}$ (ignoring linear mappings):

$$\mathbf{m}' = \mathbf{C} \operatorname{softmax}(\beta \mathbf{C}^T \mathbf{m}) \quad (\text{A9})$$

$$\mathbf{x}' = \mathbf{C} \operatorname{softmax}(\beta \mathbf{C}^T \mathbf{x}) \quad (\text{A10})$$

921 Then the similarity of the retrieved representations as measured by the dot product can be expressed
 922 in terms of covariances:

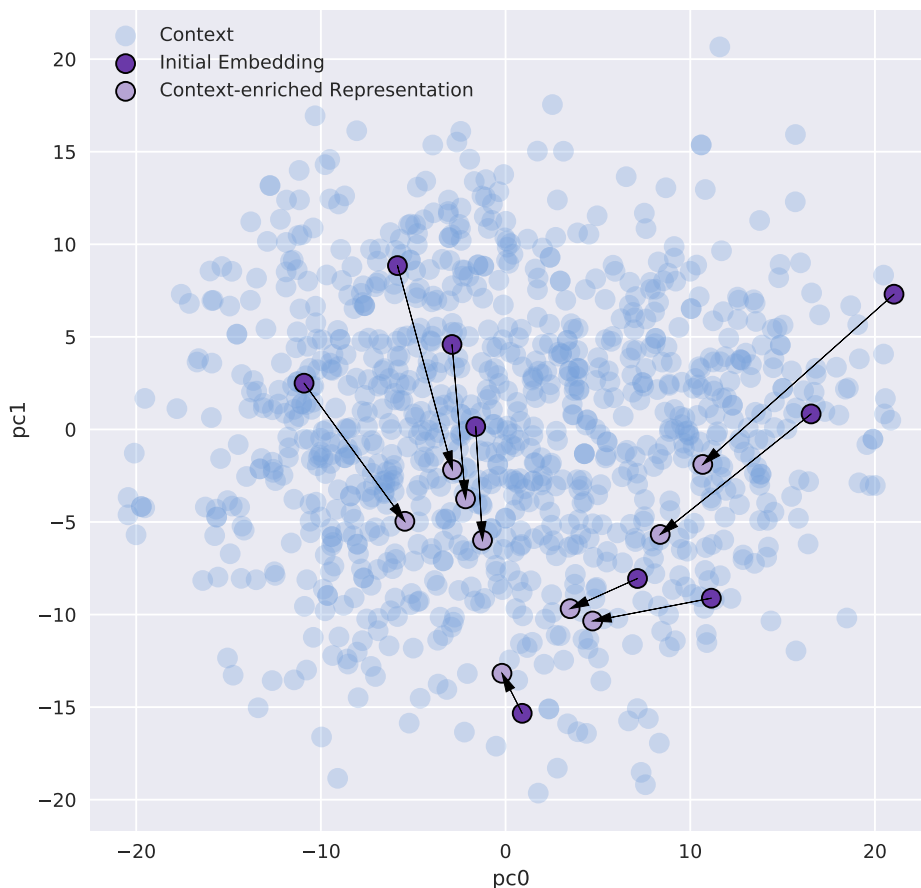


Figure A5: PCA plot of molecule embeddings. Each dot in the plot represents a molecule embedding, of which the first two principal components are displayed on the x- and y-axis. Blue dots represent context molecules. Dark purple dots represent initial embeddings for some exemplary molecules, of which some exhibit extreme characteristics and are thus located away from the center. Arrows and light purple dots represent the enriched molecule embeddings after the retrieval step. Especially molecules from extreme positions are moved stronger to the center and thus are more similar to known molecules after retrieval.

$$\mathbf{m}'^T \mathbf{x}' = \text{softmax}(\beta \mathbf{C}^T \mathbf{m})^T \mathbf{C}^T \mathbf{C} \text{softmax}(\beta \mathbf{C}^T \mathbf{x}) = \quad (\text{A11})$$

$$= (\bar{\mathbf{c}} + \text{Cov}(\mathbf{C}, \mathbf{m})^T \mathbf{m})^T (\bar{\mathbf{c}} + \text{Cov}(\mathbf{C}, \mathbf{x}) \mathbf{x}), \quad (\text{A12})$$

923 where $\bar{\mathbf{c}}$ is the row mean of \mathbf{C} and following the *weighted covariances* are used:

$$\text{Cov}(\mathbf{C}, \mathbf{m}) = \mathbf{C} \mathbf{J}^m(\beta \mathbf{C} \mathbf{m}) \mathbf{C}^T \quad \text{Cov}(\mathbf{C}, \mathbf{x}) = \mathbf{C} \mathbf{J}^m(\beta \mathbf{C} \mathbf{x}) \mathbf{C}^T. \quad (\text{A13})$$

924 $\mathbf{J}^m : \mathbb{R}^M \mapsto \mathbb{R}^{M \times M}$ is a mean Jacobian function of the softmax (Fürst et al., 2021, Eq.(A172)).

925 The Jacobian \mathbf{J} of $\mathbf{p} = \text{softmax}(\beta \mathbf{a})$ is $\mathbf{J}(\beta \mathbf{a}) = \beta (\text{diag}(\mathbf{p}) - \mathbf{p} \mathbf{p}^T)$.

$$\mathbf{b}^T \mathbf{J}(\beta \mathbf{a}) \mathbf{b} = \beta \mathbf{b}^T (\text{diag}(\mathbf{p}) - \mathbf{p} \mathbf{p}^T) \mathbf{b} = \beta \left(\sum_i p_i b_i^2 - \left(\sum_i p_i b_i \right)^2 \right), \quad (\text{A14})$$

926 this is the second moment minus the mean squared, which is the variance. Therefore, $\mathbf{b}^T \mathbf{J}(\beta \mathbf{a}) \mathbf{b}$ is
 927 times the covariance of \mathbf{b} if component i is drawn with probability p_i of the multinomial distribution
 928 \mathbf{p} . In our case the component i is context sample \mathbf{c}_i . \mathbf{J}^m is the average of $\mathbf{J}(\lambda \mathbf{a})$ over $\lambda = 0$ to $\lambda = \beta$.

929 Note that we can express the enriched representations using these covariance functions:

$$\mathbf{m}' = (\bar{\mathbf{c}} + \text{Cov}(\mathbf{C}, \mathbf{m})^T \mathbf{m}) \quad (\text{A15})$$

$$\mathbf{x}' = (\bar{\mathbf{c}} + \text{Cov}(\mathbf{C}, \mathbf{x})^T \mathbf{x}), \quad (\text{A16})$$

930 which connects retrieval from MHNs with reinforcing the covariance structure of the data.