

Cite this: DOI: 00.0000/xxxxxxxxxx

"Amide – amine + alcohol = carboxylic acid." Chemical reactions as linear algebraic analogies in graph neural networks.[†]

Amer Marwan El-Samman,^{*a} and Stijn De Baerdemacker^{a,b}

Received Date

Accepted Date

DOI: 00.0000/xxxxxxxxxx

In deep learning methods, especially in the context of chemistry, there is an increasing urgency to uncover the hidden learning mechanisms often dubbed as "black box." In this work, we show that graph models built on computational chemical data behave similar to natural language processing (NLP) models built on text data. Crucially, we show that atom-embeddings, a.k.a atom-parsed graph neural activation patterns, exhibit arithmetic properties that represent valid reaction formulas. This is very similar to how word-embeddings can be combined to make word analogies, thus preserving the semantic meaning behind the words, as in the famous example "King" – "Man" + "Woman" = "Queen." For instance, we show how the reaction from an alcohol to a carbonyl is represented by a constant vector in the embedding space, implicitly representing " H_2 ," independent from the particular carbonyl reactant and alcohol product. This reveals a highly-structured vector space, wherein the directions in the embedding space are synonymous with chemical changes (ex. the oxidation direction), and distinct chemical changes are orthogonal. In contrast to natural language processing, we can explain the observed chemical analogies using algebraic manipulations on the local chemical composition that surrounds each atom-embedding. Furthermore, the observations find applications in transfer learning, for instance in the formal structure and prediction of atomistic properties, such as $^1\text{H-NMR}$ and $^{13}\text{C-NMR}$. This work is in line with the recent push for interpretable explanations to graph neural network modeling of chemistry and uncovers a latent model of chemistry that is highly structured, consistent, and analogous to chemical syntax.

Neural networks in chemistry have gained tremendous traction in the past decade, carrying a broad range of applicability, from aiding drug and material discovery,^{1–11} to speeding up or bypassing the prediction of electronic structure properties.^{12–33} For each application, numerous approaches have been designed, including both graph-^{33–36} and text-based models^{3,37–44} where such techniques enjoy a varying degree of success.

In the context of text-based models, sophisticated text-based chemical inputs have been developed such as the Simplified Molecular Input Line Entry System (SMILES),^{43–46} the Self-Referencing Embedded Strings (SELFIES),^{47,48} and SMILES arbitrary-target specification (SMARTS).^{49,50} Though convolutional neural networks have been tried successfully on SMILES-based text for the detection of chemical motifs and on prediction of drug activity,^{39,43} the majority of so-called "linear graph models" have been fitted using recurrent neural networks (RNNs) due to their capacity of holding short- and long-term information about text.^{38,40,41,44} For example, SMILES2Vec is a deep RNN that learns important features from SMILES strings to predict toxicity, activity, solubility and solvation energy of chemical compounds.⁴⁴ Text-based models also facilitate the design of generative architectures that predict the result of chemical reactions, generating the product strings from a reactant string input, or generating molecules of a desired property.^{37,38,40}

Graph models,^{17–21,23,26–36} on the other hand, are based on the raw visual representation of a molecule as a collection of atoms in three-dimensional coordinate space. Generally speaking, the coordinates of the atoms serve as inputs to such models and the output is the target chemical property under investigation, often on energy. Permutational- and symmetry- invariant graph models have been designed successfully for the prediction of electronic energy within chemical accuracy.^{17–21,23} Notably, such graph-based methods in chemistry share important properties with text-based recurrent and transformer models making it instructive to examine the connection between these seemingly different approaches.

The central object shared in all these approaches is the "embedding" which is the feature-building quantity of the neural network

^a University of New Brunswick, Department of Chemistry. 30 Dineen Dr, Fredericton, Canada. E-mail: aelsamma@unb.ca

^b University of New Brunswick, Department of Mathematics and Statistics. 30 Dineen Dr, Fredericton, Canada. E-mail: stijn.debaerdemacker@unb.ca

in the latent space. For example, in text-based models, word-embeddings represent the latent features for each word in the context of the sentence after training.⁵¹⁻⁶⁰ In chemistry, atom-embeddings, also optimized through training, are the features representing an atom in the context of a molecular graph.^{17-21,23} These neural activations for the word/atom are tuned to hold meaningful information about the context of the data (i.e. words, molecules) and the specific input. However, they are generally high-dimensional and obscure objects to analyze on their own. In previous works, we showed that the embeddings of chemical GNN models hold valuable information about chemistry such as the ability to distinguish molecular environments and the ability to quantify molecular similarity.⁶¹ We also showed that the chemical embedding space is a readily transferable representation for a wide array of properties such as for pKa, NMR, and solubility, underscoring the completeness of these representations.⁶² In this work, we go beyond the locality of the representation, and show that graph-embeddings behave similar to text-embeddings in that they have arithmetic properties that reveal meaningful combinations, akin to how word-embeddings can be combined to make word analogies.⁵¹⁻⁵⁵ This will naturally uncover the chemical syntactical organization of the embedding space.

The surprising property of word analogies using vector arithmetic has been observed first in natural language models.⁵¹⁻⁵⁵ In trained models such as skip-gram with negative sampling (SGNS), the word analogy "King is to X as Man is to Woman?" is solved by taking the closest vector to "King – Man + Woman" which happens to be the vector for X = "Queen." The success of this is based loosely on the Pennington et al. conjecture,⁵⁵ proven in ref⁵¹, which states that such word analogies are linear iff $p(w|a)/p(w|b) \approx p(w|x)/p(w|y)$. In simple terms, if words a and b are found in the same ratios as x and y across all words w of a vocabulary, then there must be a linear analogy between, a , b , x , and y . We show in this work that the same observation can be made with graph-embeddings of molecular graphs.

For chemical statistical models this comes with promising consequences. Firstly, it means that fundamental stoichiometric reactions can be modelled with vector algebra thereby opening a new way to traverse the chemical space in an algebraically structured way. Most significantly, and in contrast to NLP, the structure of the chemical embedding space is not a consequence of the social construct of language. Rather, it relates quantitatively to the underlying chemical structure. In other words, for a chemical GNN model, we can replicate the embedding space of chemistry by involving the notion of locality and chemical environments. We demonstrate how this can be achieved using perturbational updates that are based on the composition of consecutively layered neighborhood environments. Simultaneously, this leads to a natural interpretation of how atomistic and extensive properties are deduced from GNN models based on neighborhood composition.

This paper is organized as follows. In the Methodology Section, we first recapitulate the training of our pretrained GNN model on electronic energy,^{61,63} and our transfer learning models to other properties (¹H-NMR and ¹³C-NMR).⁶² Details of the hyperparameters chosen for each architecture and the dataset used for training can all be found in Section 1.1. Following this, in Section

1.2 we discuss how we prepare our reaction datasets from the QM9 dataset⁶⁴ using an algorithm that can query any class of reactants and transform them to products via a specified reaction. This dataset creation procedure will serve our observations on chemical reaction analogies in the embedding space which we present in Sections 2.1-2.3. Following our observations, we deduce a replicate model of the embedding space based on layered atomic neighborhood information in Section 2.4 which will explain our chemical analogy observations and reveal the global framework of the embedding space. Lastly, in Section 2.5 we show how linear analogies can reveal hidden relations in chemical properties such as ¹H-NMR and ¹³C-NMR.

1 Methodology

1.1 GNN, Pretraining and Hyperparameters

We employed a pretrained graph neural network, SchNet,¹⁷⁻²⁰ on electronic energy of the QM9 dataset.⁶⁴ QM9 is a set of 134K small-sized organic molecules ($\sim 5\text{-}10 \text{ \AA}$ in size) with optimized conformations all computed using the B3LYP/6-31G(d,p) level of density-functional theory. For the SchNet GNN model, the nodal features (i.e the atom-embeddings) were chosen to have a 128-dimensional latent feature space. The edges that update the nodal features employ an initial expansion of interatomic distances using equally separated Gaussians with a cutoff of 50 Å. This provides the edges with enough parametric flexibility to update the nodal features via convolutional operation described in ref¹⁷. More details on the GNN algorithm can also be found in refs¹⁷⁻²⁰. Note that more efficient GNN training algorithms employ cutoff distances that are significantly shorter which allow for efficiently learning the neighbor interactions. However, this avenue was not chosen since a large cutoff distance is purposeful to maintain a global representation of molecules in the embedding space. We trained on 100K molecules with total electronic energy at 0 Kelvin as the target property. An additional 10,000 data points were used for validation during the training process. The rest of the set (20,000) was leftover for testing. The trained model achieves a MAE of 0.2 meV and 1 meV on training and testing set's molecular energy, respectively. The trained model, details on all the hyperparameters, as well as the extracted embeddings for QM9 molecules can be found at⁶³.

For the prediction of ¹H-NMR and ¹³C-NMR, we employed a transfer learning model as introduced earlier in⁶². Such models help to transfer the integrity of learned chemical representation from GNN models to new molecular properties and datasets. The transfer learning architecture used is a simple feed forward neural network (made up of one layer of 128 nodes, followed by "tanh" activation, followed by another linear layer of 128 nodes), that intakes energy-trained embeddings (from SchNet) as inputs to transfer learn to other properties such as ¹H-NMR, and ¹³C-NMR. This follows very closely to the transfer learning architecture used in previous works,⁶² only here the activation function has been replaced from linear or "Relu", to "tanh."

For ¹H-NMR and ¹³C-NMR data, we used the QM9NMR dataset,⁶⁵ which has gas- and solvent-phase chemical shifts computed at the mPW1PW91/6-311+G(2d,p) for all QM9 molecules

at geometry optimized conformations computed at the B3LYP/6-31G(2df,p) level. We used the hydrogen- and carbon-site embeddings as the input representation for training on the prediction of gas-phase ^1H -NMR and ^{13}C -NMR. The model was trained using an 8000 molecule randomized dataset from QM9NMR, achieving a RMSE of 2.69 ppm for carbon NMR and 0.20 ppm for proton NMR on 2000 molecule separate test set, which is comparable to full-fledged training on significantly larger datasets. For instance, the highly accurate kernel regression model that was applied at the inception of the QM9NMR dataset⁶⁵ achieves a mean error of 1.9 ppm for carbon chemical shifts. Both carbon and proton NMR results are within the accuracy of density functional methods.

1.2 Reactions Dataset Creation

After training the model, we automated a dataset creation procedure mimicking the endpoints of reaction processes such as hydrolysis, diels-alder, and more. The procedure works as follows. First, we automate the identification of reactant functional groups across the entire QM9 using a specified reactant label. The automation procedure can annotate every atom's environment up to any compositional depth (one, two, to several bonds away) in a bijective labelling representation using atomic number priority to order neighboring atoms ensuring the uniqueness of the label. This is done in a way very similar to well-established standards of ordering R/S or E/Z nomenclature.⁶⁶ Ultimately, we obtain a local-centric label for each atom in the dataset that can be queried for long-range features, for instance, straight-chain alcohols, or alpha-positioned alkynes, or any other branching motif specified, or left unspecified (ex. all alcohols in QM9). Second, once reactants are isolated, a specified reaction is automated on the dataset. For instance, if the specified reaction is a methylation, this is carried out by removing a hydrogen and adding a methyl to mimic methylation on the alcohols. Lastly, we geometry optimize using RDKit's force field, MMFF94.⁶⁷

Following the preparation of the reactant and product databases, we then run the reactants and products through the pretrained GNN model to extract their atom-embeddings at the reaction-site. For instance, in hydrolysis, this would be the carbon on the carbonyl group. This is all automated in a python workflow visualised in Figure 1. The software for this procedure is open source and is referenced in the Supplementary Information section.

In order to visualize the embeddings as they change from reactants to products, we project the high-dimensional vectors of the embedding space to a lower-dimensional space while incurring minimal data loss. This can be done with Principal Component Analysis (PCA)⁶⁸ which finds the lower-dimensional space that packs the largest variance in the data.

Ultimately, visualization techniques fall short from giving a comprehensive quantitative grasp of the chemical syntax in the embedding space. We will quantify our observations of linear analogies by means of cosine similarities in high-dimensional spaces, which will provide a measure on how well two different vectors are mutually aligned (see next Section).

2 Results & Discussion

We investigate the chemical analogies in the embedding space with key reaction processes of increasing complexity that amount to all of the basic features of chemical reactions (adding/breaking bonds, one-step/multi-step). These reactions being 1) oxidation of alcohols/alkanes/alkenes, 2) diels-alder, and 3) hydrolysis of amides to carboxylic acids. While this set of reactions is definitely not exhaustive compared to the wide variety of reactions chemistry has to offer, it nevertheless is a balanced sample and will incur observations that can be easily extended to any reaction using the proposed methodology.

2.1 Oxidation Reactions

We start with the simplest of the listed reactions, a) the oxidation reaction of alkanes, alkenes and alcohols,



Figure 2 depicts the reaction process in the embedding space in the PCA projection for the oxidation of alkanes (2a), alkenes (2b), and alcohols (2c). In each, only the embedding vector of the reactant-center carbon C_r and product-center carbon C_p are depicted in color, and commented by an arrow representing the reaction. Note that the arrows only represent the reactant and product end-points and do not represent the full reaction process, for which QM9 is not designed, however there has been some recent effort in exploring the success of transfer learning to real-space quantities of chemistry that would include any part of the chemical space.⁶⁹ The scatterplot in the background of the arrows, consists of all the carbon embeddings of QM9 labelled according to chemical moiety as was shown in ref⁶¹. Other C embeddings in QM9 consisting of distinct functional groups from the reactant's functional group are shown in gray, whereas the colored groups (annotated in the caption) represent atoms in QM9 that have same functional groups as the reactants and products of the reaction.

The first thing to notice from Figure 2a is that the alkane and alkene carbons aggregate in different and well-separated clusters, as does every other carbon-centered chemical moiety in the database. This was already observed and discussed previously in^{23,61}. The second observation is that all the vectors of the transformation appear to be equal to a large extent. The extent to which the vectors are equal can be quantified by considering the average "oxidation vector from alkane to alkene" and use this as a proxy to transform the embedding of any reactant alkane to its alkene counterpart, that is,

$$\langle \text{alkane oxidation} \rangle = \frac{1}{N} \sum_{i=1}^N \left[x_{C_{\text{alkane}_i}} - x_{C_{\text{alkene}_i}} \right]. \quad (4)$$

In this approach, $\langle \text{alkane oxidation} \rangle$ can be used to estimate any $x_{C_{\text{alkene}_i}}$ from its $x_{C_{\text{alkane}_i}}$,

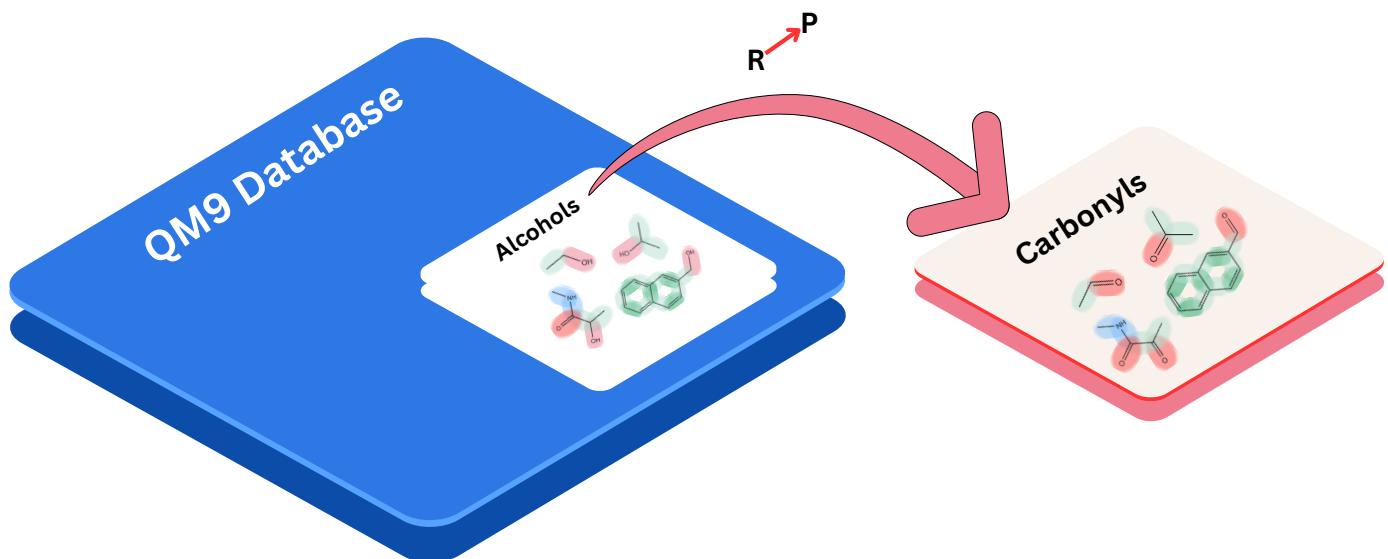


Fig. 1 Work flow for creating datasets of reactants and products. Starting with the QM9 dataset, we label all the atoms in the QM9 database according to the molecular environment surrounding them. Second, we isolate reactants of a reaction using the labelling system, in the illustration shown this is done for alcohol reactants. Third, the reactants are transformed to their product form, for instance, oxidation of alcohols, resulting in product carbonyls. Then lastly, we test the reactants and product through SchNet model to extract and calculate the embedding difference between reactant and product at the reaction site.

$$x_{C_{alkene_i}}^{\text{pred}} = x_{C_{alkane_i}} + \langle \text{alkane oxidation} \rangle. \quad (5)$$

We investigate the validity of this approximation for each reaction in the dataset by means of the "neighbor test," i.e. if the resulting $x_{C_{alkene_i}}$ of eq. (5) leads closest to the true alkene embedding or not within the total compounded set of all alkenes ($= 2192$, see Table 1). The neighbor test is analogous to the one carried out for the original word vector analogies where the word embedding for "King" – "Man" + "Woman" came out nearest to the vector for "Queen." Here, the average oxidation vector serves the same role as the vector for "– "Man" + "Woman" which transformed the word "King" to "Queen," and similarly "policeman" to "policewoman" and "boy" to "girl". In this case, the average oxidation vector can map 62.8% of reactant alkanes to be nearest their true product alkene C-embedding at the reaction center, see Table 1. This points towards a highly structured space, as more than 1 in 2 vectors are mapped exactly to the correct alkene using an average estimate reaction vector, opposed to a probability of 1 in N_{alkenes} ($1/2192$, see Table 1) if this would be random over a uniform distribution. In other words, the average $\langle \text{alkane oxidation} \rangle$ vector does not map to just any of the 2192 available alkenes, but lands exactly on the correct one 62.8% of the time. A noteworthy observation is that a 100% neighbor test can be achieved if the removal of hydrogen is performed without any subsequent force field optimization. Such discrepancy and the role of geometry optimization will be discussed later.

A different and more continuous measure to validate the performance of $\langle \text{alkane oxidation} \rangle$ is to compare the average distance between the predicted and the true value, $\frac{1}{N} \sum_i |x_{C_{alkene_i}}^{\text{pred}} - x_{C_{alkene_i}}| = 3.72$, to the average pairwise distance of all alkenes, $\frac{1}{\frac{N(N-1)}{2}} \sum_{i,j} |x_{C_{alkene_i}} - x_{C_{alkene_j}}| = 9.46$, or the average distance to the

nearest neighbor of each alkene, $\frac{1}{N} \sum_i |x_{C_{alkene_i}} - x_{C_{alkene_{\text{nearest}(i)}}}| = 1.79$. It is evident that we are well within the 'alkenes' feature space bordering on the exact subspace that the alkene product lies in. It is also worth mentioning that in contrast to language models, the analogies '– "Man" + "Woman"' have not been learned explicitly as "H₂" does not represent any single atomic embedding and is therefore categorically excluded from the dataset.

Similar results hold for the oxidation of alkenes to alkynes, and alkanes to alcohols, see Table 1 and Figures 2b and 2c. Visually, it also appears that the average oxidation vectors, across reaction classes, i.e for alkane, alkene, and alcohol oxidation in the PCA space appear to be largely colinear from the Figures. To confirm that this colinearity is not a coincidence of the 2D projection space, we can measure the cosine similarity between the average vectors in the original high-dimensional space, see Table 2. In the Table, the off-diagonal elements show the average cosine similarity between reaction classes (ex. between alkane oxidation and alkene oxidation) taken by considering the cosine similarity between the average reaction vectors of those classes,

$$\cos \Theta_{ij} = \frac{\langle \text{reaction}_i \rangle \cdot \langle \text{reaction}_j \rangle}{\sqrt{|\langle \text{reaction}_i \rangle|^2 |\langle \text{reaction}_j \rangle|^2}}. \quad (6)$$

It is apparent from the Table that all oxidations share a high degree of cosine similarity, especially when compared to the other reactions studied. This is indeed significant, considering that in a high-dimensional space (128-D) it is increasingly likely that any two random vectors are orthogonal.⁷⁰ It can be shown that the dot product of normally distributed vectors in D dimensions are strongly centered at $\cos(\theta) = 0$, with a standard deviation of $\sigma = 1/D$. For our 128-dimensional latent space, deviation from orthogonality of a normal distribution is 0.008. Thus the co-

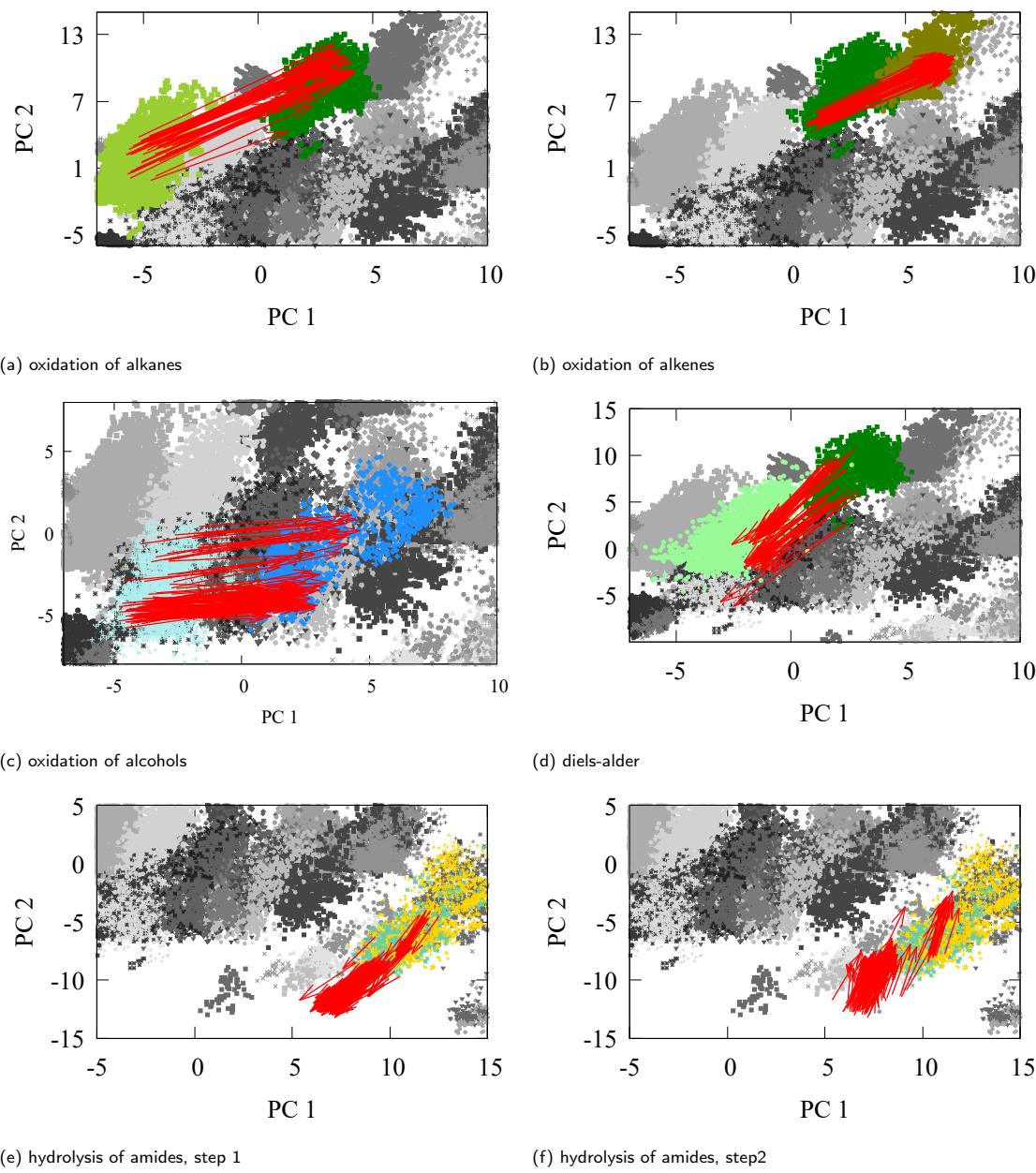


Fig. 2 Transformation vectors from reactant to product embedding for a) alkanes, b) alkenes, c) alcohols, d) diels-alder reaction, e) hydrolysis of amides step 1, f) hydrolysis of amides step 2, at the reaction center after geometry optimizing the products with MMFF94 force field. The scatterplot around the arrows comes from QM9's carbon embeddings which naturally separate based on functional groups and have been greyed out except for the embeddings that resemble the reaction center embeddings for reactant and product, ex. all other alkanes and alkenes in QM9. The colors represent the functional groups for **alkanes**, **alkenes**, **alkynes**, **alcohols**, **carbonyls**, **methines** of the diels-alder product, **amides**, and **carboxylic acid**.

reaction	QM9	new	total	uniform density	% match
$\text{RCH}_2\text{CH}_2\text{R} \rightarrow \text{RCHCHR}$	2062	118	2192	5×10^{-4}	62.8
$\text{RCHCHR} \rightarrow \text{RCCR}$	2285	80	2365	4×10^{-4}	71.1
$\text{CH}_2\text{OH} \rightarrow \text{CHO}$	1749	228	1977	5×10^{-4}	62.7
dieneophile + diene \rightarrow cyclohexene	7896	112	8008	1×10^{-4}	73.2
$\text{RCONH}_2 \rightarrow \text{RCONH}_2\text{OH}$	0	127	127	8×10^{-3}	60.9
$\text{RCONH}_2\text{OH} \rightarrow \text{RCOOH}$	1787	127	1914	5×10^{-4}	60.2

Table 1 Results of the neighbor test using the average reaction vector to transform reactants to products. If the product obtained using the average reaction vectors estimate is indeed nearest to the true product's embedding from the GNN space in a Euclidean sense, then it counts as a success to the neighbor test.

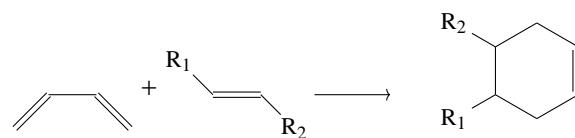
sine similarity of 0.72, as shown in the Table, between alkane and alkene oxidation for instance, is an appreciable breach to orthogonality and implies significant colinearity in the 128-D latent space.

The significance of this result is two-fold. First, it means that oxidation is always in a similar direction in the embedding space. Explicitly, that means its opposite, reduction, must be in the exact opposite direction. Second, it is apparent that the embedding space is organized based on chemical composition; The removal of hydrogen is colinear in the embedding space. Later, we will quantify that observation using a proposed replicate model of the embedding space based on neighborhood composition. This replicate model will map all the trajectories of compositional changes for the embedding (adding carbons, removing nitrogens, ... etc), up to several bonds away.

We mentioned previously that the neighbor test yields near perfect results when hydrogens are removed without the subsequent optimization step for any oxidation reaction. The difference in performance between the optimized and non-optimized results is chemically meaningful, and can be best illustrated in the case of oxidation of an alkane. In such a reaction, it is possible to have either the cis or the trans product, which are stereoisomers of each other. A closer look at the oxidation vectors to both isomers we find a difference in the cosine similarity of the reaction vector going to trans vs going to cis. For the cis isomers, the high-dimensional cosine similarity of the embedding is consistently larger at an average cosine similarity of 0.95, whereas for the trans isomers, an average cosine similarity of 0.91 is obtained with respect to the mean reaction vector embedding. The difference between the cis and trans cosine similarities with respect to the mean reaction vector is significant according to an independent samples t-test which gave a p-value of 1×10^{-4} . The neighbor test using the average reaction vector yields slightly better results on cis (67.1%) than on trans (54.9%), which implies that it is biased towards cis geometry, as QM9 has many alkanes inside rings. This explains part of the discrepancy in the neighbor test going from non-optimized to optimized, as product alkenes determine which shape they will hold implies that a single average reaction vector approximation for all possible optimized states becomes a biased assumption. Additionally, we performed MMFF94 geometry optimization on our new products for efficiency, whereas the QM9 dataset, and correspondingly the SchNet model, is optimized at the B3LYP/6-31G(d,p) level of theory. Both of these discrepancies may be slightly affecting the results, nonetheless, even with these approximations it is clear that geometry optimization plays only a minor role to that of chemical composition in mapping out the chemical embedding space.

2.2 One-step Reactions: Diels-Alder

The oxidation reactions discussed in the previous section demonstrate leaving group reactions whereby atoms leave the reaction center. The rational next question is about incoming groups making a bond at the reaction center. Can we still find linear analogies for this slightly increased complexity? The answer is affirmative. A class of such a reaction is the diels-alder reaction:

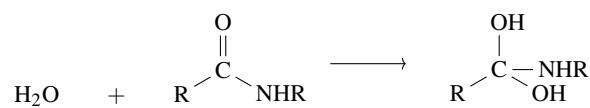


Once again, we find a strong linear analogy for the reaction in the embedding space, whereby the neighbor test yields a 73% of the transformed diels-alder, using the average reaction vector estimate. See Figure 2d for the PC projected reaction embedding vectors of the diels-alder reaction.

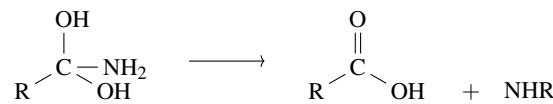
Diels-alder also brings to light an additional corroborating observation. Diels-alder shares similarities with reduction, because the double bond at the reaction center is reduced to a single bond as the ring closes to make the resulting product adduct. Evidently this makes the reaction vector for diels-alder face the opposite direction to oxidation reaction vectors as can be seen from Table 2, where the cosine similarity between diels-alder and for instance alkane oxidation is -0.73 . This can also be seen in the PC projection of the reaction vectors in Figure 2d when compared with that for oxidation, Figure 2a. Therefore, even though diels-alder is not technically a reduction via adding hydrogens, it is highly colinear with reduction which once again points to a highly organized space. It has been observed in past work⁷¹ that SchNet's modeling of chemistry is interpretable based on chemical bond-order. Our results put their findings in a larger framework based in the implicit chemical syntactical relationship between the embedding's various subspaces.

2.3 Multi-step Reactions: Hydrolysis

The hydrolysis of amides to make carboxylic acids is a two-step reaction process, first making a tetrahedral intermediate of a carboxylic acid through the imission of a water molecule,



After a quick proton transfer step, an ammonia then leaves the tetrahedral intermediate,



Each of these processes comes with its own distinct embedding transformation, see Figures 2e and 2f for the PC projected vectors for each step. The average reaction vector (in the original 128-D space), for each step, is once again a good proxy to transform any reactant as shown in Table 1.

Analyzing cosine similarity also proves insightful in the case of amide hydrolysis. This can be seen at the second step of the reaction when the ammonia leaves the reaction center and forms the carboxylic acid. This step has similar characteristics to the oxidation of alcohols, because as the ammonia leaves, a double bond is formed at the reaction center making a carbonyl product. Cosine similarity between these two processes, Table 2 row

	alkane ox	alkene ox	alcohol ox	diels-alder	hydro step 1	hydro step 2
alkane ox	1.00					
alkene ox	0.72	1.00				
alcohol ox	0.35	0.29	1.00			
diels-alder	-0.73	-0.54	0.00	1.00		
hydro step 1	-0.17	-0.10	-0.70	0.09	1.00	
hydro step 2	0.14	0.13	0.43	-0.01	-0.64	1.00

Table 2 The high-dimensional cosine similarity between the various average reaction vectors studied. Most reaction embedding vectors represent distinct reaction trajectories, and are thus close to orthogonal, however where orthogonality is breached it signifies similarity between reaction trajectories in the embedding space.

6, is 0.43, well beyond the prevalent orthogonality window of high-dimensional space ($\sigma = 0.008$). Reiterated, this points to a highly organized space with implicit relationship between its various subspaces.

We have restricted ourselves to just three types of elementary reactions, but similar conclusions can be drawn for other processes using the provided methodology. Our query software has been written with sufficient generality in mind, to quickly query any reactant (with any short- and long-range features), remove leaving groups, build any specified functional group, add/remove bonds at the reactions site, and extract embeddings at the reaction site, in a fully automated manner. The software is open source and the link can be found in the Supplementary Section.

2.4 Perturbational Replicates

Regardless of the multiple examples shown in the previous section alluding to a highly organized embedding space in terms of chemical analogies, this structure may still be perceived as a coincidence similar to the analogies found in natural language. However, in contrast to NLP, we can explain the overall structure of the embedding space from chemical principles. We will show how it is possible to replicate the embedding space by introducing consecutively layered neighborhood composition to model the atom-embedding, which we call perturbational replicates. Moreover, we are able to explain the nearly constant nature of reaction vectors found in the previous Section. We shall also see how this leads naturally to the implicit syntactical relationships found in the previous Section, and how the replicates reveal a valid interpretation of decision-making of such models, shedding light on its "black box" nature.

The embedding replicates are based on introducing the chemical neighborhood composition around each atom in layers. Thus, formally, at the very first layer, the zeroth order perturbation, the embedding replicates for atom i are defined as just the average embedding in the entire dataset, \bar{x}_i , that is,

$$x_i^{(0)} = \bar{x}_i. \quad (7)$$

So, the 0th order replicate of the oxygen atom is just the gross average of all oxygen embedding vectors in the set and does not contain any functional group representation. Following this, at the first order perturbation, the replicates are updated using the embeddings corresponding to their direct neighborhood. This process is repeated for all the atoms in the molecule, and generally for any atom-embedding, the first perturbation is given by

$$x_i^{(1)} = \bar{x}_i^{(0)} + \sum_{j \in m_i[1]} c_j \bar{x}_j^{(0)}, \quad (8)$$

where $m_i[1]$ is the set of neighbors that are only one bond away from atom i . The perturbational update uses the average embedding of the element for neighbor j , $\bar{x}_j^{(0)}$. The c_j are the linear coefficients that fit the resulting replicate at first perturbation, for each atom-embedding, $x_i^{(1)}$, with the exact embedding vector. From a practical point of view, it is crucial to design a learning algorithm that labels chemical neighborhoods uniquely without redundancy in the representation. This requires an ordering of the elements of the neighborhood environment for which the practical implementation of the above equation can be found in the supplementary material.

We repeat our approach for a second order perturbational update,

$$x_i^{(2)} = x_i^{(1)} + \sum_{j \in m_i[1]} c_j x_j^{(1)}. \quad (9)$$

The only difference now is that we no longer rely on using the average neighbor-type embedding but rather here, $x_j^{(1)}$ are the results of the previous perturbational update which took into account the direct neighborhood. Since the direct neighboring atoms also incorporated their own neighborhood from the first order perturbation, then the atom-embedding for atom i is now recognizing the indirect neighborhood layer that is two bonds away. If repeated one more time, for the third order perturbational replicate, then atom i is indirectly incorporating the neighborhood layer that is three bonds away, and so on. Again, this is repeated for each atom in the molecule to get all of their embeddings replicates based on the third layer of neighborhood composition.

Note that at each update, the number of features being used is the same, $D_{\text{embedding}} \times D_{\text{elements}}$, and only involves the embedding of the direct neighbors, either from the previous replicate or the initial average. However, in contrast to the true GNN which produces the embeddings, the replicates effectively reduce that model's correlations to be based purely on neighborhood composition. While this does mean that the exact graph layout of the neighborhood is ignored from our feature space, it provides us with the majority of the embedding space.

Table 3 show how well the perturbative replicates reproduce the true embedding vector for alkenes only using both the neighbor test and the mean distance to the true embedding. The lat-

perturbation	linear	non-linear 1	non-linear 2
0	0.26% 12.6	0.26% 12.6	0.26% 12.6
1	0.26% 8.22	0.26% 7.96	0.29% 7.00
2	3.52% 6.93	4.11% 7.22	5.28% 3.87
3	29.5% 3.60	20.8% 4.82	41.6% 2.42
4	47.5% 3.01	40.2% 3.58	71.6% 2.04
5	57.4% 2.99	54.8% 3.05	82.1% 1.84
6	56.6% 2.96	59.8% 2.96	82.7% 1.70

Table 3 Results of the neighbor test and mean distance to true embedding for alkene embeddings in the QM9 test set (341 molecules), after successive application of the perturbational updates to form the embedding replicates. The updates were fitted with both linear and non-linear regression. For non-linear we used neural networks with two architectures, one with a non-linear ‘tanh’ function on the 128-feature replicate (‘non-linear 1’), and another one with an additional linear layer (‘non-linear 2’). If the embedding lies nearest to its true GNN embedding then that counts as a success to the neighbor test. The table also compares the mean distance to the true GNN embedding, which can be compared to the mean minimum distance between alkene embeddings (1.79) or compared to the mean distance between any two alkene embeddings (9.46).

ter needs to be compared with mean distance between neighboring embeddings (1.79) and the average distance between embeddings of the same class (alkenes: 9.46). Additionally, the Table shows results for both the linear model mentioned before, and an added non-linearity (non-linear 1) and one with additional extra linearity (non-linear 2). For the non-linear regression, we used a neural network that consists of a single linear layer which inputs neighborhood embedding feature space ($D_{\text{embedding}} \times D_{\text{elements}}$), and outputs 128 features, the same size as the embedding space, which was then followed by a tanh non-linearity, that is

$$x_i^{(n)} = x_i^{(n-1)} + \tanh\left(\sum_{j \in m_i[1]} c_j x_j^{(n-1)}\right). \quad (10)$$

Using linear regression, the embedding replicate space is a near 56% replicate of the true embedding space based on the neighbor test after the sixth neighborhood layer for the alkenes embeddings. Whereas the non-linear coefficients can provide a replicate of up to 82% success on the neighbor test by the fifth order replicate. Additionally, the mean distance to the true embedding (1.70) is near the minimum distance between embeddings (1.79). This means that it is nearing the neighborhood density of the true embedding model underlining the limitations of the neighbor test. Comparing this distance with the average distance between any two alkene embeddings (9.46), we can see that we are well-within the alkene predictions, making fine-tuned replicates based on multiple bonds away.

The perturbational model can also be generalized further by including beyond the direct neighbors for prediction. That is by increasing the depth of neighbors allowed for defining the chemical composition, $m_i[2]$, $m_i[3]$, and so on, at each perturbative update. This added freedom contributes to significantly lowering the error on the embedding replication as summarized in Table 4 using the RMSE loss metric on the original GNN embedding space.

To test how the replicates of the embeddings perform on real chemical properties we employed our pre-trained transfer learning model on ^1H -NMR and ^{13}C -NMR, see section 1.1 for details

pert/depth	1	2	3	4	5
1	0.594	0.514	0.473	0.458	0.454
2	0.451	0.327	0.241	0.210	0.191
3	0.345	0.238	0.208	0.188	0.172
4	0.296	0.237	0.203	0.184	0.170
5	0.287	0.237	0.201	0.184	0.170
6	0.285	0.237	0.201	0.183	0.170

Table 4 RMSE between perturbational replicates and the true GNN embedding space using various neighborhood depths in the input neighborhood feature representation that predict the replicates.

on the transfer learning model. If it is true that the replicates are getting closer to the true embedding space, they should then perform successively more accurately on such properties as the replicates incorporate further neighborhood layers into the prediction. Using a neighborhood depth of one and the ‘non-linear 2’ model for replicating the embedding space based on neighborhood composition, Figures 3a and 3b for ^1H -NMR and ^{13}C -NMR confirm this effect revealing how the embedding space works to make finer predictions based on layered chemical composition.

With the perturbational replicates at hand, we are now in a position to prove the linear analogies up to first order perturbation. We first define a formal reaction vector, under any given perturbation. At a perturbation of m , for example, we can isolate a reaction vector as

$$\Delta X_{rxm_{r \rightarrow p}}^{(m)} = x_r^{(m)} - x_p^{(m)}. \quad (11)$$

For an oxidation of an alcohol at a perturbation of 1, this would give the following

$$\Delta X_{ox,O}^{(1)} = c_C \bar{x}_C^{(0)} - (c_C \bar{x}_C^{(0)} + c_H \bar{x}_H^{(0)}), \quad (12)$$

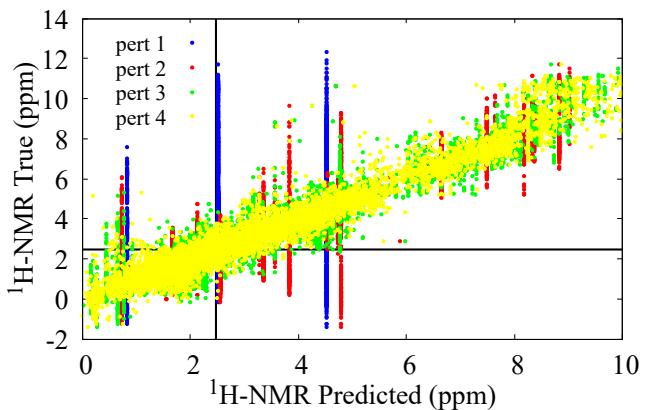
for the oxygen, where the left-hand term represents the product, and the right-hand term represent the reactant oxygen of alcohol having both a hydrogen and a carbon neighbor embedding. The difference amounts to only the 0th order embedding of the hydrogen, which was removed at that site,

$$\Delta X_{ox,O}^{(1)} = -c_H \bar{x}_H^{(0)}. \quad (13)$$

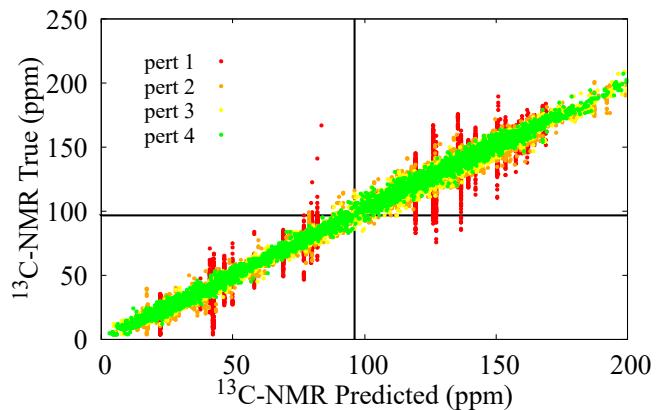
The key observation is that the reaction vector is proportional to the average $\bar{x}_H^{(0)}$, a vector which is independent of the local neighborhood of the oxygen. Similarly, the reaction vector for the oxidation from alkanes to alkenes is

$$\Delta X_{ox,C}^{(1)} = -c_H \bar{x}_H^{(0)}, \quad (14)$$

which is also proportional to $\bar{x}_H^{(0)}$ and independent of the C environment. This explains how the perturbational replicates can reduce to the linear analogies found at the first order perturbation, and how similar changes whether on alkane or alcohol are colinear. As Table 3 shows, we obtain greater accuracy with greater perturbations on the neighborhood embedding replicates. Of course, including the higher order corrections, will integrate long-range effects and provide a more fine-tuned reaction vector that will explain deviations from the exact constant vector.



(a)



(b)

Fig. 3 The accuracy of perturbational replicates across the 2000 molecule test set for gas-phase ^{13}C -NMR and ^1H -NMR compared to the DFT-computed values in QM9NMR, see Section 1.1. The black line represents the value of prediction for perturbational replicate of zero (i.e. the gross average hydrogen or carbon embedding). For proton NMR, the first five perturbational replicates achieve a RMSE of 2.23 ppm, 2.14 ppm, 0.98 ppm, 0.64 ppm, and 0.51 ppm, respectively. For carbon NMR, the first five perturbations replicates give a RMSE of 20.37 ppm, 11.27 ppm, 5.45 ppm, 3.47 ppm, and 2.82 ppm, respectively.

2.5 Applications of Linear Analogies

An important question is to what extent these linear analogies in abstract latent space will provide imprints on real chemical observables. For instance, one can intuitively expect that chemical observables that are reasonably well described by a quasi-linear model in this organized latent space will display similar behaviour. Remarkably, it is possible to extract tangible relations for chemical observables based on fairly general assumptions. Consider an atom-wise chemical observable that can be obtained as a function $f(x)$ of its chemical composition or atomic embedding vector, e.g. ^{13}C -NMR shifts, either from an end-to-end⁷² or transfer learning⁶² point of view. Our first assumption is that this function $f(x)$ should be sufficiently smooth over each individual functional group class, so one can expand the function around the average \bar{x} , or 0th order embedding vector $x^{(0)}$ for that functional group

$$f(x) = f(\bar{x} + \delta) = f(\bar{x}) + \delta \cdot \nabla f(\bar{x}) + \mathcal{O}(\delta^2). \quad (15)$$

Our second assumption resides on the validity of the linear analogy, stating that there exists a constant reaction vector $\langle \text{reaction} \rangle = \Delta$ that brings us from the reactant (r) to the product (p) for all reaction pairs

$$x_r + \Delta = x_p, \quad (16)$$

which includes the averages \bar{x} , by construction. Furthermore, it implies that each δ specifying the individual end points of the reaction are necessarily equal

$$\delta_r = \delta_p =: \delta. \quad (17)$$

As a result, one can write for both end points

$$f(x_r) = f(\bar{x}_r) + \delta \cdot \nabla f(\bar{x}_r) + \mathcal{O}(\delta^2), \quad (18)$$

$$f(x_p) = f(\bar{x}_p) + \delta \cdot \nabla f(\bar{x}_p) + \mathcal{O}(\delta^2). \quad (19)$$

Solving δ formally from the first equation, we obtain the simple linear relation

$$f(x_p) = f(\bar{x}_p) + [f(x_r) - f(\bar{x}_r)] \nabla f(\bar{x}_r)^{-1} \cdot \nabla f(\bar{x}_p) + \mathcal{O}(\delta^2), \quad (20)$$

$$= \alpha f(x_r) + \beta. \quad (21)$$

with the constants α and β only dependent on the overall functional group class embedding vectors \bar{x}_r and \bar{x}_p , and not on the individual reaction end points. The relation becomes even stronger when considering a global linear transfer model

$$f(x) = a \cdot x + b, \quad (22)$$

for which a quick insightful re-derivation of (21) yields

$$f(x_p) = a \cdot x_p + b = a \cdot (x_r + \Delta) + b = f(x_r) + a \cdot \Delta, \quad (23)$$

so $\alpha = 1$ and $\beta = a \cdot \Delta$. In Figures 4, we investigate relation (21) on ^{13}C -NMR shifts and atomization energies for the alkane and alcohol reactions, extracted from QM9NMR and SchNet respectively. In previous work,^{61,62} we showed that ^{13}C -NMR shifts are reasonably well reproduced from a global linear regression model in the latent space. For our oxidation reactions, see Figures 4a and 4c, despite the fluctuations [$R^2 = 0.41$, $R^2 = 0.26$], a linear fit on the data is in line with the $\alpha = 1$ relation, reconfirming the linear analogies in embedding space between the reaction pairs. To the best of our knowledge, this is the first observation of the simple linear relation of NMR-shifts between reaction pairs, entirely facilitated by the linear analogies in the underlying embedding

space.

A fit on the atomization energies also yield linear relations (21) with $\alpha \approx 1.8$ [$R^2 = 0.79$, $R^2 = 0.77$] (see Figures 4b and 4d). This deviation from $\alpha = 1$ is to be expected from the architecture of SchNet, in which the final atomization energies are obtained from a fully connected feed-forward neural network. Furthermore, a closer scrutiny of Figures 4b and 4d reveals a substructure into bands for which each individual α is closer to 1, which may point towards a more locally fine-tuned sub-classification of the alcohol/aldehyde and alkane/alkene groups for atomization energy.

3 Conclusions

In this work, we uncover a latent space of graph neural network models that maintains a high degree of structural integrity. The structure of the graph neural network latent space is largely based on a fine chemical syntax organization. We demonstrate this via the use of linear analogies, constant vectors that help transform from one chemical formula to another, a.k.a reaction vectors. We observed how linear analogies themselves form a coherent structure, in that similar reactions (ex. oxidations) are colinear in the latent space. Thus, the structure of the embedding space can be thought of to have two levels of organization. The first level is that of molecular substructure composition; similar substructures are placed next to each other (alkanes and alkenes vs aldehydes and ketones). The second level is that of changes to molecular substructure, similar chemical changes are in the same direction in the latent space. This is in line with previous observations in natural language models in which word analogies can be found in using vector arithmetic, such as ‘King’ – ‘Man’ + ‘Woman’ = ‘Queen’. In a similar vein, a chemist can write: ‘Amide’ – ‘Amine’ + ‘Alcohol’ = ‘Carboxylic Acid,’ only in a quantitative chemical compositional sense.

Our observations were largely explained by a replicate model of the latent space based on perturbational updates that integrate neighborhood chemical compositions in layers. This perturbational model relates the structure of the latent space to chemical composition. We showed how such a model can explain the approximately constant reaction vectors. We then showed how linear analogies in the latent space carry over to chemical properties such as NMR and atomization energies. In other words, the integrity of linear analogies lies beyond just the latent space and can be used to explain quasi linear changes in chemical properties, such as constant NMR shifts, and near-constant changes in atomization energies.

To the extent of our knowledge, this is the first work that uncovers a quantitative global explanatory framework from chemistry’s deep learning models, such as GNN models. In the past, explainability has been done locally, using extrinsic tools, and observing how the model responds to limited examples. However, our linear analogies structure, and perturbational replications sets the framework for a global explanation to the latent space of GNN chemistry. While the main frame is set, there is still much room to explore. One direction is to map out the trajectories of reactions in the latent space, and uncover what is happening between end-points of reactions, i.e. transition states. Another direction is

to study how this global explanatory framework reduces to a local one and thus can give us explanations on a case-by-case basis. For example, we can study how the replicates improve prediction of chemical quantities (or changes in chemical quantities) based on finely tuning the neighborhood composition. This also sets the stage for a global evaluation of model generalizability and transferability. Lastly, there is room for improving the replicates of the GNN model. For instance, the perturbational replicate model uses a neighborhood composition feature space that ignores the exact layout of the graph, by incorporating the exact layout, and using the same geometry optimization as the GNN model, can lead to a finer replicate. Nonetheless, the findings in this paper set the stage for global explorations in GNN modeling of chemistry in a single coherent framework.

4 Supplementary Information

The reaction creation and analysis code can be found at: <https://github.com/QuNB-Repo/DLChem/tree/master>. This automates reaction dataset creation. It is able to query a dataset for all reactants specified (with any long-range feature), remove leaving groups, build any functional group, and extract embeddings at the reaction site for analysis.

4.1 Representing Neighborhood Feature Spaces

Practically, to represent equation 2.4, we must avoid problems of having a variant number of neighbors as we cannot handle data of various sizes without reverting back to graph neural networks. Additionally, we must ensure our representation for the neighborhood is unique to each chemical neighborhood. To avoid these problems we use a data structure whereby the chemical neighborhood feature space is described by embedding-sized placeholders for the existence of H, C, N, O, and/or F neighbor, in that order. For instance, being an oxygen atom in an alcohol, will fill the carbon and the hydrogen embedding placeholders for each update (while leaving the rest as zero), with either some previous update or with the initial average for that neighboring element. If, for instance, two carbons are found as neighbors, then their embeddings are added onto the same placeholder for the carbon neighbor. This practical approach avoids issues of variant neighborhood sizes, as now at each update the number of features is the same, $D_{\text{embedding}} \times D_{\text{elements}}$. If there are multiple neighborhood depths included in the representations, then the order of the neighbor elements is repeated for each depth such that we have $D_{\text{embedding}} \times D_{\text{elements}} \times D_{\text{depth}}$. This ensures uniqueness in the representation even across multiple layers of neighborhood depth.

Author Contributions

SDB & AES contributed equally to the conceptualization and methodology development of the project. AES implemented the methodology, data curation, and produced results. SDB & AES contributed to the analysis of the results. AES wrote the manuscript and SDB edit it. SDB provided funding acquisition.

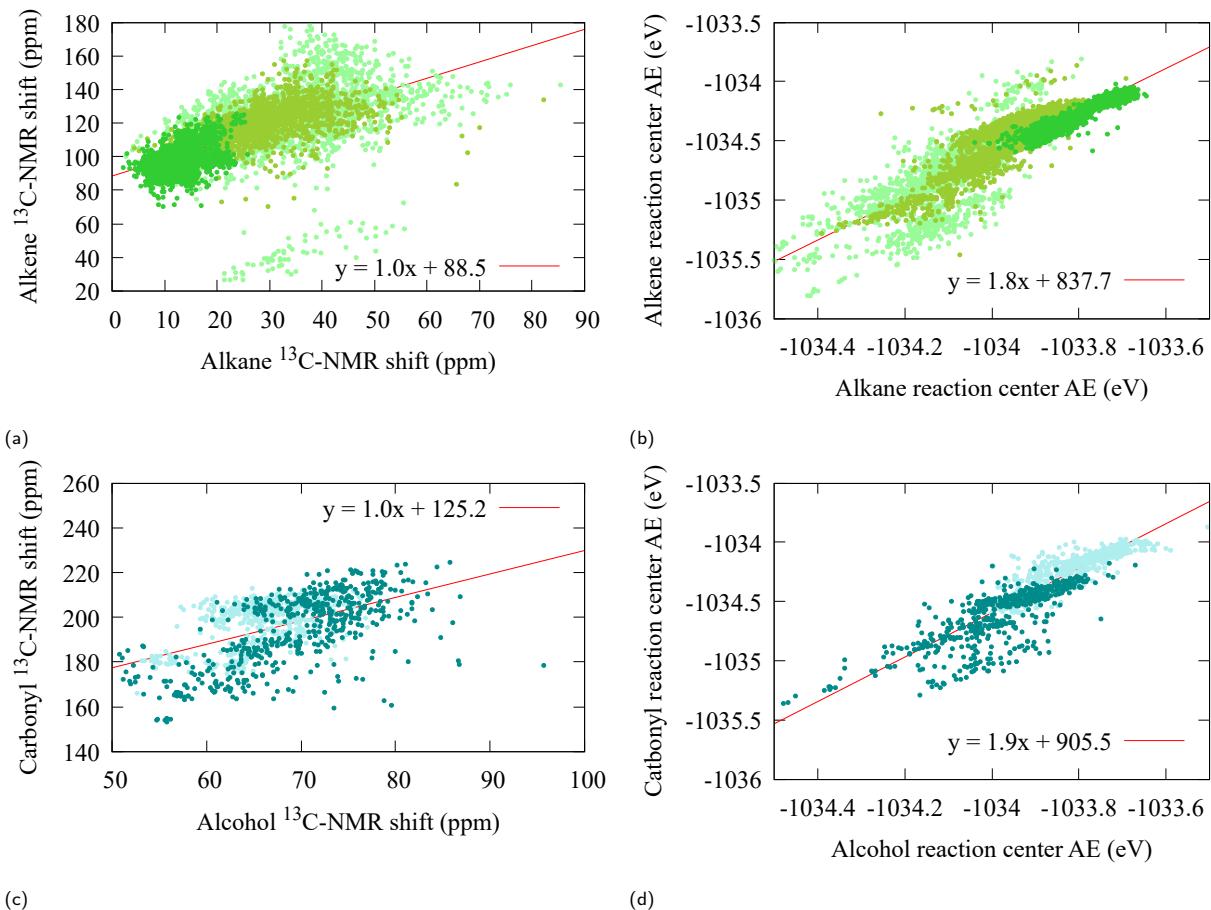


Fig. 4 The ^{13}C -NMR chemical shift (a, c), and the change in atomization energy (b, d), associated with the alkane and alcohol oxidation, at the carbon reaction center. The colors represent the functional group of the reactant that was involved in the oxidation, methyls, methylenes, methines, alcohols, carbonyls

Conflicts of Interest

"There are no conflicts to declare".

Acknowledgements

Financial support from NSERC, under the discovery grant program, CFI program and NBIF RAI program are acknowledged. SDB thanks the Canada Research Chair program for financial support.

Notes and references

- 1 A. C. Mater and M. L. Coote, *Journal of chemical information and modeling*, 2019, **59**, 2545–2559.
- 2 G. B. Goh, N. O. Hodas and A. Vishnu, *Journal of computational chemistry*, 2017, **38**, 1291–1307.
- 3 M. Vogt, *Expert Opinion on Drug Discovery*, 2022, **17**, 297–304.
- 4 A. Mayr, G. Klambauer, T. Unterthiner and S. Hochreiter, *Frontiers in Environmental Science*, 2016, **3**, 80.
- 5 J. Ma, R. P. Sheridan, A. Liaw, G. E. Dahl and V. Svetnik, *Journal of Chemical Information and Modeling*, 2015, **55**, 263.
- 6 T. Unterthiner, A. Mayr, G. Klambauer, M. Steijaert, J. K. Wegner, H. Ceulemans and S. Hochreiter, Proceedings of the Deep Learning Workshop at NIPS, 2014, p. 1.
- 7 G. E. Dahl, N. Jaitly and R. Salakhutdinov, *arXiv preprint arXiv:1406.1231*, 2014.
- 8 A. Korotcov, V. Tkachenko, D. P. Russo and S. Ekins, *Molecular Pharmaceutics*, 2017, **14**, 4462.
- 9 T. Unterthiner, A. Mayr, G. Klambauer and S. Hochreiter, *arXiv preprint arXiv:1503.01445*, 2015.
- 10 J. Wenzel, H. Matter and F. Schmidt, *Journal of Chemical Information and Modeling*, 2019, **59**, 1253.
- 11 M. Li, H. Zhang, B. Chen, Y. Wu and L. Guan, *Scientific Reports*, 2018, **8**, 1.
- 12 K. Mills, M. Spanner and I. Tamblyn, *Physical Review A*, 2017, **96**, 042113.
- 13 K. Yao and J. Parkhill, *Journal of Chemical Theory and Computation*, 2016, **12**, 1139.
- 14 R. T. McGibbon, A. G. Taube, A. G. Donchev, K. Siva, F. Hernández, C. Hargus, K.-H. Law, J. L. Klepeis and D. E. Shaw, *The Journal of Chemical Physics*, 2017, **147**, 161725.
- 15 S. Lorenz, A. Groß and M. Scheffler, *Chemical Physics Letters*, 2004, **395**, 210.
- 16 T. B. Blank, S. D. Brown, A. W. Calhoun and D. J. Doren, *The Journal of Chemical Physics*, 1995, **103**, 4129.
- 17 K. T. Schütt, P.-J. Kindermans, H. E. Sauceda, S. Chmiela, A. Tkatchenko and K.-R. Müller, *arXiv preprint arXiv:1706.08566*, 2017.
- 18 K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller and A. Tkatchenko, *Nature Communications*, 2017, **8**, 1.
- 19 K. Schutt, P. Kessel, M. Gastegger, K. Nicoli, A. Tkatchenko and K.-R. Müller, *Journal of Chemical Theory and Computation*, 2018, **15**, 448.
- 20 K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko and K.-R. Müller, *The Journal of Chemical Physics*, 2018, **148**, 241722.
- 21 O. T. Unke and M. Meuwly, *Journal of Chemical Theory and Computation*, 2019, **15**, 3678.
- 22 J. S. Smith, O. Isayev and A. E. Roitberg, *Chemical Science*, 2017, **8**, 3192.
- 23 R. Zubatyuk, J. S. Smith, J. Leszczynski and O. Isayev, *Science Advances*, 2019, **5**, eaav6490.
- 24 J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals and G. E. Dahl, International conference on machine learning, 2017, p. 1263.
- 25 J. Jo, B. Kwak, H.-S. Choi and S. Yoon, *Methods*, 2020, **179**, 65.
- 26 J. Han, H. Kang, S. Kang, Y. Kwon, D. Lee and Y.-S. Choi, *Physical Chemistry Chemical Physics*, 2022, **24**, 26870.
- 27 Y. Kwon, D. Lee, Y.-S. Choi, M. Kang and S. Kang, *Journal of Chemical Information and Modeling*, 2020, **60**, 2024.
- 28 H. Rull, M. Fischer and S. Kuhn, *arXiv preprint arXiv:2304.03361*, 2023.
- 29 J. Xiong, Z. Li, G. Wang, Z. Fu, F. Zhong, T. Xu, X. Liu, Z. Huang, X. Liu, K. Chen et al., *Bioinformatics*, 2022, **38**, 792.
- 30 D. Zhang, S. Xia and Y. Zhang, *Journal of Chemical Information and Modeling*, 2022, **62**, 1840.
- 31 Y. Pathak, S. Mehta and U. D. Priyakumar, *Journal of Chemical Information and Modeling*, 2021, **61**, 689.
- 32 K. Low, M. L. Coote and E. I. Izgorodina, *Journal of Chemical Information and Modeling*, 2022, **62**, 5457.
- 33 L. David, A. Thakkar, R. Mercado and O. Engkvist, *Journal of Cheminformatics*, 2020, **12**, 1.
- 34 J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li and M. Sun, *AI Open*, 2020, **1**, 57.
- 35 O. Wieder, S. Kohlbacher, M. Kuenemann, A. Garon, P. Ducrot, T. Seidel and T. Langer, *Drug Discovery Today: Technologies*, 2020, **37**, 1.
- 36 Y. Wang, Z. Li and A. B. Farimani, *arXiv preprint arXiv:2209.05582*, 2022.
- 37 P. Schwaller, T. Gaudin, D. Lanyi, C. Bekas and T. Laino, *Chemical science*, 2018, **9**, 6091–6098.
- 38 R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, J. Aguilera-Iparraguirre, T. D. Hirzel, R. P. Adams and A. Aspuru-Guzik, *ACS central science*, 2018, **4**, 268–276.
- 39 S. Jastrzębski, D. Leśniak and W. M. Czarnecki, *arXiv preprint arXiv:1602.06289*, 2016.
- 40 M. J. Kusner, B. Paige and J. M. Hernández-Lobato, International conference on machine learning, 2017, pp. 1945–1954.
- 41 E. J. Bjerrum, *arXiv preprint arXiv:1703.07076*, 2017.
- 42 M. H. Segler, T. Kogej, C. Tyrchan and M. P. Waller, *ACS central science*, 2018, **4**, 120–131.
- 43 M. Hirohara, Y. Saito, Y. Koda, K. Sato and Y. Sakakibara, *BMC bioinformatics*, 2018, **19**, 83–94.
- 44 G. B. Goh, N. O. Hodas, C. Siegel and A. Vishnu, *arXiv preprint arXiv:1712.02034*, 2017.

- 45 D. Weininger, *Journal of chemical information and computer sciences*, 1988, **28**, 31–36.
- 46 D. Weininger, A. Weininger and J. L. Weininger, *Journal of chemical information and computer sciences*, 1989, **29**, 97–101.
- 47 M. Krenn, Q. Ai, S. Barthel, N. Carson, A. Frei, N. C. Frey, P. Friederich, T. Gaudin, A. A. Gayle, K. M. Jablonka *et al.*, *Patterns*, 2022, **3**, 10.
- 48 M. Krenn, F. Häse, A. Nigam, P. Friederich and A. Aspuru-Guzik, *Machine Learning: Science and Technology*, 2020, **1**, 045024.
- 49 J. N. Wei, D. Duvenaud and A. Aspuru-Guzik, *ACS central science*, 2016, **2**, 725–732.
- 50 C. W. Coley, R. Barzilay, T. S. Jaakkola, W. H. Green and K. F. Jensen, *ACS central science*, 2017, **3**, 434–443.
- 51 K. Ethayarajh, D. Duvenaud and G. Hirst, *arXiv preprint arXiv:1810.04882*, 2018.
- 52 A. Gittens, D. Achlioptas and M. W. Mahoney, Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), 2017, pp. 69–76.
- 53 S. Arora, Y. Li, Y. Liang, T. Ma and A. Risteski, *Transactions of the Association for Computational Linguistics*, 2016, **4**, 385–399.
- 54 A. Drozd, A. Gladkova and S. Matsuoka, Proceedings of coling 2016, the 26th international conference on computational linguistics: Technical papers, 2016, pp. 3519–3530.
- 55 J. Pennington, R. Socher and C. D. Manning, Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), 2014, pp. 1532–1543.
- 56 R. Bamler and S. Mandt, International conference on Machine learning, 2017, pp. 380–389.
- 57 M. Kusner, Y. Sun, N. Kolkin and K. Weinberger, International conference on machine learning, 2015, pp. 957–966.
- 58 M.-E. Brunet, C. Alkalay-Houlihan, A. Anderson and R. Zemel, International conference on machine learning, 2019, pp. 803–811.
- 59 R. Petrolito and F. Dell'Orletta, *Turin, Italy*, 2018.
- 60 S. Wang, W. Zhou and C. Jiang, *Computing*, 2020, **102**, 717–740.
- 61 A. M. El-Samman, I. A. Husain, M. Huynh, S. De Castro, B. Morton and S. De Baerdemacker, *Digital Discovery*, 2024, **3**, 544–557.
- 62 A. M. El-Samman, S. De Castro, B. Morton and S. De Baerdemacker, *Canadian Journal of Chemistry*, 2023, **102**, 4.
- 63 A. M. El-Samman, *SchNet Model Embedding Vectors of QM9 Atoms Labelled According to Functional Groups Designation*, 2023, <https://doi.org/10.25545/EK1EQA>.
- 64 R. Ramakrishnan, P. O. Dral, M. Rupp and O. A. Von Lilienfeld, *Scientific Data*, 2014, **1**, 1.
- 65 A. Gupta, S. Chakraborty and R. Ramakrishnan, *Machine Learning: Science and Technology*, 2021, **2**, 035010.
- 66 I. Hunt, *Basic IUPAC Organic Nomenclature: E- and Z-nomenclature of alkenes*, <https://www.chem.ucalgary.ca/courses/350/WebContent/orgnom/alkenes/alkenes-03.html>, Accessed: 06 24, 2024.
- 67 P. Tosco, N. Stiefl and G. Landrum, *Journal of cheminformatics*, 2014, **6**, 1–4.
- 68 H. Abdi and L. J. Williams, *Wiley Interdisciplinary Reviews: Computational Statistics*, 2010, **2**, 433.
- 69 M. Gallegos, V. Vassilev-Galindo, I. Poltavsky, Á. Martín Pendás and A. Tkatchenko, *Nature Communications*, 2024, **15**, 4345.
- 70 T. Hastie, R. Tibshirani, J. H. Friedman and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*, Springer, 2009, vol. 2.
- 71 S. Letzgus, P. Wagner, J. Lederer, W. Samek, K.-R. Muller and G. Montavon, *IEEE Signal Processing Magazine*, 2022, **39**, 40.
- 72 Y. Guan, S. S. Sowndarya, L. C. Gallegos, P. C. S. John and R. S. Paton, *Chemical Science*, 2021, **12**, 12012–12026.