

Euclidean Fast Attention: Machine Learning Global Atomic Representations at Linear Cost

J. Thorben Frank^{1, 2, 3}, Stefan Chmiela^{2, 3}, Klaus-Robert Müller^{1, 2, 3, 4, 5} and Oliver T. Unke¹

¹Google DeepMind, Berlin, Germany, ²Machine Learning Group, Technische Universität Berlin, 10587 Berlin, Germany, ³Berlin Institute for the Foundations of Learning and Data – BIFOLD, Germany, ⁴Max Planck Institute for Informatics, Stuhlsatzenhausweg, 66123 Saarbrücken, Germany, ⁵Department of Artificial Intelligence, Korea University, Anam-dong, Seongbuk-gu, Seoul 02841, Korea

Long-range correlations are essential across numerous machine learning tasks, especially for data embedded in Euclidean space, where the relative positions and orientations of distant components are often critical for accurate predictions. Self-attention offers a compelling mechanism for capturing these global effects, but its *quadratic* complexity presents a significant practical limitation. This problem is particularly pronounced in computational chemistry, where the stringent efficiency requirements of machine learning force fields (MLFFs) often preclude accurately modeling long-range interactions. To address this, we introduce Euclidean fast attention (EFA), a *linear*-scaling attention-like mechanism designed for Euclidean data, which can be easily incorporated into existing model architectures. A core component of EFA are novel Euclidean rotary positional encodings (ERoPE), which enable efficient encoding of spatial information while respecting essential physical symmetries. We empirically demonstrate that EFA effectively captures diverse long-range effects, enabling EFA-equipped MLFFs to describe challenging chemical interactions for which conventional MLFFs yield incorrect results.

Introduction

Many applications of machine learning (ML) are characterized by a complex interplay of short-range (local) and long-range (global) correlations. For instance, in natural language processing, the meaning of a word in a sentence sometimes can be determined solely by its immediate neighbors, whereas in other cases, it is significantly influenced by contextual information appearing several paragraphs earlier in the text.¹ Analogously, in computer vision, local visual features like texture may suffice to identify certain objects,² but fully understanding the content of a picture may necessitate integrating globally distributed visual cues.³ Modeling global correlations is particularly challenging for data embedded in Euclidean space, where the relationships between distant components can be governed by both their relative positions *and* orientations. For applications of ML in physics and computational chemistry^{4–7}, e.g., for performing molecular dynamics (MD) simulations, matters are further complicated by the fact that atomic interactions obey certain symmetry relations (e.g., translational/rotational invariance/equivariance), and

violating these constraints can severely degrade performance.⁸

MD simulations allow to study the movement of individual atoms over time and provide insights into the structure, function, and thermodynamics of molecular systems and materials.⁹ The accuracy of such simulations crucially depends on the quality of the description of the forces between atoms, which drive their motion. An increasingly popular method for modelling these interactions are machine learning force fields (MLFFs),^{6,7,10–13} which can reach the accuracy of *ab initio* electronic structure calculations at a fraction of the computational cost. MLFFs are now routinely used to study a variety of complex systems and phenomena, with applications spanning from protein dynamics^{14,15} to the discovery of new materials.¹⁶ However, despite remarkable progress in the development of MLFFs in recent years, accurately modelling long-range effects with MLFFs remains a persistent challenge. This shortcoming often stems from stringent computational efficiency requirements: Many systems of practical interest (e.g., bio-molecules) consist of hundreds of thousands of atoms, so linear scaling

w. r. t. the number of atoms is a prerequisite. The lack of an accurate treatment of long-range interactions is problematic, because although they are usually comparatively weak (in magnitude), they can play a crucial role for the stability, long time-scale dynamics, structure, and response properties of a variety of chemical and biological systems.^{17–19}

In principle, the self-attention mechanism underlying the transformer architecture²⁰ offers a compelling mechanism for capturing such global effects. It has revolutionized natural language processing and, due to its flexibility and generality, also found applications in many other fields, including computer vision,²¹ graph learning,²² and the natural sciences.²³ Unfortunately, standard self-attention has quadratic time and memory cost in the number of inputs, hindering its wide-scale adoption in MLFFs. While the memory cost of self-attention can be reduced to scale linearly by clever implementations of the algorithm,^{24,25} the time complexity remains quadratic. Fully linear-scaling variants of attention also do exist,^{26,27} but face a fundamental issue when being applied to data embedded in Euclidean space: When modelling interatomic forces, the spatial arrangement of atoms is crucially important. Encoding this information is trivial in standard quadratic-scaling attention, but it is less clear how to achieve this in linear-scaling formulations.

In this work, we address these challenges by proposing the linear-scaling Euclidean fast attention (EFA) mechanism. It enables learning global representations that encode the spatial structure of a chemical system while respecting all relevant physical symmetries. This is achieved via our novel Euclidean rotary positional encodings (ERoPE), which allow a description of the relative positions and orientations of atoms with linear complexity. EFA can be incorporated into existing MLFFs with minimal architectural modifications, enabling them to accurately model global correlations. We empirically demonstrate that EFA-augmented models are able to describe various long-range interactions and non-local effects, while MLFFs without EFA yield incorrect results as they are unable to capture long-range structure.

Context and Challenges

The aim of this section is to provide the necessary context to appreciate the difficulty of modelling long-range interactions with linear complexity and to point out challenges when working with data embedded in Euclidean space. It starts with a description of message passing neural networks (MPNNs),²⁸ a popular model class for constructing MLFFs,^{29–34} which are used throughout this work to demonstrate the difference between conventional and EFA-augmented models. In particular, we highlight how MPNNs achieve linear scaling w. r. t. to the number of atoms, and why this necessarily limits their ability of modelling long-range interactions. This is followed by a brief review of (quadratically-scaling) self-attention, a compelling mechanism for capturing global correlations, and linear-scaling variants of attention. We also describe a straightforward method to include information about the spatial structure of Euclidean data in standard self-attention, and give arguments why this is (seemingly) incompatible with linear-scaling formulations.

MPNNs represent chemical structures as graphs embedded in Euclidean space, where individual atoms correspond to nodes with associated features $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N \mid \mathbf{x}_m \in \mathbb{R}^H\}$ and positions $\mathcal{R} = \{\vec{r}_1, \dots, \vec{r}_N \mid \vec{r}_m \in \mathbb{R}^3\}$. Starting from the initial features $\mathcal{X}^{[0]}$, which are set to (learned) embeddings depending only on atomic numbers, the node representations are updated iteratively via T MessagePassing layers

$$\mathcal{X}^{[t+1]} = \text{MessagePassing}(\mathcal{X}^{[t]}, \mathcal{R}). \quad (1)$$

At each layer, all nodes connected by an edge “pass messages” to each other, typically containing information about the values of the current features modulated by the distances (or displacement vectors) between nodes. To construct MLFFs, the final representations $\mathcal{X}^{[T]}$ are used to predict atom-wise energy contributions $\{E_1, \dots, E_N \mid E_m \in \mathbb{R}\}$, and forces can be obtained by automatic differentiation of the total energy $\sum_m E_m$ w. r. t. the positions \mathcal{R} . If chemical structures were modelled as fully-connected graphs within MPNNs, so that each node could pass messages to all other nodes at every iteration of Eq. 1, their evaluation cost would scale as $O(N^2)$ with

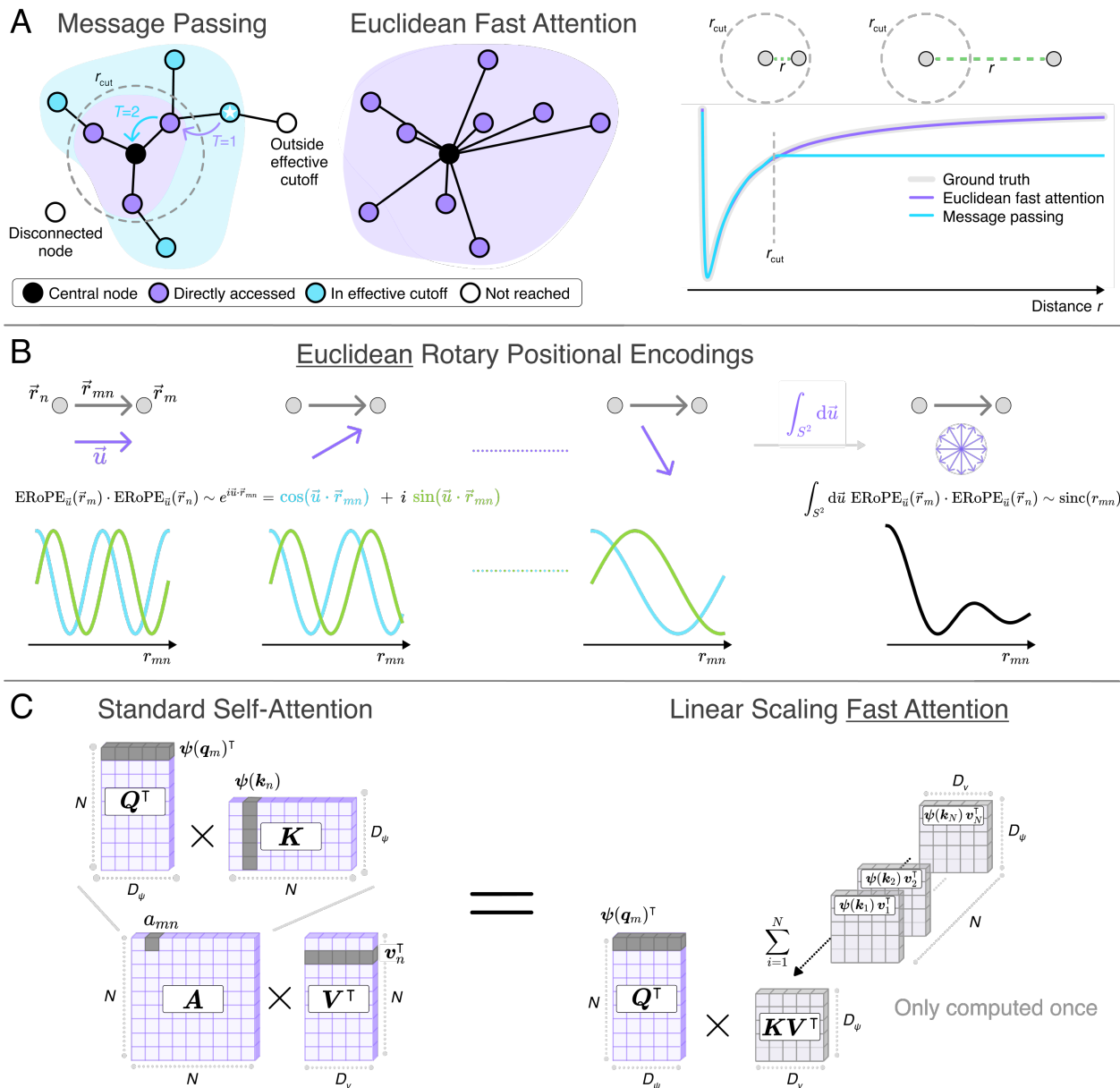


Figure 1 | Overview of the central concepts of this work. (A) Information about the overall structure of a graph/molecule accessible by a central node/atom. With message passing (MP), only nodes within the local cutoff r_{cut} (green) can directly pass information to the central node (black). With $T = 2$ update blocks (Eq. 1), information about more distant nodes (blue) also becomes available (the arrows illustrate how information “hops” from the highlighted node (\star) to the central node). Unreachable nodes (white) are either outside the effective cutoff $T \cdot r_{\text{cut}}$, or disconnected (there is no “hopping path”). In contrast, with Euclidean fast attention (EFA), all nodes can be accessed directly – irrespective of distance. This enables EFA to faithfully capture long-range effects, illustrated here for a simple pairwise potential $V(r)$, where MP fails to model $V(r)$ when r is larger than the cutoff. (B) EFA encodes geometric information using Euclidean rotary positional encodings (ERoPE), see Eq. 9. ERoPE effectively projects the displacement vectors \vec{r}_{mn} between nodes m and n onto a unit vector \vec{u} , and expands the result in a complex exponential, i.e., $e^{i\vec{u} \cdot \vec{r}_{mn}}$. The frequency of the resulting “wave” depends on the angle between \vec{r}_{mn} and the chosen \vec{u} . By averaging over all possible choices of \vec{u} (integration over the unit sphere S^2), a rotationally invariant encoding (independent of \vec{u}) of the distance r_{mn} can be obtained. (C) Schematic representation of (quadratically-scaling) standard self-attention (Eq. 7a) vs. a linear-scaling formulation²⁶ (Eq. 7b). Here, the operation is shown in matrix form, which computes attention for all inputs at once.

the number of atoms. As mentioned above, this is prohibitive for many systems of practical interest, which often consist of hundreds of thousands of atoms. To achieve $O(N)$ scaling, most MPNNs therefore introduce a cutoff r_{cut} , and two nodes m and n are only considered connected to each other if their distance is below the cutoff, i.e., $\|\vec{r}_m - \vec{r}_n\| < r_{\text{cut}}$.

While desirable from an efficiency perspective, this also means that interactions beyond a certain distance, which we refer to as *effective cutoff*, cannot be modelled by construction. Most other MLFF model architectures employ similar cutoff strategies,^{35–37} but MPNNs are special in that the effective cutoff can exceed r_{cut} . This is because nodes can also gather information from other nodes they are not directly connected to, as long as this information could first propagate to one of their direct neighbors in previous iterations of Eq. 1. The maximally possible effective cutoff is $T \cdot r_{\text{cut}}$, but if there is no “hopping path” between two nodes to act as “relay” over multiple updates, they cannot exchange information, even if their distance is below $T \cdot r_{\text{cut}}$ (see “disconnected node” in Fig. 1A). Further, this indirect information transfer between nodes only captures some “mean-field effect”, because information from multiple nodes is aggregated during each update, and afterwards cannot be attributed to a specific source anymore. As our experiments show (see Results), this is often insufficient to accurately capture long-range interactions. Since an increased effective cutoff is unique to MPNNs, we use them as base model architecture for all experiments in this work. Other MLFFs employing local cutoffs behave qualitatively similar to MPNNs with $T = 1$, so they do not need to be considered separately to show general properties of local MLFFs and their shortcomings when modelling long-range interactions.

We now describe self-attention²⁰ and highlight why it has not found wide-scale adoption in MLFFs, despite its success in modelling global correlations in other domains.^{21,23} Given a set of N features $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N \mid \mathbf{x}_m \in \mathbb{R}^H\}$, self-

attention calculates

$$\text{ATT}(\mathcal{X})_m = \frac{\sum_{n=1}^N \text{sim}(\mathbf{q}_m, \mathbf{k}_n) \mathbf{v}_n}{\sum_{n=1}^N \text{sim}(\mathbf{q}_m, \mathbf{k}_n)} = \sum_{n=1}^N a_{mn} \mathbf{v}_n, \quad (2)$$

where

$$a_{mn} = \frac{\text{sim}(\mathbf{q}_m, \mathbf{k}_n)}{\sum_{n=1}^N \text{sim}(\mathbf{q}_m, \mathbf{k}_n)}$$

are so-called attention coefficients. The vectors $\mathbf{q}, \mathbf{k} \in \mathbb{R}^{D_{qk}}$, and $\mathbf{v} \in \mathbb{R}^{D_v}$ are called query, key, and value, respectively, and are obtained from the features \mathbf{x} , typically via linear transformations

$$\mathbf{q} = \mathbf{W}_q \mathbf{x} \quad \mathbf{k} = \mathbf{W}_k \mathbf{x} \quad \mathbf{v} = \mathbf{W}_v \mathbf{x}$$

with trainable weight matrices $\mathbf{W}_q, \mathbf{W}_k \in \mathbb{R}^{D_{qk} \times H}$ and $\mathbf{W}_v \in \mathbb{R}^{D_v \times H}$. The similarity kernel is usually chosen as

$$\text{sim}(\mathbf{q}, \mathbf{k}) = \exp\left(\frac{\mathbf{q}^\top \mathbf{k}}{\sqrt{D_{qk}}}\right). \quad (3)$$

To apply self-attention to structures embedded in Euclidean space, Eq. 2 needs to be modified to also include information about the positions $\mathcal{R} = \{\vec{r}_1, \dots, \vec{r}_N \mid \vec{r}_m \in \mathbb{R}^3\}$ associated with the features \mathcal{X} . A straightforward way to achieve this is to define a geometric version of self-attention

$$\text{ATT}_{\text{Geom}}(\mathcal{X}, \mathcal{R})_m = \frac{\sum_{n=1}^N \text{sim}(\mathbf{q}_m, \mathbf{k}_n, r_{mn}) \mathbf{v}_n}{\sum_{n=1}^N \text{sim}(\mathbf{q}_m, \mathbf{k}_n, r_{mn})}, \quad (4)$$

which uses a modified similarity kernel that explicitly depends on the pairwise distances $r_{mn} = \|\vec{r}_m - \vec{r}_n\|$. The use of distances to encode spatial information is convenient, because the resulting operation is naturally invariant to rigid translations and rotations, and therefore respects physical symmetries. A possible choice would be

$$\text{sim}(\mathbf{q}, \mathbf{k}, r) = \exp\left(\frac{\mathbf{q}^\top \mathbf{k} \cdot f(r)}{\sqrt{D_{qk}}}\right), \quad (5)$$

which is similar to Eq. 3, except that the dot product $\mathbf{q}^\top \mathbf{k}$ is modulated by a (possibly learned) function f of the distance r .

Note that Eq. 4 can also be thought of as a type of MessagePassing layer (see Eq. 1) acting on a fully-connected graph, where the “message” sent

from node n to node m is given by $a_{mn}\mathbf{v}_n$. Given this analogy, it follows directly that the overall complexity of evaluating Eq. 4 is $\mathcal{O}(N^2)$. It would be possible to introduce a cutoff distance r_{cut} and recover $\mathcal{O}(N)$ scaling, but then Eq. 4 would not be able to model global correlations anymore (similar to ordinary MP with a local cutoff), defeating the purpose.

However, it is also possible to derive a linear-scaling version of self-attention²⁶ in an alternative manner by re-writing the similarity kernel (Eq. 3) as a scalar product in an associated (implicit) feature space^{38–43} as

$$\text{sim}(\mathbf{q}, \mathbf{k}) = \boldsymbol{\psi}(\mathbf{q})^\top \boldsymbol{\psi}(\mathbf{k}) \quad (6)$$

via a feature map $\boldsymbol{\psi}$ taking values in a separable Hilbert space of dimensions $D_\psi \leq \infty$. Inserting Eq. 6 into Eq. 2 leads to

$$\text{ATT}(\mathcal{X})_m = \frac{\sum_{n=1}^N \boldsymbol{\psi}(\mathbf{q}_m)^\top \boldsymbol{\psi}(\mathbf{k}_n) \mathbf{v}_n}{\sum_{n=1}^N \boldsymbol{\psi}(\mathbf{q}_m)^\top \boldsymbol{\psi}(\mathbf{k}_n)}, \quad (7a)$$

which can equivalently be written as

$$\text{ATT}_{\text{Lin}}(\mathcal{X})_m = \frac{\boldsymbol{\psi}(\mathbf{q}_m)^\top \sum_{n=1}^N \boldsymbol{\psi}(\mathbf{k}_n) \mathbf{v}_n^\top}{\boldsymbol{\psi}(\mathbf{q}_m)^\top \sum_{n=1}^N \boldsymbol{\psi}(\mathbf{k}_n)}. \quad (7b)$$

In this formulation, it is evident that the sums $\sum_{n=1}^N \boldsymbol{\psi}(\mathbf{k}_n) \mathbf{v}_n^\top$ and $\sum_{n=1}^N \boldsymbol{\psi}(\mathbf{k}_n)$ are identical for each query vector \mathbf{q}_m and therefore only need to be computed once. Thus, the overall complexity of computing self-attention for every input via Eq. 7b is only $\mathcal{O}(N)$ (see Fig. 1C for an illustration that contrasts the evaluation of Eq. 7a with Eq. 7b).

To apply linear-scaling attention to Euclidean data, we would need to re-write the geometric version of self-attention (Eq. 4) in the form of Eq. 7b. This would require re-writing the modified similarity kernel as

$$\text{sim}(\mathbf{q}_m, \mathbf{k}_n, r_{mn}) = \boldsymbol{\psi}(\mathbf{q}_m, \vec{r}_m)^\top \boldsymbol{\psi}(\mathbf{k}_n, \vec{r}_n), \quad (8)$$

but it is non-obvious how to design a feature map $\boldsymbol{\psi}$ that achieves this. In fact, $\mathcal{O}(N)$ scaling seems inherently incompatible with encoding information about all pairwise distances r_{mn} , because computing them directly has quadratic complexity. A possible alternative could be to encode geometric

information in some other way that does not rely on pairwise distances, but then it is unclear how to ensure that physical symmetries, such as translational and rotational invariance, are respected by the resulting mechanism.

Results

In the following, we address the challenges outlined above by proposing Euclidean Fast Attention (EFA). We then contrast EFA with MP on idealized model systems to highlight important differences in their properties and their ability to model global correlations under controlled conditions. Finally, we apply EFA-augmented and standard MPNN architectures to several challenging realistic chemical systems. We show that a correct description of long-range effects, which is only achieved by EFA-augmented models, is crucial for a qualitatively correct description of many practically relevant chemical interactions.

Euclidean fast attention

Inspired by rotary positional encodings (RoPE),⁴⁴ we propose a mechanism to encode positions $\vec{r} \in \mathbb{R}^3$ in Euclidean space into feature vectors \mathbf{x} , which we call Euclidean RoPE (ERoPE). It is given by

$$\text{ERoPE}_{\vec{u}}(\mathbf{x}, \vec{r}) := \mathbf{x} \cdot e^{i\omega\vec{u}\cdot\vec{r}}, \quad (9)$$

where i is the imaginary unit, $\omega \in \mathbb{R}$ is a frequency coefficient and $\vec{u} \cdot \vec{r}$ is the dot product of \vec{u} and \vec{r} . Here, $\vec{u} \in S^2$ is a three-dimensional unit vector (S^2 refers to the set of all points on the surface of a unit sphere centered at the origin). To see the effect of this encoding, consider the scalar product $\langle \mathbf{q}, \mathbf{k} \rangle := \mathbf{q}^\top \mathbf{k}$ between two (real) vectors \mathbf{q} and \mathbf{k} (a crucial component of similarity kernels for attention mechanisms, see Eq. 3). When the vectors have complex entries, the scalar product is instead defined as $\langle \mathbf{q}, \mathbf{k} \rangle := \mathbf{q}^\top \bar{\mathbf{k}}$, where $\bar{\mathbf{k}}$ denotes the entry-wise complex conjugate of \mathbf{k} . After the positions \vec{r}_m and \vec{r}_n have been encoded into \mathbf{q} and \mathbf{k} with Eq. 9, their scalar product is

$$\begin{aligned} \langle \mathbf{q} \cdot e^{i\omega\vec{u}\cdot\vec{r}_m}, \mathbf{k} \cdot e^{i\omega\vec{u}\cdot\vec{r}_n} \rangle &= \left(\mathbf{q} \cdot e^{i\omega\vec{u}\cdot\vec{r}_m} \right)^\top \left(\overline{\mathbf{k} \cdot e^{i\omega\vec{u}\cdot\vec{r}_n}} \right) \\ &= \mathbf{q}^\top \bar{\mathbf{k}} \cdot e^{i\omega\vec{u}\cdot\vec{r}_m} \overline{e^{i\omega\vec{u}\cdot\vec{r}_n}} \\ &= \langle \mathbf{q}, \mathbf{k} \rangle \cdot e^{i\omega\vec{u}\cdot(\vec{r}_m - \vec{r}_n)}. \end{aligned}$$

Since $\vec{u} \cdot (\vec{r}_m - \vec{r}_n)$ is the projection of the vector $\vec{r}_{mn} = \vec{r}_m - \vec{r}_n$ onto \vec{u} , the scalar product now contains geometric information about the relative displacement vector \vec{r}_{mn} .

Unfortunately, while displacement vectors are naturally invariant with respect to translations, the value of the projection $\vec{u} \cdot \vec{r}_{mn}$ depends on the choice of \vec{u} (or equivalently, the choice of coordinate system) and is therefore *not* invariant with respect to rotations. This issue can be resolved by considering the average over all possible choices of \vec{u} , obtained by integrating over S^2 (see Fig. 1B):

$$\frac{1}{4\pi} \int_{S^2} e^{i\omega\vec{u}\cdot\vec{r}_{mn}} d\vec{u} = \frac{\text{sinc}(\omega r_{mn})}{\omega r_{mn}} = \text{sinc}(\omega r_{mn}). \quad (10)$$

The resulting expression only depends on the relative distance $r_{mn} = \|\vec{r}_{mn}\|$ and is therefore invariant with respect to rotations (see [Analytic Solution of the Surface Integral](#) in the Supplementary Information for a derivation of Eq. 10). Instead of using a single fixed frequency ω , it is also possible to encode different components with varying frequencies (corresponding to a mixture of multiple sinc functions), which improves the expressive power of the operation (see [Methods](#) for details).

We now combine ERoPE, integration over S^2 , and ideas from linear-scaling attention (Eq. 7b) to arrive at Euclidean fast attention (EFA). It is given by

$$\text{EFA}(\mathcal{X}, \mathcal{R})_m = \frac{1}{4\pi} \int_{S^2} \phi_{\vec{u}}(\mathbf{q}_m, \vec{r}_m)^\top \sum_{n=1}^N \overline{\phi_{\vec{u}}(\mathbf{k}_n, \vec{r}_n)} \mathbf{v}_n^\top d\vec{u}, \quad (11)$$

where we define the short-hand $\phi_{\vec{u}}(\mathbf{x}, \vec{r}) := \text{ERoPE}_{\vec{u}}(\psi(\mathbf{x}), \vec{r})$ for conciseness and ψ can be any feature map (see Eq. 6). Note that instead of the linear-scaling attention mechanism in Eq. 7b, we base EFA on an *attention-like* mechanism given by

$$\widetilde{\text{ATT}}_{\text{Lin}}(\mathcal{X})_m = \psi(\mathbf{q}_m)^\top \sum_{n=1}^N \psi(\mathbf{k}_n) \mathbf{v}_n^\top, \quad (12)$$

which omits the denominator. This is done because our aim is to model long-range interactions in chemical structures and normalization is not

meaningful in this case: If an atom interacts with many other atoms, we want the effects to be additive (size extensive), whereas with normalization, we would instead compute a (weighted) average effect.

Finally, there is an additional generalisation we can make, which further increases the geometric expressiveness of Eq. 11. So far, we have implicitly assumed that the representations from which queries, keys, and values are computed are themselves rotationally invariant, as is the case in many applications. However, recent MLFFs often use *equivariant features*,^{31,45} which can be thought of as containing additional ‘‘directional information’’ (see Ref. 46 for an in-depth introduction). We can generalize the EFA mechanism as

$$\text{EFA}(\mathcal{X}, \mathcal{R})_m = \frac{1}{4\pi} \int_{S^2} \phi_{\vec{u}}(\mathbf{q}_m, \vec{r}_m)^\top \sum_{n=1}^N \overline{\phi_{\vec{u}}(\mathbf{k}_n, \vec{r}_n)} \mathbf{v}_n^\top \otimes Y(\vec{u}) d\vec{u}, \quad (13)$$

so it is also applicable to equivariant features. Here, ‘ \otimes ’ denotes a tensor product and $Y(\vec{u})$ is a $(\ell_{\max} + 1)^2$ -vector containing all *spherical harmonics* up to degree ℓ_{\max} . While Eq. 11 can only resolve information about pairwise distances, Eq. 13 can also resolve information about the relative *orientation* between the inputs. For invariant input features and $\ell_{\max} = 0$, Eq. 13 simplifies to Eq. 11 (up to a constant factor that depends on the normalization convention used for the spherical harmonics).

The resulting atomic representations are of the form (\hat{r}_{mn} is the unit vector in direction of \vec{r}_{mn})

$$\mathbf{x}_m \simeq \sum_{n=1}^N \mathbf{f}(r_{mn}) \circ \mathbf{x}_n \otimes Y(\hat{r}_{mn}), \quad (14)$$

where \mathbf{f} is a vector-valued radial function of r_{mn} and ‘ \circ ’ denotes element-wise multiplication (see [Analytic Solution of the Surface Integral](#) in the Supplementary Information for a derivation). As such, Eq. 14 resembles the structure of SO(3) convolutions – the basic MP building block found in many equivariant MPNN architectures^{31,47} – but now *without* needing to introduce a local cutoff to achieve linear scaling.

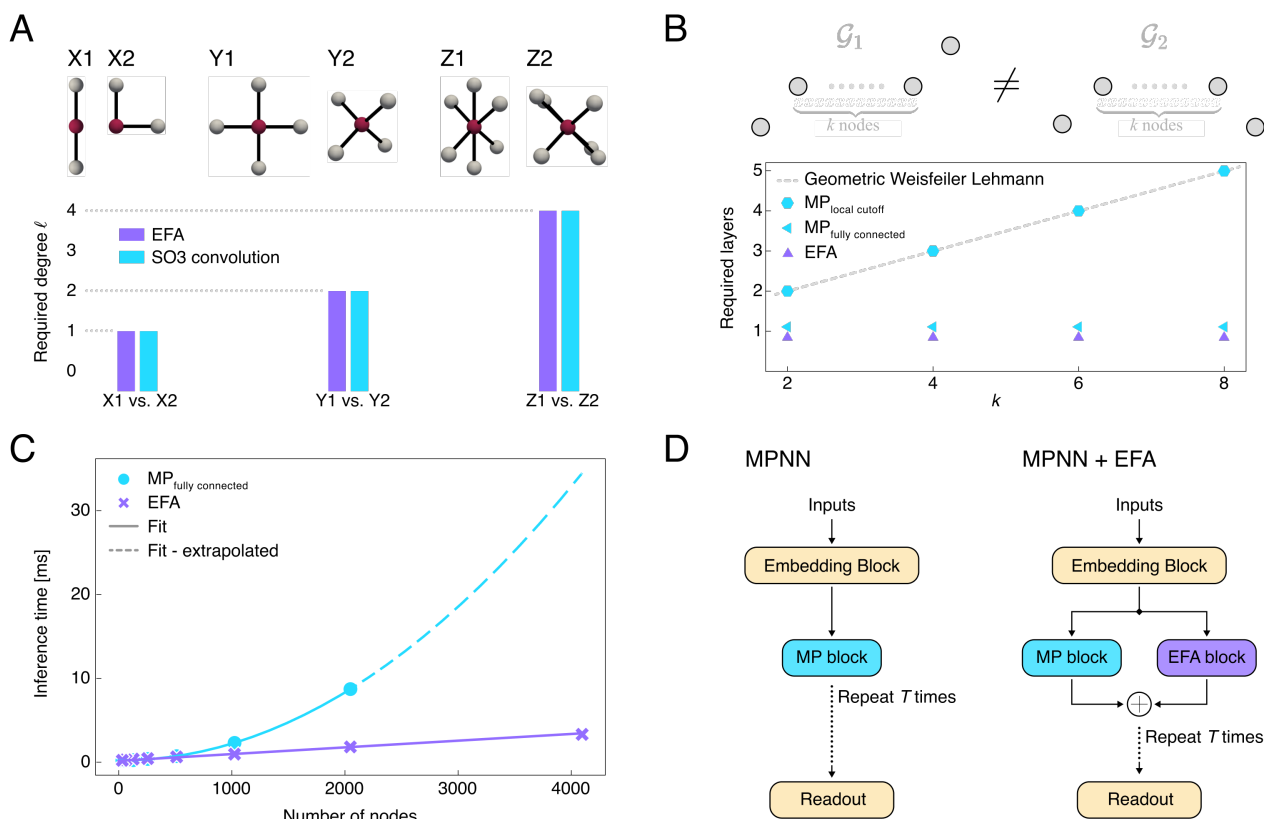


Figure 2 | Geometric expressiveness of Euclidean fast attention (EFA). (A) Smallest degree ℓ of spherical harmonics Y required for EFA (Eq. 13) and SO(3) convolutions (Eq. 14) to distinguish pairs (X1/X2, Y1/Y2, and Z1/Z2) of distinct environments⁴⁸ of the central red node (in a single MP/update). None of the pairs can be distinguished by using invariant ($\ell = 0$) information only (unless multiple updates are used), because all neighbors are equidistant from the central node. (B) Number of MP layers required to distinguish the graphs \mathcal{G}_1 and \mathcal{G}_2 for different chain lengths k .⁴⁹ The graphs are non-isomorphic according to a geometric Weisfeiler-Leman⁵⁰ (GWL) test, due to different orientations of the terminal nodes w. r. t. each other. With MP (using a local cutoff), at least $T = \lfloor k/2 \rfloor + 1$ layers are required to solve the task, following the theoretically expected behavior from the GWL test (grey dashed line). In contrast, with EFA, \mathcal{G}_1 and \mathcal{G}_2 can always be distinguished with a single update (similar to MP on a fully connected graph). (C) Inference time for EFA and MP on a fully connected graph (without using a local cutoff) as function of the number of nodes. MP has quadratic complexity and runs out of memory for $N > 2048$ nodes, whereas EFA has linear complexity and enables scaling to tens of thousands of nodes (see Fig. S8). (D) Comparison of a standard MPNN and an MPNN augmented with an EFA block. The design of EFA enables plug-in integration into existing models with only minimal architectural modifications.

Idealized systems

Before applying EFA to realistic molecular systems, we evaluate and contrast it with standard MP on idealized model systems with known properties. This is done to empirically demonstrate some of the shortcomings of local models in a controlled environment, and to verify that the EFA mechanism mitigates these issues without sacrificing performance.

Geometric expressiveness First, we investigate the geometric expressiveness of EFA and MP, i.e., their ability to resolve geometric information. When distinguishing pairs of local atomic neighborhoods,⁴⁸ we empirically confirm the similarity between EFA and SO(3) convolutions (see Eq. 14) and find that increasing the degree ℓ in Y_ℓ has the same effect for both operations, allowing to resolve differences in increasingly complex structures (see Fig. 2A, details of the experimental

setup can be found in [Methods](#)). While both EFA and MP are equally expressive in this setting, the difference between them becomes apparent when investigating their ability to distinguish molecular graphs as a whole. To this end, a recent work⁴⁹ proposed to test whether models can distinguish two graphs that are non-isomorphic according to a geometric Weisfeiler-Leman⁵⁰ (GWL) test. Both graphs consist of a chain with k nodes and only differ in the relative orientation of additional nodes at both ends of the chain (see [Fig. 2B](#) for an illustration). As expected, when nodes are far enough apart that information can only “hop” between direct neighbors (see also [Fig. 1A](#)), the minimum number of MP layers required to distinguish both graphs is $T = \lfloor k/2 \rfloor + 1$ (following the theoretically expected behavior from the GWL test). In contrast, when MP is augmented with EFA, both graphs can always be distinguished with only a single update, irrespective of chain length and distance between nodes. Of course, this is also possible with a single MP layer (without EFA) by getting rid of the local cutoff (so the graph becomes fully-connected), but results in quadratic complexity. The crucial distinction is that EFA achieves this feat with linear complexity, allowing to scale the method to tens of thousands of nodes (see [Fig. 2C](#)). Note that the comparisons between MP and EFA use essentially the same underlying MPNN backbone in both cases, with the only difference being an additional EFA block that gets applied during each update (see [Fig. 2D](#) for an illustration and [Methods](#) for details). Because adding an EFA block to an existing model architecture requires only minimal modifications, we expect it will be straightforward to extend other local models with EFA in a similar manner.

Pairwise potentials Next, we consider model systems that only interact via pairwise potentials of the form

$$V(r) = \frac{1}{r^3} - \frac{c}{r^b}, \quad (15)$$

where the first term is a short-ranged repulsion and the latter a long-ranged interaction. The coefficient c determines the relative strength of the two terms (and whether the long-range term is repulsive or attractive) and the exponent b controls the decay behavior of the potential. The

total energy of a system of N atoms interacting via [Eq. 15](#) is given by

$$E(\vec{r}_1, \dots, \vec{r}_N) = \sum_m^N \sum_{n>m}^N V(r_{mn}), \quad (16)$$

where $r_{mn} = \|\vec{r}_m - \vec{r}_n\|$ is the distance between atoms m and n .

We start with the simplest case of two particles ($N = 2$) and set $c = b = 1$. With this choice for b , the long-range decay behavior is reminiscent of (attractive) charge-charge interactions. As expected, a standard MPNN predicts a qualitatively wrong, constant energy profile as soon as the particle separation exceeds the local cutoff. When the model is augmented with EFA, the pairwise potential can be accurately described over the full interaction length (see [Fig. 1A](#)). For systems consisting of more than two atoms, the additional atoms can potentially act as “information relays”, and the resulting increased effective cutoff (see [Fig. 1A](#)) may allow MPNNs to model the potential accurately, even without EFA. To test this hypothesis, we consider systems of size $N = 16$ and $N = 32$, where atoms are randomly distributed within spheres of diameter $d = 5 \text{ \AA}$ and $d = 10 \text{ \AA}$, respectively ([Fig. 3A](#)). Since we use a cutoff of $r_{\text{cut}} = 5 \text{ \AA}$, all nodes are direct neighbors of each other in the smaller system ($N = 16$) by construction. For the larger system ($N = 32$), we verify that there are no disconnected sub-graphs for the given cutoff ([Fig. S4](#)), such that information from each node can reach any other node with a sufficiently large number of MP updates. The $N = 16$ system acts as a baseline here: Since all nodes are direct neighbors, a standard MPNN is expected to be able to learn this task without issues, even with just a single MP layer. For $N = 32$, although a large fraction of nodes are still direct neighbors of each other (see [Fig. 3A](#)), there are also many nodes separated by two or three “hops” that can only interact after the corresponding number of MP updates (once the effective cutoff is sufficiently large). It seems reasonable to assume that an MPNN would become more accurate on this system when increasing the number of MP layers, because it should be able to (indirectly) model an increasing number of interactions between nodes beyond the local

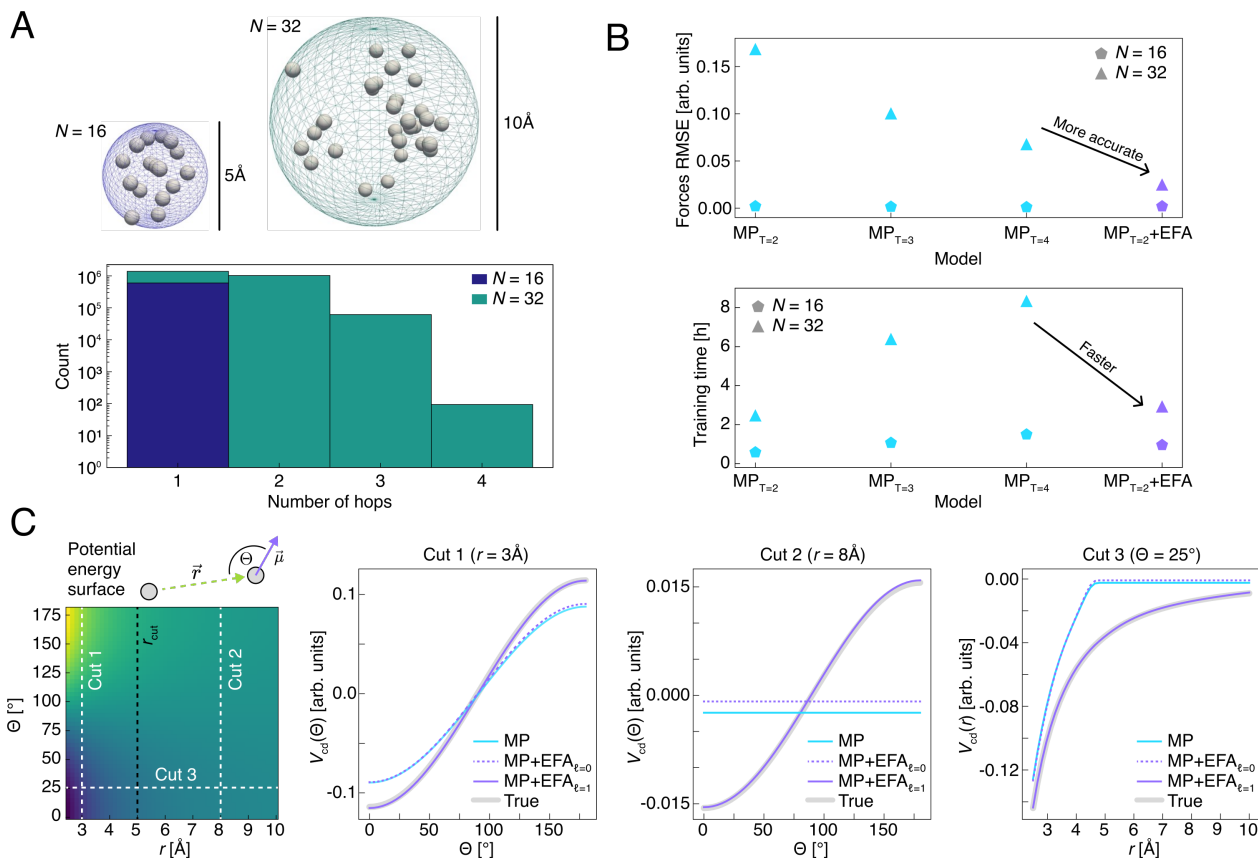


Figure 3 | Euclidean fast attention (EFA) overcomes limitations of message passing (MP). (A) Visualization of idealized model systems with $N = 16$ and $N = 32$ atoms randomly distributed in spheres with a diameter of $d = 5\text{Å}$ and $d = 10\text{Å}$, respectively. The histogram shows the distribution of the number of “hops” between all pairs of nodes in the data (2.5k structures per system) for a cutoff of $r_{\text{cut}} = 5\text{Å}$. (B) Root mean square error (RMSE) for force predictions and total training times for standard MPNNs with different numbers of layers ($T = 2, 3, 4$) and an MPNN ($T = 2$) augmented with EFA. For the $N = 32$ system, the model with EFA is significantly more accurate than standard MPNNs that are up to twice as deep, while at the same time being faster to train. For $N = 16$ differences between the models are negligible. (C) Standard MP and MP+EFA _{ℓ} applied to a non-isotropic interaction (ℓ refers to the maximal degree of the spherical harmonics vector \mathbf{Y} in Eq. 13). Energy predictions along different cuts of the two-dimensional potential energy surface (see leftmost panel for reference) are shown. Only MP+EFA _{$\ell=1$} has access to information about both relevant degrees of freedom (distance and orientation) and can therefore faithfully describe the interaction beyond the cutoff r_{cut} .

cutoff (but within the effective cutoff). We compare the performance of MPNNs with different numbers of layers ($T = 2, 3, 4$) and an MPNN augmented with EFA ($T = 2$) by analysing the root mean squared error (RMSE) between predicted and ground truth forces. As expected, all models yield highly accurate results on the $N = 16$ system. For the $N = 32$ system, standard MPNNs indeed become more accurate when the number of MP layers is increased. However, the most accurate model is the MPNN augmented with EFA, despite using only $T = 2$ layers, which is the same depth

as the shallowest (and worst-performing) of the non-augmented MPNNs we tested (see Fig. 3B).

These results suggest that “mean-field interactions” between node neighborhoods are not sufficient to capture all relevant details of interatomic potentials (not even in this simple toy model) and augmenting models with EFA is beneficial in practice – even if the effective cutoff of a local model is large enough to capture all interactions in theory. As an additional remark, we note that the overall effect of increasing the number of MP

layers goes beyond merely increasing the size of the effective cutoff. For example, the number of node pairs separated by four hops makes up less than 0.003 % of the dataset for the $N = 32$ system (see Fig. 3A). Consequently, if increasing the number of MP layers had no effect other than increasing the effective cutoff, we would expect nearly identical performance for the MPNNs with three ($T = 3$) and four ($T = 4$) layers. However, the MPNN with $T = 4$ achieves a significantly lower RMSE, demonstrating that more MP layers also improve representations for interactions that are already within the effective cutoff. Increasing the accuracy by adding more MP layers is not for free though, as it also increases inference cost and training time (see Fig. 3B). Thus, adding an EFA block to MPNNs does not only allow them to model long-ranged interactions, rather, it generally increases their accuracy, but does so more efficiently than adding more MP layers: The MPNN+EFA model with $T = 2$ is more accurate than a standard MPNN that is twice as deep, while at the same time being roughly four times faster to train on the system with $N = 32$ atoms (see Fig. 3B). We expect this trend to be even more pronounced in larger and more complex real-world systems.

Given the similarity of the pairwise potential used in the previous tests to charge-charge interactions, we also consider a variation of the $N = 32$ system, where atoms are randomly assigned a “charge” value $q \in (-2, -1, +1, +2)$. The interaction coefficient in Eq. 15 is then made charge-dependent (i.e., $c = q_m q_n$), so the strength of each pairwise interaction (and even whether it is attractive or repulsive) is non-uniform across different atom combinations. In this more complicated setup, MP+EFA ($T = 2$) still outperforms standard MP for all numbers of layers ($T = 2, 3, 4$), demonstrating that EFA is also able to learn interaction patterns that are determined by atom-type specific properties (see Supplementary Fig. S4C).

Finally, we design an experiment to reveal the shortcomings of purely distance-based descriptions and show the utility of the equivariant version of the EFA mechanism (Eq. 13). We go back to a two particle system, but make the pairwise potential V (Eq. 15) non-isotropic, i.e., it now not

only depends on distance, but also orientation. For this, we modify V to mimic charge-dipole interactions: One atom is assigned a “charge” $q = -1$ and the other a unit “dipole” vector $\vec{\mu} \in \mathbb{R}^3$. The interaction coefficient in Eq. 15 is then given by $c = \cos(\Theta)$, where Θ is the angle between $\vec{\mu}$ and the normalized displacement vector \vec{r}_{mn}/r_{mn} (the decay parameter is set to $b = 2$, reflecting the physical behavior of charge-dipole interactions). We then compare the performance of a standard MPNN with two variants augmented with EFA, one using degree $\ell = 0$ and the other using $\ell = 1$ in the spherical harmonics vector Y_ℓ (referred to as EFA $_{\ell=0}$ and EFA $_{\ell=1}$ in the following). As expected, all models are able to accurately describe the energy profile up to the local cutoff. Beyond the cutoff, however, both MP and MP+EFA $_{\ell=0}$ predict a constant value, and only MP+EFA $_{\ell=1}$ successfully captures the true behavior of the underlying potential (Fig. 3C). The reason why both MP and MP+EFA $_{\ell=0}$ fail here is because they are effectively “blind” to the relevant degrees of freedom. For MP, no geometric information at all is accessible beyond the cutoff, whereas for MP+EFA $_{\ell=0}$, only distance information is visible. This is not sufficient, because the interaction potential depends on the angle Θ . Therefore, both of these models essentially become mean predictors at long range. Only MP+EFA $_{\ell=1}$ has access to information about both distance *and* orientation beyond the cutoff and is therefore able to learn the correct long-range behavior.

Molecular Systems

Our results on Idealized systems provide evidence that standard MPNNs have fundamental shortcomings which can be mitigated by augmenting them with EFA. In the following, we investigate whether this also translates to increased performance on realistic data. We study several molecular systems exhibiting various long-ranged and non-local effects, which are representative for common interaction patterns found throughout chemistry. We demonstrate that EFA systematically improves the performance of local MPNNs and that an accurate description of long-range interactions can be crucial for predicting the dynamic properties of molecular systems.

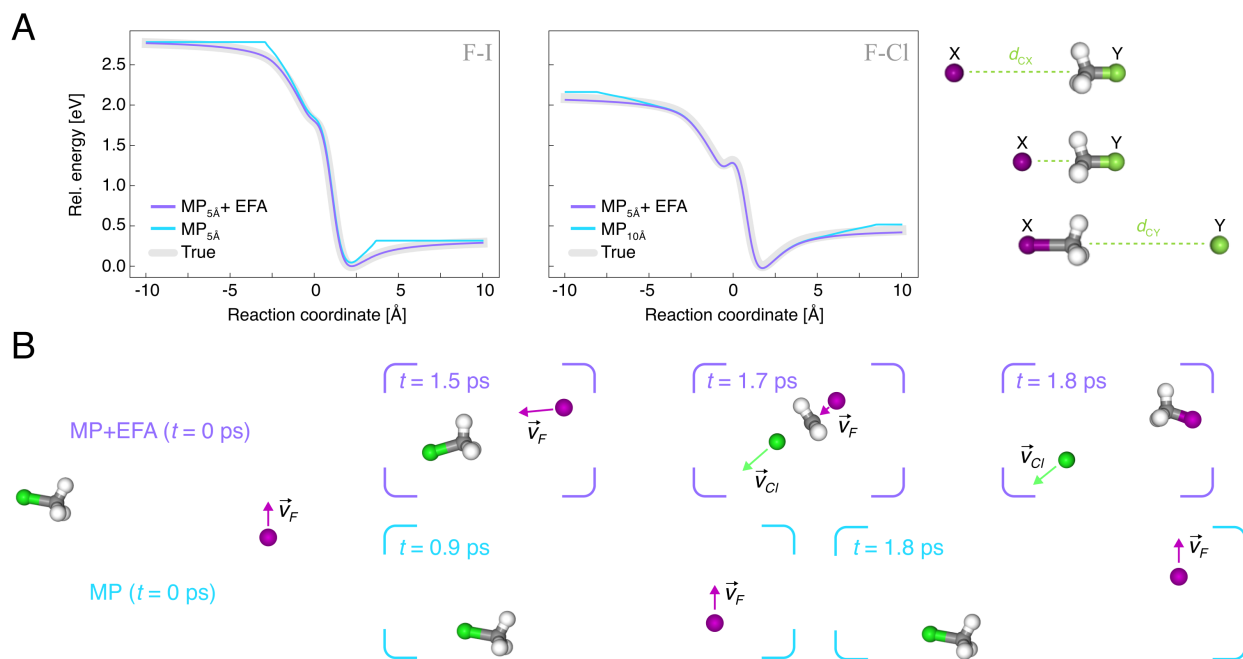


Figure 4 | **Euclidean fast attention (EFA) for reactions.** (A) Illustration of an S_N2 reaction of the form $X^- + H_3C-Y \rightarrow X-CH_3 + Y^-$. The panels show a one-dimensional cut of the potential energy surface (PES) along the reaction coordinate ($d_{CY} - d_{CX}$) for two systems ($X=F, Y=I$ and $X=F, Y=Cl$). A model using standard message passing (MP) with a local cutoff of $r_{cut} = 5 \text{ \AA}$ shows qualitatively wrong asymptotic behavior, which also cannot be fixed by increasing the cutoff to $r_{cut} = 10 \text{ \AA}$. In contrast, models augmented with EFA reproduce the correct energy profile for all particle separations. (B) Snapshots at different times t of two MD trajectories starting from identical initial conditions simulated with the MP+EFA model (top) and MP model (bottom), respectively. The velocities of fluorine (\vec{v}_F) and chlorine (\vec{v}_{Cl}) atoms are highlighted by arrows to indicate their motion. Since the MP model predicts no forces between ion and molecule beyond the cutoff, the reactants fly past each other undisturbed and no reaction occurs. In contrast, for MP+EFA, the reactants are attracted towards each other and the methyl halide molecule is being re-oriented and properly positioned for backside attack of the ion, facilitating the reaction.

S_N2 Reactions For many chemical reactions to proceed, the involved molecules must come into close contact. In addition, the reactants often need to be oriented in a specific way when they meet, so that the reaction can occur. As an example, we focus on prototypical S_N2 reactions of the type $X^- + H_3C-Y \rightarrow X-CH_3 + Y^-$, where X and Y are halogens (F, Cl, Br, I).⁵¹ Such systems exhibit strong long-ranged electrostatic interactions between the methyl halide molecule (which has a large dipole moment) and the halide ion (carrying a negative charge), which makes predicting energies and forces challenging for local models. We find that adding an EFA block to a local MPNN architecture ($r_{cut} = 5 \text{ \AA}$) yields a 34× and 8× reduction in mean absolute errors (MAEs) for energy and force predictions, respec-

tively (Tab. S1). Interestingly, merely increasing the cutoff distance of the MPNN does not lead to satisfactory improvements: Even with a cutoff distance of $r_{cut} = 10 \text{ \AA}$, a standard MPNN has 18× and 2× larger MAEs for energies and forces compared to the EFA-augmented MPNN with $r_{cut} = 5 \text{ \AA}$ (Tab. S1). The reason for the substantially lower MAEs when using EFA becomes apparent when visualizing one-dimensional cuts of the potential energy surface along the reaction coordinate (see Fig. 4A). Both models without EFA predict the wrong asymptotic behavior and show unphysical artefacts for large distances between ion and molecule. As soon as their separation exceeds the cutoff, local MPNNs predict a constant energy value. In contrast, the model with EFA accurately describes the energy profile for the

full range of particle separations. The unphysical artefacts in the long-range region of the potential energy surface also lead to qualitatively incorrect dynamic behavior. To demonstrate this, we run MD simulations of $\text{CH}_3\text{Cl} + \text{F}^-$ using both the MP+EFA and MP models starting from identical initial conditions. We find that only the trajectory driven by the MP+EFA model leads to a reaction resulting in $\text{CH}_3\text{F} + \text{Cl}^-$, whereas the reactants erroneously fly past each other undisturbed in the trajectory predicted by the MP model without EFA (see Fig. 4B).

Electronic delocalization So far, our analyses have focused primarily on long-range interactions that can be described as functions of atomic (or molecular) distances r , usually decaying proportionally to some power law r^{-b} . While these kinds of effects are arguably among the most important long-ranged interactions in molecules, some chemical systems exhibit *non-local* effects, which are more complicated in nature. They are typically due to electronic delocalization and cannot easily be described as a function of distance. Instead, they are usually characterized by a strong dependence of the energy on the relative orientation of substructures in a molecule. As implied by their name, non-local effects are particularly challenging to describe for local models, and contrary to most long-ranged interactions, it is also difficult to account for them with empirical correction terms. Here we consider the strong non-local effects in cumulene molecules, which have been found to be particularly challenging for both, global and local MLFFs.^{6,45}

In cumulenes, the energy strongly depends on the dihedral angle Θ between the hydrogen rotors at the ends of a (in theory almost arbitrarily long) chain of carbon atoms (see Fig. 5A for an illustration). We compare the energy profile of a $T = 3$ layer MPNN augmented with EFA using different maximum degrees $\ell = 0, 1, 2$ for the spherical harmonics with regular MPNNs with $T = 3$ and $T = 5$ layers. EFA $_{\ell=0}$ uses only invariant features, whereas EFA $_{\ell=1}$ uses equivariant representations up to degree $\ell = 1$ (equivalent to vector representations as, e.g., used in PaiNN⁵²) and EFA $_{\ell=2}$ goes up to degree $\ell = 2$. As expected

from our investigations on **Idealized systems**, the MP+EFA $_{\ell=0}$ model cannot describe the energy profile, because it only leverages pairwise distances, which cannot resolve the change in the dihedral angle sufficiently well.⁴⁵ As soon as equivariant representations are allowed to enter the EFA update ($\ell > 0$), the energy is modelled faithfully, with the prediction of the energy barrier increasing in accuracy when transitioning from $\ell = 1$ to $\ell = 2$. MP models without EFA always fail (even when they are made equivariant) and predict a flat energy profile as soon as the distance between hydrogen rotors exceeds the effective cutoff. Of course, it is possible to increase the effective cutoff by increasing the number of layers to $T = 5$, but even then the barrier height is still underestimated (inset Fig. 5A), indicating over-squashing of information.⁵³ Even some models without a local cutoff, such as sGDML⁵⁴ or a fully connected SchNet¹² have problems with this task, as they are unable to resolve relative orientations over large distances (see Fig. S5B). Modelling the energy barrier correctly is crucially important for predicting the correct dynamical behavior during MD simulations. For example, models that predict a flat energy profile incorrectly predict that all possible dihedral angles between the terminal CH_2 rotors are sampled equally likely (see Fig. 5B). This in turn leads to wrong predictions of physical observables (see Fig. 5C).

We remark that all tested models are unable to reproduce the sharp cusp in the energy profile at a dihedral angle of 180° . However, the cusp is an artefact due to the inability of the reference *ab initio* method to describe the non-adiabatic couplings at the conical intersection between two potential energy surfaces. A cusp like this leads to discontinuous forces, which would result in unstable MD simulations. For this reason, most MLFFs (including all models tested here) make discontinuous forces impossible by design: “Smoothing” of the problematic region in the energy profile is actually a desirable feature, as it mimics the appearance of a (physically correct) diabatic surface.

Dimers As a final representative system, we probe non-covalent interactions in dimer systems

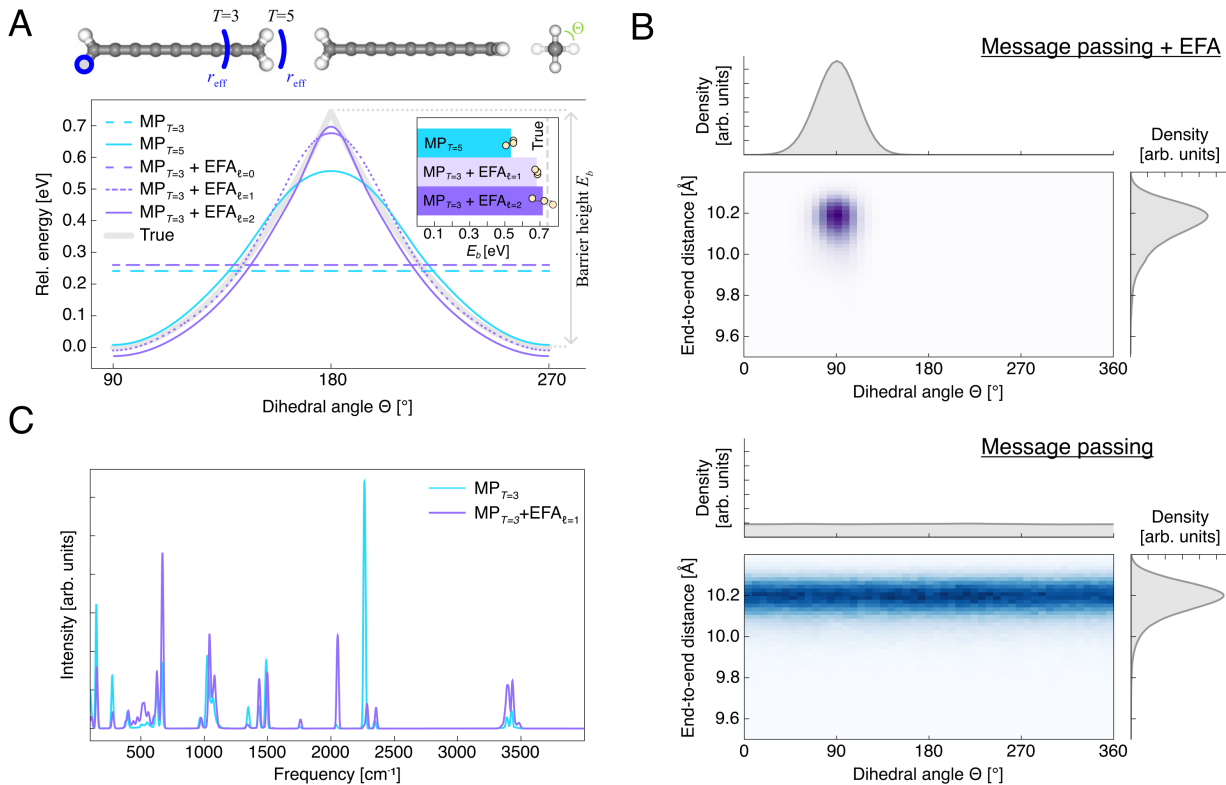


Figure 5 | **Euclidean fast attention (EFA) for electronically de-localized effects.** (A) Cumene molecule $C_{k+2}H_4$ with chain length $k = 7$ (the effective cutoff r_{eff} from the atom marked in blue is indicated for different numbers of MP layers T). The panel shows the energy profile as function of the dihedral angle Θ between the terminal CH₂ rotors for MP ($T = 3, 5$) and MP+EFA ($T = 3$) with different maximal degrees $\ell = 0, 1, 2$. The inset shows the predicted energy barrier height (only for models that do not predict a flat energy profile), yellow dots corresponds to models trained with different random seeds. (B) Visualization of the space spanned by the end-to-end distance (measured as the distance between the outer carbons) and the dihedral angle between the CH₂ rotors visited during 2 ns molecular dynamics (MD) simulations at 300 K in the NVT ensemble. Due to the large energy barrier in the ground truth potential energy surface at a dihedral angle of 180°, physically accurate trajectories starting from the minimum at 90° are expected to only deviate from this value slightly during simulations (there is not enough energy to cross the barrier at a temperature of 300 K). The MP+EFA model (top) follows this expectation and the dihedral angle stays around 90°, whereas the MP model predicts wrong dynamic behavior (all dihedral angles are sampled equally likely, consistent with the prediction of a flat energy profile). (C) Predicting the wrong dynamic behavior has direct consequences on experimental observables: For example, the power spectrum extracted from the dynamics driven by the local MP model exhibits a strong spurious peak at a frequency of ~ 2300 cm⁻¹.

sampled from the DES370K data set.⁵⁵ Here, the interaction energy can be dominated by contributions from electrostatics, induction, dispersion, or mixtures thereof, and also depend on the relative orientation between the molecules. To perform well on this challenging benchmark, ML models need to generalize across the functional form of different long-range interactions and molecular structures. Consistent with our previous findings, MP fails to model dimer interaction energies

as soon as the separation between the individual molecules exceeds the local cutoff, whereas MP+EFA captures the energy profile faithfully (Fig. 6A).

To probe the generalization capabilities of the EFA mechanism, we also evaluate the trained MP+EFA model on four completely unseen dimers (not in the training data) and find that the predicted energy profiles agree closely with

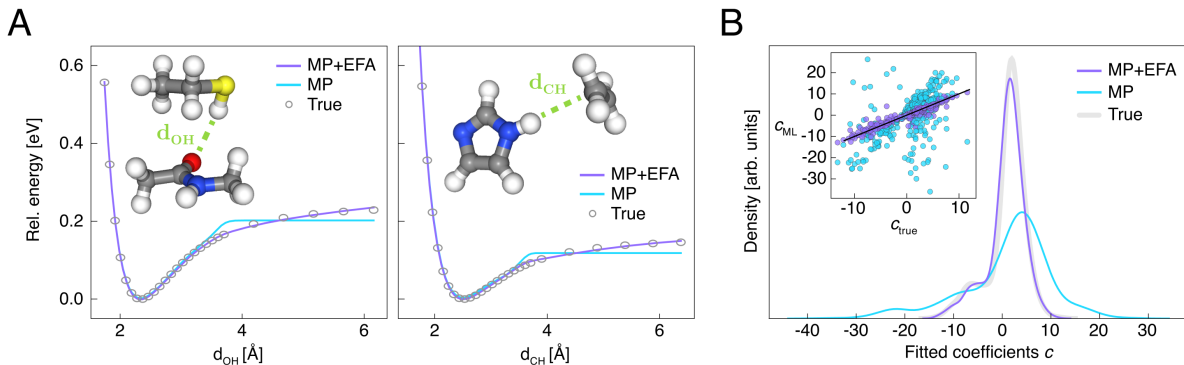


Figure 6 | **Dimers.** (A) Binding curves for different dimers from the DES370K data set⁵⁵ for a standard MP model and for a model augmented with EFA (ground truth reference values are shown with circles). (B) Correlation (inset) and distribution of coefficients c_1, \dots, c_6 for long-range dimer interactions (see Eq. 17). The Pearson correlation w. r. t. the coefficients fitted to the ground truth is $s = 0.56$ for MP and $s = 0.95$ for MP+EFA. Coefficient distributions show a similar trend, indicating an improved agreement to the ground truth for MP+EFA compared to MP.

the ground truth (Fig. S7A). This indicates that EFA enables learning general long-range interaction patterns that are transferable across different chemical systems.

To further quantify the agreement of different models w. r. t. to the ground truth reference, we fit the long-range behaviour of the inter-molecular interactions to the functional form

$$V_{lr}(r_{AB}) = \sum_{i=1}^6 c_i r_{AB}^{-i}, \quad (17)$$

where c_i are expansion coefficients and r_{AB} is the distance between molecules A and B . Eq. 17 can capture the behavior of many physically relevant interaction types, e.g., charge-charge (r_{AB}^{-1}), charge-dipole (r_{AB}^{-2}), or dispersion (r_{AB}^{-6}). Which terms are dominant for a particular dimer interaction depends on the molecules involved and their relative orientation. We can then compare the fitted coefficients for predicted energy curves with those for the ground truth to quantify the agreement in long-range decay behavior. This measure is chosen because (especially in the asymptotic tail) inter-molecular interactions can be relatively weak, so standard aggregate metrics like MAE or RMSE are not adequate to resolve small qualitative differences, which are nevertheless important for the correct description of many chemical systems.¹⁹ We calculate the Pearson correlation s of coefficients for model predictions with those

for the ground truth and find $s = 0.56$ for the standard MP model vs. $s = 0.95$ for the MP+EFA model. Further, we find excellent agreement between the coefficient distribution for the ground truth and the MP+EFA model, whereas the MP model yields a qualitatively wrong distribution (Fig. 6B). We also investigate correlations and distributions for individual c_i with similar trends (Fig. S7C). These results indicate that augmenting local models with EFA enables them to describe a wide range of molecular interactions with different decay behaviors accurately and reliably, whereas standard models fail to do so.

Discussion and Conclusion

Many existing machine learning force fields (MLFFs) introduce a local cutoff radius when modelling atomic interactions to ensure linear scaling w. r. t. the number of atoms. While a cutoff is desirable from a computational efficiency standpoint, it prevents the correct description of *global* interactions, which extend beyond the chosen cutoff distance, by construction. Although such long-ranged contributions (e.g., van der Waals forces) are often weak in magnitude compared to short-ranged interactions (e.g., covalent bonds), they can nonetheless be crucial for the correct treatment of large structures such as proteins, biomolecules, or solids.¹⁹ New methods for treating

global interactions, preferably without compromising computational efficiency, are thus a necessary requirement for successfully scaling MLFFs to larger chemical systems of practical interest.

Different approaches have been proposed to overcome the limitations of strictly local MLFFs. They can be broadly categorized into two classes: Methods that (i) augment otherwise local models with global correction terms or (ii) abandon the concept of strict locality altogether by instead learning global representations. The first category typically relies on a physically motivated description of specific long-range effects.^{30,51,56–64} These approaches usually involve learning intermediate, local quantities, which are then used as parameters in fixed interaction terms. The computational complexity of evaluating these terms often scales (at least) quadratically with the number of atoms, but has a relatively small pre-factor, which allows their application to moderately large systems. Examples include learning partial charges, which are used to calculate electrostatic interactions,^{30,51,58,59} the prediction of atomic electronegativities, which are used as inputs to a charge equilibration scheme and correct for non-local charge transfer,⁶⁰ or trainable polarizability volumes used to correct for dispersion interactions.^{56,57} All of these methods define fixed functional forms, implicitly assuming that all relevant long-range interactions can be modelled this way. While this may be enough for some systems with known types of interactions, recent studies highlight that complex systems, such as proteins, can display long-range information exchange with unique and non-trivial interaction strengths.⁶⁵ In contrast, approaches in the second category instead aim to learn the correct interaction patterns directly from data without relying on domain knowledge or making assumptions about the nature of the underlying interactions. However, achieving this without sacrificing computational efficiency is difficult and existing approaches are typically tailored to very specific systems and interaction types^{30,45,66} or assume a fixed reference frame,^{67,68} which limit their general applicability.

To address this challenge, our work proposes Euclidean fast attention (EFA), which allows to learn

global atom/node representations for molecules (and other graphs embedded in Euclidean space) with linear time and memory complexity. A core component of EFA are Euclidean rotary positional encodings (ERoPE), which are combined with a linear-scaling attention-like mechanism (and integration over the unit sphere S^2) to arrive at the EFA method, which respects all relevant physical symmetries, such as translational invariance and rotational equivariance.

As a representative local model for constructing MLFFs, we compare standard (equivariant) message passing neural networks (MPNNs) with models that combine message passing (MP) with our EFA mechanism. The required architectural modifications are minimal, as EFA blocks can be readily added to other state-of-the-art local MLFFs.

A natural question to ask is whether it is possible to build accurate models purely out of EFA blocks – without also relying on methods specialized on local interactions, such as MP with a cutoff. We believe this to be difficult within the current implementation of EFA, and that further algorithmic improvements would be necessary to make this feasible. The reason is that the particular functional form of the ERoPE mechanism (Eq. 9) induces oscillatory Bessel functions as a radial basis (see [Analytic Solution of the Surface Integral](#) in the Supplementary Information for details). In practice, because we use a Lebedev quadrature⁶⁹ to perform the integration over the unit sphere S^2 (see [Methods](#)), the frequency with which these Bessel functions oscillate must be limited to maintain high computational efficiency of the EFA mechanism. Consequently, EFA is best suited to model “low frequency” functions that vary slowly with distance. Our results demonstrate that this seems to be no fundamental issue for modelling non-local and long-ranged interactions can be observed, but it is likely to cause problems when describing strong short-ranged interactions, which can change rapidly for small changes in distance (e.g., exchange repulsion at small nuclear separations). Thus, for best results, the current form of EFA should always be coupled with other methods specialized in modelling short-range effects. Nonetheless, we believe advancements in alternative integration methods

could allow removing this frequency limitation from EFA, which in turn would ultimately allow building ML architectures based on a unified mechanism for modelling both short- and long-ranged interactions. Future work will thus focus on augmenting existing local architectures with EFA to increase their accuracy for long-ranged interactions and on developing faster numerical integration methods for the special structure of the surface integral in Eq. 13.

In our work we have demonstrated on idealized chemical model systems that ordinary MPNNs have systematic shortcomings, which limit their accuracy and prevent them from capturing global and long-ranged correlations. When MPNNs are augmented with EFA, we find that these shortcomings are successfully mitigated and models now become enabled to learn global interactions. EFA is furthermore applied to representative more realistic chemistry in order to demonstrate its broad utility. The systems studied encompass reactions, electronic delocalization and dimers. For all cases studied, we find that an accurate modeling of global interactions is indispensable. Specifically, we could show that EFA leads to significantly improved performance, both by increasing overall accuracy, and by fixing artefacts in the potential energy surface, all of which otherwise leads to qualitatively incorrect physical behavior.

Finally, we would like to stress once more the high computational efficiency of the EFA framework: it provides a versatile linearly scaling technique (plug-in) for modeling and exploring complex short and long-ranged (global) interactions throughout chemistry and beyond.

Methods

Euclidean Fast attention (EFA)

In the following, we only describe the implementation of EFA in its most general form (Eq. 13), which assumes equivariant features (the mechanism for invariant features in Eq. 11 can be recovered as a special case). We first describe all individual components and then show how they are combined to the full EFA mechanism.

Equivariant features On a high level, equivariant features can be thought of as containing additional “directional information”. Under transformations of the coordinate system, the numerical values of equivariant features may change in a complicated manner, but they still encode “the same” directional information (but transformed accordingly). Here, we briefly describe a particular type of equivariant features introduced in Ref. 46 (for a more detailed overview, we refer the reader to Ref. 46).

The equivariant features we consider consist of irreducible representations (irreps) of the orthogonal group in three dimensions $O(3)$ (rotations and reflections). Features of degree ℓ are denoted as $\mathbf{x}^{(\ell)} \in \mathbb{R}^{P \times (2\ell+1) \times H}$, where P is the size of the “parity axis” (either 1 or 2) and H is the feature space dimension. The parity of an irrep can be either even (+1) or odd (−1) and determines how it behaves under reflections (it either changes sign or not), whereas the degree ℓ determines the behavior under rotations (an irrep of degree ℓ has $2\ell+1$ components/is described by $2\ell+1$ numbers). For example, irreps with $\ell = 0$ are described by a single number and do not change (are invariant) under rotations, whereas irreps with $\ell = 1$ consist of three components and rotate similar to ordinary three-dimensional vectors. When irreps have parity +1 for even ℓ and −1 for odd ℓ , we refer to them as “tensors”, and when they have the opposite parity, we call them “pseudotensors”. Features either consist only of tensors (in which case $P = 1$), or of both tensors and pseudotensors (in which case $P = 2$ and irreps of even/odd parity are stored in separate slices, i.e., indices 0 and 1, of the parity axis).

We are usually working with the concatenation of features starting from the lowest degree $\ell = 0$ up to some maximum degree L , which we denote as $\mathbf{x} \in \mathbb{R}^{P \times (L+1)^2 \times H}$ (each degree ℓ contributes a slice of size $2\ell + 1$, giving a total size $(L + 1)^2$ for the “degree axis”). We use \mathbf{x}_m to refer to the representation of the m -th atom and $\mathbf{x}^{(\ell_{\pm})}$ to refer to the “slice” of irreps of degree ℓ with parity $p = \pm 1$. Invariant features, for example those used in SchNet¹² or PhysNet,⁵¹ can be considered as a special case with $L = 0$ and $P = 1$, i.e., they only consist of irreps with degree 0 and even

parity.

Two irrep representations \mathbf{x} and \mathbf{y} can be ‘‘coupled’’ via tensor product contractions⁷⁰ to produce new features \mathbf{z} . The irreps of degree c and parity γ of the new features \mathbf{z} are given by

$$\mathbf{z}^{(c_\gamma)} = \sum_{(a_\alpha, b_\beta)} \mathbf{x}^{(a_\alpha)} \otimes^{(c_\gamma)} \mathbf{y}^{(b_\beta)}, \quad (18)$$

where the sum runs over all combinations of irreps with degrees $a + b = c$ and parities $\alpha \cdot \beta = \gamma$ of the input features \mathbf{x} and \mathbf{y} . Evaluating the ‘ $\otimes^{(c_\gamma)}$ ’ operation in Eq. 18 involves a summation over components of the tensor (outer) product weighed with so-called Clebsch-Gordan coefficients (see Ref. 46 for implementation details).

Performing all possible tensor product contractions between irreps of maximal degree L_x and L_y up to a maximal output degree $L_z \leq L_x + L_y$ is written as

$$\mathbf{z} = \mathbf{x} \bigotimes_{L_z}^{L_x L_y} \mathbf{y}, \quad (19)$$

such that $\mathbf{z} \in \mathbb{R}^{P \times (L_z + 1)^2 \times H}$. In other words, this is just a compact notation for the concatenation of all tensor product contractions (Eq. 18) for $c = 0, \dots, L_z$ and $\gamma = \pm 1$.

Euclidean rotary positional encodings (ERoPE)

Following typical implementations of the RoPE mechanism,⁴⁴ ERoPE is implemented in a slightly different (but mathematically equivalent) manner than what is suggested in Eq. 9. Consider a complex number $c = a + bi$ ($c \in \mathbb{C}$ and $a, b \in \mathbb{R}$). Recall that, to encode a position $\vec{r} \in \mathbb{R}^3$ in Euclidean space into c with ERoPE, we calculate

$$\text{ERoPE}_{\vec{u}}(c, \vec{r}) = c \cdot e^{i\omega \vec{u} \cdot \vec{r}}.$$

To avoid computations with complex numbers, we can instead collect the real and imaginary parts of c into a vector $\mathbf{x} = [a \ b]^\top \in \mathbb{R}^2$ and perform an equivalent encoding of \vec{r} into \mathbf{x} as

$$\text{ERoPE}_{\vec{u}}(\mathbf{x}, \vec{r}) = \mathbf{M}\mathbf{x},$$

where \mathbf{M} is the 2×2 rotation matrix

$$\begin{bmatrix} \cos(\omega \vec{u} \cdot \vec{r}) & -\sin(\omega \vec{u} \cdot \vec{r}) \\ \sin(\omega \vec{u} \cdot \vec{r}) & \cos(\omega \vec{u} \cdot \vec{r}) \end{bmatrix}.$$

This formulation can be extended to higher-dimensional vectors $\mathbf{x} \in \mathbb{R}^H$ with $H = 2K$, where \mathbf{M} is now a $(2K) \times (2K)$ block-diagonal matrix consisting of 2×2 rotation matrices and we allow different coefficients ω_k for each of the K rotations matrices (the motivation for this is explained below). For high-dimensional vectors \mathbf{x} , instead of evaluating $\mathbf{M}\mathbf{x}$ via matrix multiplication, it is more efficient (due to the sparsity of \mathbf{M}) to compute

$$\begin{aligned} \text{ERoPE}_{\vec{u}}(\mathbf{x}, \vec{r}) = & \begin{bmatrix} \cos(\omega_1 \vec{u} \cdot \vec{r}) \\ \cos(\omega_1 \vec{u} \cdot \vec{r}) \\ \vdots \\ \cos(\omega_K \vec{u} \cdot \vec{r}) \\ \cos(\omega_K \vec{u} \cdot \vec{r}) \end{bmatrix} \odot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{2K-1} \\ x_{2K} \end{bmatrix} \\ & + \begin{bmatrix} \sin(\omega_1 \vec{u} \cdot \vec{r}) \\ \sin(\omega_1 \vec{u} \cdot \vec{r}) \\ \vdots \\ \sin(\omega_K \vec{u} \cdot \vec{r}) \\ \sin(\omega_K \vec{u} \cdot \vec{r}) \end{bmatrix} \odot \begin{bmatrix} -x_2 \\ x_1 \\ \vdots \\ -x_{2K} \\ x_{2K-1} \end{bmatrix}, \end{aligned} \quad (20)$$

where ‘ \odot ’ denotes element-wise multiplication and x_j is the j -th entry of \mathbf{x} . For vectors with an odd number of dimensions, i.e., $H = 2K - 1$, the missing entry x_{2K} can be simply replaced by zero (\mathbf{x} is zero-padded to an even number of dimensions).

Expressed in words, if we interpret a vector $\mathbf{x} \in \mathbb{R}^{2K}$ as consisting of the concatenated real and imaginary parts of the (complex) entries of a vector $\tilde{\mathbf{x}} \in \mathbb{C}^K$, then applying Eq. 20 to \mathbf{x} (with $\omega_k = \omega$) is equivalent to applying Eq. 9 to $\tilde{\mathbf{x}}$. The motivation for allowing different ω_k for each entry k is that it increases the expressivity of the operation: Recall that after integration over S^2 , encoding positions \vec{r}_m and \vec{r}_n into vectors \mathbf{q}_m and \mathbf{k}_n with ERoPE and computing their scalar product essentially corresponds to scaling $\langle \mathbf{q}_m, \mathbf{k}_n \rangle$ with $\text{sinc}(\omega r_{mn})$ (see Eq. 10 and [Analytic Solution of the Surface Integral](#) in the Supplementary Information for more details). When different ω_k are used for each entry, instead of applying a uniform scaling factor, the contribution of different entries to the scalar product can be scaled independently by factors of $\text{sinc}(\omega_k r_{mn})$, which allows to express a more complicated radial dependence.

The extension of ERoPE to equivariant features $\mathbf{x} \in \mathbb{R}^{P \times (L+1)^2 \times H}$ is straightforward: We simply consider \mathbf{x} as a collection of $P \times (L+1)^2$ vectors of dimension H , to which Eq. 20 is applied separately. In other words, ERoPE is applied along the “feature axis” (of size H), whereas the “parity axis” (of size P) and the “degree axis” (of size $(L+1)^2$) are treated as “batch dimensions”.

Spherical Harmonics The spherical harmonics are functions defined on the surface of the unit sphere $S^2 := \{\vec{u} \in \mathbb{R}^3 : \|\vec{u}\| = 1\}$. We follow the conventions used in Ref. 46 and define them as real and vector-valued, i.e., $Y_\ell : S^2 \mapsto \mathbb{R}^{2\ell+1}$. For example, the first two of these functions (using the Racah normalization convention) are given by $Y_0(\vec{u}) = 1$ and $Y_1(\vec{u}) = \vec{u}$. In general, the individual components of Y_ℓ are polynomials of degree ℓ in the x, y , and z components of the vector \vec{u} (see Ref. 46 for a general definition for arbitrary degrees ℓ).

Lebedev quadrature A Lebedev quadrature⁶⁹ is a numerical approximation to the surface integral of a function f over the unit sphere S^2

$$\int_{S^2} f(\vec{u}) d\vec{u} \approx 4\pi \sum_{j=1}^{N_\Omega} \lambda_j f(\vec{u}_j), \quad (21)$$

where $\vec{u}_j \in S^2$ are grid points and $\lambda_j \in \mathbb{R}$ the corresponding quadrature weights. The grid can be made denser, i.e., the number of grid points N_Ω can be increased, to make the approximation more accurate. Similar to one-dimensional Gaussian quadratures, Lebedev quadratures have the important property that integration of polynomials up to a certain degree is *exact*. For example, a grid with just $N_\Omega = 6$ points integrates polynomials of up to degree 3 (such as the spherical harmonics with $\ell \leq 3$) without error.⁶⁹

Implementation of EFA Given atomic features $\mathbf{x}_m^{(\ell_p)} \in \mathbb{R}^{P \times 2\ell+1 \times H}$, queries, keys and values are calculated as

$$\begin{aligned} \mathbf{q}_m^{(\ell_p)} &= \mathbf{x}_m^{(\ell_p)} W_q^{(\ell_p)}, \\ \mathbf{k}_m^{(\ell_p)} &= \mathbf{x}_m^{(\ell_p)} W_k^{(\ell_p)}, \\ \text{and } \mathbf{v}_m^{(\ell_p)} &= \mathbf{x}_m^{(\ell_p)} W_v^{(\ell_p)}. \end{aligned}$$

Here, $W_q^{(\ell_p)} \in \mathbb{R}^{H \times D_{qk}}$, $W_k^{(\ell_p)} \in \mathbb{R}^{H \times D_{qk}}$ and $W_v^{(\ell_p)} \in \mathbb{R}^{H \times D_v}$ are separate weight matrices for each combination of degree ℓ and parity p and matrix multiplication is performed along the “feature axis”, e.g., $\mathbf{x}_m^{(\ell_p)} W_v^{(\ell_p)} \in \mathbb{R}^{P \times 2\ell+1 \times D_v}$. Atomic positions are encoded into query and key vectors with ERoPE

$$\begin{aligned} \tilde{\mathbf{q}}_{m,\vec{u}}^{(\ell_p)} &:= \text{ERoPE}_{\vec{u}}\left(\vec{r}_m, \boldsymbol{\psi}\left(\mathbf{q}_m^{(\ell_p)}\right)\right), \\ \tilde{\mathbf{k}}_{m,\vec{u}}^{(\ell_p)} &:= \text{ERoPE}_{\vec{u}}\left(\vec{r}_m, \boldsymbol{\psi}\left(\mathbf{k}_m^{(\ell_p)}\right)\right), \end{aligned}$$

where $\boldsymbol{\psi}$ can be any feature map that preserves equivariance (see Eq. 6). Here we choose a simple gated⁷¹ GELU⁷² non-linearity, which is applied element-wise along the “feature axis” of size D_{qk} .

The individual $\tilde{\mathbf{q}}_m^{(\ell_p)}$, $\tilde{\mathbf{k}}_m^{(\ell_p)}$, and $\mathbf{v}_m^{(\ell_p)}$ for all of the N atoms are concatenated row-wise to give

$$\begin{aligned} \tilde{\mathbf{Q}}_{\vec{u}} &\in \mathbb{R}^{N \times P_{qk} \times (L_{qk}+1)^2 \times D_{qk}}, \\ \tilde{\mathbf{K}}_{\vec{u}} &\in \mathbb{R}^{N \times P_{qk} \times (L_{qk}+1)^2 \times D_{qk}}, \\ \mathbf{V} &\in \mathbb{R}^{N \times P_v \times (L_v+1)^2 \times D_v}. \end{aligned}$$

Note that queries $\tilde{\mathbf{Q}}_{\vec{u}}$ and keys $\tilde{\mathbf{K}}_{\vec{u}}$ must have the same maximal degree L_{qk} , parity size P_{qk} , and feature dimension D_{qk} , whereas the corresponding choices for the values \mathbf{V} may be different. Additionally, we include \vec{u} as a subscript for $\tilde{\mathbf{Q}}_{\vec{u}}$ and $\tilde{\mathbf{K}}_{\vec{u}}$ to make explicit that their entries depend on the choice of \vec{u} used for the ERoPE encoding of the atomic positions.

Using slight abuse of notation (which is motivated below), we can now “translate” the linear-scaling attention-like mechanism introduced in Eq. 12 to

$$\mathbf{B}_{\vec{u}} := \hat{\mathbf{Q}}_{\vec{u}} \left(\hat{\mathbf{K}}_{\vec{u}}^\top \mathbf{V} \right), \quad (22)$$

where $\mathbf{B}_{\vec{u}} \in \mathbb{R}^{N \times P_v \times (L_v+1)^2 \times D_v}$. Here, the “outer product” $\hat{\mathbf{K}}_{\vec{u}}^\top \mathbf{V}$ of keys and values can be computed with time complexity $\mathcal{O}(N)$ and should be thought of as evaluating to an array of shape $[P_{qk} \times (L_{qk}+1)^2 \times D_{qk} \times P_v \times (L_v+1)^2 \times D_v]$. The subsequent multiplication with $\hat{\mathbf{Q}}_{\vec{u}}$ from the left side should be understood as reducing over (element-wise multiplication followed by summation) the last three axes of $\hat{\mathbf{Q}}_{\vec{u}}$ and the first three axes of $\hat{\mathbf{K}}_{\vec{u}}^\top \mathbf{V}$. This operation also has a time complexity

of $O(N)$, which makes the full operation (Eq. 22) scale linearly in the number of atoms. We use the suggestive “matrix-like” notation in Eq. 22, because when $P_{qk} = P_v = 1$ and $L_{qk} = L_v = 0$ (special case for “ordinary” invariant features), many axes of the involved arrays have size 1 and can be “squeezed away” for simplicity. In this special case, $\tilde{\mathbf{Q}}_{\vec{u}}, \tilde{\mathbf{K}}_{\vec{u}} \in \mathbb{R}^{N \times D_{qk}}$ and $V \in \mathbb{R}^{N \times D_v}$, so Eq. 22 is correctly described with conventional notation for matrix operations.

The full equivariant Euclidean fast attention (EFA) update (Eq. 13) is implemented as

$$\text{EFA}(\mathcal{X}, \mathcal{R}) = \sum_{j=1}^{N_\Omega} \lambda_j \mathbf{B}_{\vec{u}_j} \bigotimes_{L_{\text{out}}}^{L_v L_Y} \mathbf{Y}_{\vec{u}_j}, \quad (23)$$

where we replace the integral over S^2 with a Lebedev quadrature and $\mathbf{Y}_{\vec{u}} \in \mathbb{R}^{1 \times 1 \times (L_Y+1)^2 \times 1}$ is the concatenation of spherical harmonics vectors $Y_\ell(\vec{u})$ up to maximal degree L_Y with additional “dummy axes” of size 1 that are broadcasted over the corresponding axes of $\mathbf{B}_{\vec{u}_j}$, so that the tensor product contractions (see Eq. 19) are well-defined. When $L_Y = 0$ (and $\mathbf{B}_{\vec{u}_j}$ contains only invariant information, i.e., $P_v = 1$ and $L_v = 0$), the invariant form of EFA (Eq. 11) is recovered.

Numerical accuracy Since the integration over S^2 is performed numerically in Eq. 23, the accuracy of the rotational invariance/equivariance of the EFA update is directly related to the precision of the Lebedev quadrature. To better understand the limitations of this approach, we point out that the function that is integrated in Eq. 23 basically corresponds to a linear combination of terms of the form $\sin(\omega_k \vec{u} \cdot \vec{r}_{mn}) \cdot \text{poly}_\ell(\vec{u})$ and $\cos(\omega_k \vec{u} \cdot \vec{r}_{mn}) \cdot \text{poly}_\ell(\vec{u})$, where $\text{poly}_\ell(\vec{u})$ is some polynomial of degree $\ell \leq L_Y$ in the x , y , and z components of the vector \vec{u} (the sine/cosine components stem from ERoPE, whereas the polynomials come from the spherical harmonics). As mentioned above, polynomials up to a certain degree are integrated exactly when using Lebedev quadratures. While, the sine/cosine components are not polynomials, it is still possible to estimate a sort of “pseudo-degree” for them: We can approximate these terms by a Taylor series in $b := \omega_k \vec{u} \cdot \vec{r}_{mn}$ around zero. Clearly, the larger b is,

the more terms are required in the Taylor series for a good approximation, which corresponds to a larger pseudo-degree. The largest “total degree” of the terms in the integral is then given by the sum of the highest degree L_Y of the spherical harmonics and the estimated pseudo-degree of the sine/cosine terms. Assuming a fixed budget of grid points N_Ω , the integration will be sufficiently accurate up to some maximal input value b_{max} which yields the inequality $\omega_k r_{nm} \leq b_{\text{max}}$. This can be used to define a maximal frequency

$$\omega_{\text{max}} = \frac{b_{\text{max}}}{r_{\text{max}}} \quad (24)$$

up to which (almost) exact rotational invariance/equivariance is ensured. Here, r_{max} is the maximal expected distance between atoms (which can be estimated from the training data), and the value of ω_{max} can be chosen accordingly. Alternatively, the number of grid points N_Ω can be increased to make b_{max} larger, which in turn allows choosing a larger ω_{max} . Because the analytic solution of the integral is known (see [Analytic Solution of the Surface Integral](#) in the Supplementary Information), the value of b_{max} (for a given N_Ω) can be easily determined numerically. In this work, we define b_{max} as the largest value for which the absolute deviation between numerical and analytic solution stays below a value of 10^{-5} (which roughly corresponds to the numerical precision of single precision floating point arithmetic). We confirm empirically that our implementation of EFA (Eq. 23) is rotationally invariant/equivariant up to the desired numerical precision in Fig. S2 and provide precomputed values for b_{max} as function of N_Ω in Tab. S3.

Neural Network Implementation

For the experiments performed in this work, we choose an equivariant MPNN in the spirit of tensor field networks⁷³ or NequIP.³¹ This base model is then optionally augmented with EFA blocks, and we compare the performance of the different architectural variants. Because our analyses do not depend on particular model details, but rather investigate general limitations inherent to all strictly local MLFFs, we expect our results to also transfer to other choices of local models

(which can be augmented with EFA blocks in the same manner). MPNNs are merely chosen as a representative local model due to their popularity, and because an effective cutoff (extending beyond the local cutoff) only exists for MPNN-like models, and needs to be carefully considered during the analysis of results (other local models, such as Behler-Parrinello neural networks,¹⁰ can often be regarded as an MPNN-like model with $T = 1$ in our considerations).

For the implementations of equivariant operations we use the E3x library⁴⁶ which is build on top of JAX⁷⁴ and FLAX.⁷⁵ An implementation of the Euclidean fast attention block is publicly available at https://github.com/thorben-frank/euclidean_fast_attention.

Message Passing Neural Networks We follow the typical design of MPNNs discussed in [Context and Challenges](#) and iteratively update initial embeddings using repeated application of [Eq. 1](#) (the MessagePass block is described below). The atomic representations with $\ell = 0$ and even parity are initialized via learned atom type embeddings dependent on the atomic numbers $\mathcal{Z} := \{z_1, \dots, z_N \mid z_m \in \mathbb{N}_+\}$. In case of the charged cluster experiments (see [Idealized systems](#)), we instead use the atomic charges $\mathcal{Q} := \{q_1, \dots, q_N \mid q_m \in \mathbb{Z}\}$ to assign the embeddings. All other degree and parity channels are always initialized to zero with the exception of the charge-dipole experiment, where we use the atomic dipole vector with an additional parity and feature axes of size 1, i.e., $\vec{\mu} \in \mathbb{R}^{1 \times 3 \times 1}$ to initialize the features with $\ell = 1$ and odd parity as $\mathbf{x}^{(\ell-)} = \vec{\mu}W$ ($W \in \mathbb{R}^{1 \times H}$ is a learnable weight matrix that maps the dipole vector to the feature space).

The $\ell = 0$ components with even parity (invariant parts) of the final atomic features $\mathbf{X}^{[T]}$ are used to predict per atom energies which are summed to give the total energy

$$E_{\text{ML}} = \sum_{m=1}^N w^\top \mathbf{x}_m^{[T](0_+)} + E_{m,\text{shift}}, \quad (25)$$

where $w \in \mathbb{R}^H$ is a trainable vector, $E_{m,\text{shift}} \in \mathbb{R}$ is a trainable atom type dependent shift and

$\mathbf{x}^{[T](0_+)} \in \mathbb{R}^H$ is the invariant part of the final features at layer T (parity and degree axes of size 1 are removed here for clarity). In the presence of atomic dipoles, the invariant part is obtained via a tensor product which maps the atomic features of all degrees to invariant features. Forces are calculated as the negative gradient w. r. t. the atomic positions as $\vec{F}_m = -\nabla_{\vec{r}_m} E \in \mathbb{R}^3$, which can be done efficiently using automatic differentiation.

Equivariant Message Passing Block Throughout this work we use a prototypical equivariant message passing neural network (MPNN). It is build around equivariant continuous convolutions which are employed in many state-of-the-art MPNNs like NequIP,³¹ TFN,⁷³ and Equiformer.⁴⁷ The message sent from atom n to atom m is given by

$$\mathbf{m}_{mn}^{[t]} = (\mathbf{W}(r_{mn}) \circ \mathbf{Y}(\hat{r}_{mn})) \bigotimes_{L_{\text{MP}}}^{L_Y L_X} \mathbf{x}_n, \quad (26)$$

where $\mathbf{W}(r_{mn}) \in \mathbb{R}^{1 \times (L_Y+1)^2 \times H}$ is a learned radial filter function containing H features for each degree ℓ up to L_Y . Importantly, the ‘‘subslices’’ of size $2\ell+1$ for each degree ℓ are constrained to contain identical entries (otherwise the operation would not preserve equivariance). The spherical harmonics $\vec{Y}(\hat{r}_{mn}) \in \mathbb{R}^{1 \times (L_Y+1)^2 \times 1}$ use ‘‘dummy axes’’ of size 1 (similar to the use of spherical harmonics in [Eq. 23](#)) and the element-wise multiplication ‘ \circ ’ is broadcasted along these axes. The messages from all neighbors n in the neighborhood \mathcal{N}_m of atom m are aggregated as $\mathbf{m}_m^{[t]} = \sum_{n \in \mathcal{N}_m} \mathbf{m}_{mn}^{[t]}$. The neighborhood \mathcal{N}_m contains all atoms within the chosen cutoff distance, i.e., those that satisfy $r_{mn} < r_{\text{cut}}$. Note that the radial filter $\mathbf{W}(r_{mn})$ in [Eq. 26](#) is enforced to smoothly decay to zero at r_{cut} , so that the learned representations vary continuously when atoms enter or leave \mathcal{N}_m . Finally, the updated atom representations are obtained as

$$\mathbf{x}_m^{[t+1]} = \text{MLP}[\mathbf{x}_m^{[t]} + \mathbf{m}_m^{[t]}], \quad (27)$$

where MLP is a multi-layer perceptron network with two equivariant dense layers and gated SiLU⁷⁶ non-linearity.

For standard MPNNs, each MP update consists of an equivariant MP Block. For an MPNN aug-

mented with EFA, the atomic representations are passed to the equivariant MP block *and* to the EFA block (see below) and their outputs are summed up afterwards (Fig. 2D).

Euclidean Fast Attention Block Non-local atomic representations are calculated using the EFA update (Eq. 23), such that

$$\mathbf{m}_{m,\text{nl}}^{[t]} = \text{EFA}(\mathcal{X}^{[t-1]}, \mathcal{R})_m. \quad (28)$$

As for the MP block, per-atomic embeddings are refined via an equivariant MLP, such that

$$\mathbf{x}_{m,\text{nl}}^{[t+1]} = \text{MLP}[\mathbf{x}_m^{[t]} + \mathbf{m}_{m,\text{nl}}^{[t]}]. \quad (29)$$

Training

Models are trained by minimizing a combined loss of energy and forces

$$\begin{aligned} \mathcal{L} = & \frac{\lambda_E}{B} \sum_{b=1}^B (E_{b,\text{true}} - E_{b,\text{ML}})^2 \\ & + \frac{\lambda_F}{B} \sum_{b=1}^B \frac{1}{N_b} \sum_{m=1}^{N_b} \|\vec{F}_{m,\text{true}} - \vec{F}_{m,\text{ML}}\|_2^2, \end{aligned} \quad (30)$$

where B is the number of molecules per batch, N_b is the number of atoms for molecule in batch b and λ_E and λ_F are scaling parameters for the energy and force components of the loss.

We use Adam⁷⁷ for parameter optimization and an initial learning rate of $\mu = 10^{-3}$. The learning rate is decayed to 10^{-5} at the end of training via an exponential decay schedule. For optimization the optax⁷⁸ library is used. In Tab. S2 we report full information on optimization settings for each experiment.

Although we did not apply this in the reported experiments, we find, that re-scaling individual components of the EFA update can help to improve the training dynamics (see Fig. S9 and section Degree Re-Scaling in the SI).

Model Hyperparameters

Default EFA blocks employ a query and key dimension of $D_{qk} = 16$ and a value dimension of $D_v = 32$. The number of Lebedev grid points

is set to $N_\Omega = 50$. The maximal atomic separation varies between data sets and is reported in Tab. S2. For simplicity, we assume the degrees for queries, keys and values to always be equal and define the shorthand $L_{\text{EFA}} = L_{qk} = L_v$. Other model hyperparameters are described below.

Geometric Expressiveness For the k -chains experiment, we use the equivariant MP block and the EFA block as standalone component. For the MP block, we follow the original publication⁴⁹ and use a cutoff of $r_{\text{cut}} = 10 \text{ \AA}$ and no radial cutoff function. For the EFA block we use the hyperparameters outlined above. The final output has two neurons whose values are passed through softmax to predict a probability for each class. The models are trained by minimizing the cross-entropy to the true class labels. The number of features for both blocks is $H = 32$ and the maximal degree in the MP block is $L_{\text{MP}} = 2$. The number of MP updates is increased until it is able to correctly classify the graphs. Because the maximal separation r_{max} between atoms increases with k , the maximal frequency for ERoPE (see Eq. 24) must be decreased accordingly, leading to slower training convergence for the longest ($k = 8$) chain. Based on theoretical analysis of positional encoding techniques⁷⁹ we can speed up convergence by increasing $b_{\text{max}} = 2\pi$ and adjust the number of Lebedev grid points to $N_\Omega = 86$ accordingly. Tab. S3 provides a mapping between b_{max} and N_Ω . For the neighborhood distinguish-ability experiments we use an EFA update as described in Eq. 13 and produce per degree invariants, by taking the L2-norm per degree in the equivariant representation \mathbf{x} of the central node. We compare the values for each degree ℓ for each pair of distinct neighborhoods and label the smallest degree for which the invariants are different.

Pairwise Potentials For the toy systems with a pairwise potential we use a feature dimension of $H = 128$. A maximal degree of $L_{\text{MP}} = 1$ is used in the MP block and $L_{\text{EFA}} = 0$ in the EFA block. The spherical harmonics vector has maximal degree $L_Y = 0$ for all models except for the one applied to the charge-dipole system, for which we additionally train a model with $L_Y = 1$ to show that

the task can only be solved by including directional information (see main text). The number of MP layers is $T = 2$ and is increased up to $T = 4$ for the local MP models in the 3D cluster data experiment. All models employ a local MP cutoff of $r_{\text{cut}} = 5 \text{ \AA}$.

Molecular Systems For the molecular systems we use $T = 2$ for the $S_{\text{N}2}$ experiments and $T = 3$ for the dimer and cumulene results. The local cutoff is varied between $r_{\text{cut}} = 5 \text{ \AA}$ and $r_{\text{cut}} = 10 \text{ \AA}$ for the $S_{\text{N}2}$ experiments as described in the main text and is $r_{\text{cut}} = 4 \text{ \AA}$ for the dimers and $r_{\text{cut}} = 3 \text{ \AA}$ for cumulene. The feature dimension for the models with EFA block is $H = 128$ for $T = 2$ and $H = 64$ for $T = 3$. For models without EFA block, the feature dimension is increased to $H = 162$ for $T = 2$ and $H = 84$ for $T = 3$, such that models with and without EFA block have approximately ($\leq 1\%$ deviation) the same number of parameters. The maximal degree in the MP block is $L_{\text{MP}} = 2$ and $L_{\text{EFA}} = 0$ in the EFA block, except for the cumulene structures where we additionally train models with $L_{\text{EFA}} = 1$ and $L_{\text{EFA}} = 2$ (see main text). The maximal degree of the spherical harmonics vector in the EFA update is always $L_Y = 0$.

Dimer Data Subset Creation

The SPICE data set⁸⁰ serves as starting point for the dimer data. SPICE contains energies and forces for the dimer geometries in the 370K data set⁵⁵ re-calculated using density functional theory (DFT). The resulting data set has 4612 entries and covers a variety of biologically relevant dimers. Each dimer consists of two molecules (also called monomers) which interact with each other via non-covalent (long-range) interactions. However, the generated data set is unbalanced in the sense that certain monomers appear much more often within dimers than others, biasing the model towards a correct description of dimers with these monomers. To alleviate this problem, we create a curated subset of the dimers data which only uses the most frequent 9 monomers. This results in a new data set of 76 dimers. Almost all combinations of monomers present in the original data set are also present in the curated version (76 vs. 81 possible combinations).

Acknowledgements

This work was in part supported by the German Ministry for Education and Research (BMBF) under Grants 01IS14013A-E, 01GQ1115, 01GQ0850, 01IS18025A, 031L0207D, and 01IS18037A. K.R.M. was partly supported by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grants funded by the Korea government (MSIT) (No. 2019-0-00079, Artificial Intelligence Graduate School Program, Korea University and No. 2022-0-00984, Development of Artificial Intelligence Technology for Personalized Plug-and-Play Explanation and Verification of Explanation). We thank Stefan Blücher and Hartmut Maennel for helpful comments on the manuscript.

References

- [1] Daniel Jurafsky. Speech and language processing, 2000.
- [2] Robert Geirhos, David HJ Janssen, Heiko H Schütt, Jonas Rauber, Matthias Bethge, and Felix A Wichmann. Comparing deep neural networks against humans: object recognition when the signal gets weaker. *arXiv preprint arXiv:1706.06969*, 2017.
- [3] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.
- [4] Matthias Rupp, Alexandre Tkatchenko, Klaus-Robert Müller, and O Anatole Von Lilienfeld. Fast and accurate modeling of molecular atomization energies with machine learning. *Phys. Rev. Lett.*, 108(5): 058301, 2012.
- [5] O Anatole von Lilienfeld, Klaus-Robert Müller, and Alexandre Tkatchenko. Exploring chemical compound space with quantum-based machine learning. *Nat. Rev. Chem.*, 4(7):347–358, 2020.
- [6] Oliver T Unke, Stefan Chmiela, Huziel E Sauceda, Michael Gastegger, Igor Poltavsky,

- Kristof T Schütt, Alexandre Tkatchenko, and Klaus-Robert Müller. Machine learning force fields. *Chem. Rev.*, 121(16):10142–10186, 2021.
- [7] John A Keith, Valentin Vassilev-Galindo, Bingqing Cheng, Stefan Chmiela, Michael Gastegger, Klaus-Robert Müller, and Alexandre Tkatchenko. Combining machine learning and computational chemistry for predictive insights into chemical systems. *Chem. Rev.*, 121(16):9816–9872, 2021.
- [8] Xiang Fu, Zhenghao Wu, Wujie Wang, Tian Xie, Sinan Keten, Rafael Gomez-Bombarelli, and Tommi Jaakkola. Forces are not enough: Benchmark and critical evaluation for machine learning force fields with molecular simulations. *arXiv preprint arXiv:2210.07237*, 2022.
- [9] Martin Karplus and J Andrew McCammon. Molecular dynamics simulations of biomolecules. *Nature structural biology*, 9(9):646–652, 2002.
- [10] Jörg Behler and Michele Parrinello. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett.*, 98(14):146401, 2007.
- [11] Albert P Bartók, Mike C Payne, Risi Kondor, and Gábor Csányi. Gaussian Approximation Potentials: the accuracy of quantum mechanics, without the electrons. *Phys. Rev. Lett.*, 104(13):136403, 2010.
- [12] Kristof T Schütt, Huziel E Saucedo, Pieter-Jan Kindermans, Alexandre Tkatchenko, and Klaus-Robert Müller. SchNet – a deep learning architecture for molecules and materials. *J. Chem. Phys.*, 148(24):241722, 2018.
- [13] Kristof T Schütt, Stefan Chmiela, O Anatole Von Lilienfeld, Alexandre Tkatchenko, Koji Tsuda, and Klaus-Robert Müller. Machine learning meets quantum physics. *Lecture Notes in Physics* 968, 2020.
- [14] Oliver T Unke, Martin Stöhr, Stefan Ganschä, Thomas Unterthiner, Hartmut Maennel, Sergii Kashubin, Daniel Ahlin, Michael Gastegger, Leonardo Medrano Sandonas, Joshua T Berryman, et al. Biomolecular dynamics with machine-learned quantum-mechanical force fields trained on diverse chemical fragments. *Science Advances*, 10(14):eadn4397, 2024.
- [15] A. Kabylda, J. T. Frank, S. S. Dou, A. Khabibrakhmanov, L. M. Sandonas, O. T. Unke, S. Chmiela, K.R. Müller, and A. Tkatchenko. Molecular simulations with a pretrained neural network and universal pairwise force fields. *ChemRxiv*, 2024. doi: 10.26434/chemrxiv-2024-bdfr0.
- [16] Amil Merchant, Simon Batzner, Samuel S Schoenholz, Muratahan Aykol, Gwooon Cheon, and Ekin Dogus Cubuk. Scaling deep learning for materials discovery. *Nature*, 624(7990):80–85, 2023.
- [17] LM Woods, Diego Alejandro Roberto Dalvit, Alexandre Tkatchenko, P Rodriguez-Lopez, Alejandro W Rodriguez, and R Podgornik. Materials perspective on casimir and van der waals interactions. *Reviews of Modern Physics*, 88(4):045003, 2016.
- [18] Jan Hermann, Robert A DiStasio Jr, and Alexandre Tkatchenko. First-principles models for van der waals interactions in molecules and materials: Concepts, theory, and applications. *Chemical Reviews*, 117(6):4714–4758, 2017.
- [19] Martin Stöhr, Troy Van Voorhis, and Alexandre Tkatchenko. Theory and practice of modeling van der waals interactions in electronic-structure calculations. *Chemical Society Reviews*, 48(15):4118–4154, 2019.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- [21] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16

- words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [22] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- [23] Ingrid von Glehn, James S Spencer, and David Pfau. A self-attention ansatz for ab-initio quantum chemistry. *arXiv preprint arXiv:2211.13672*, 2022.
- [24] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- [25] Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023.
- [26] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are RNNs: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pages 5156–5165. PMLR, 2020.
- [27] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- [28] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272. Pmlr, 2017.
- [29] Linfeng Zhang, Jiequn Han, Han Wang, Wisam Saidi, Roberto Car, et al. End-to-end symmetry preserving inter-atomic potential energy model for finite and extended systems. *Advances in neural information processing systems*, 31, 2018.
- [30] Oliver T Unke, Stefan Chmiela, Michael Gastegger, Kristof T Schütt, Huziel E Sauceda, and Klaus-Robert Müller. SpookyNet: Learning force fields with electronic degrees of freedom and nonlocal effects. *Nat. Commun.*, 12:7273, 2021.
- [31] Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P Mailoa, Mordechai Kornbluth, Nicola Molinari, Tess E Smidt, and Boris Kozinsky. E (3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nat. Commun.*, 13(1):2453, 2022.
- [32] J Thorben Frank, Oliver T Unke, Klaus-Robert Müller, and Stefan Chmiela. A euclidean transformer for fast and stable machine learned force fields. *Nature Communications*, 15(1):6539, 2024.
- [33] Igor Poltavsky, Anton Charkin-Gorbunin, Mirela Puleva, Gregory Cordeiro Fonseca, Ilyes Batatia, Nicholas J Browning, Stefan Chmiela, Mengnan Cui, J Thorben Frank, Stefan Heinen, et al. Crash testing machine learning force fields for molecules, materials, and interfaces: Model analysis in the tea challenge 2023. *ChemRxiv*, 2024.
- [34] Igor Poltavsky, Mirela Puleva, Anton Charkin-Gorbunin, Gregory Cordeiro Fonseca, Ilyes Batatia, Nicholas J Browning, Stefan Chmiela, Mengnan Cui, J Thorben Frank, Stefan Heinen, et al. Crash testing machine learning force fields for molecules, materials, and interfaces: Molecular dynamics in the tea challenge 2023. *ChemRxiv*, 2024.
- [35] Justin S Smith, Olexandr Isayev, and Adrian E Roitberg. ANI-1: an extensible neural network potential with DFT accuracy at force field computational cost. *Chem. Sci.*, 8(4):3192–3203, 2017.
- [36] Anders S Christensen, Lars A Bratholm, Felix A Faber, and O Anatole von Lilienfeld. FCHL revisited: faster and more accurate quantum machine learning. *J. Chem. Phys.*, 152(4):044107, 2020.

- [37] Albert Musaelian, Simon Batzner, Anders Johansson, Lixin Sun, Cameron J Owen, Mordechai Kornbluth, and Boris Kozinsky. Learning local equivariant representations for large-scale atomistic dynamics. *Nat. Commun.*, 14(1):579, 2023.
- [38] J Mercer. Functions of positive and negative type and their connection with theory of integral equations. *Philosophical Transactions of the Royal Society A*, 209(441–458): 415–446, 1909.
- [39] Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950.
- [40] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.
- [41] Bernhard Scholkopf, Sebastian Mika, Chris JC Burges, Philipp Knirsch, K-R Muller, Gunnar Ratsch, and Alexander J Smola. Input space versus feature space in kernel-based methods. *IEEE transactions on neural networks*, 10(5):1000–1017, 1999.
- [42] K.-R. Müller, S. Mika, G. Ratsch, K. Tsuda, and B. Schölkopf. An introduction to kernel-based learning algorithms. *IEEE Transactions on Neural Networks*, 12(2):181–201, 2001.
- [43] Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- [44] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [45] Thorben Frank, Oliver Unke, and Klaus-Robert Müller. So3krates: Equivariant attention for interactions on arbitrary length-scales in molecular systems. *Advances in Neural Information Processing Systems*, 35: 29400–29413, 2022.
- [46] Oliver T Unke and Hartmut Maennel. E3x: E(3)-equivariant deep learning made easy. *arXiv preprint arXiv:2401.07595*, 2024.
- [47] Yi-Lun Liao and Tess Smidt. Equiformer: Equivariant graph attention transformer for 3d atomistic graphs. *arXiv preprint arXiv:2206.11990*, 2022.
- [48] Sergey N Pozdnyakov, Michael J Willatt, Albert P Bartók, Christoph Ortner, Gábor Csányi, and Michele Ceriotti. Incompleteness of atomic structure representations. *Phys. Rev. Lett.*, 125(16):166001, 2020.
- [49] Chaitanya K Joshi, Cristian Bodnar, Simon V Mathis, Taco Cohen, and Pietro Lio. On the expressive power of geometric graph neural networks. In *International Conference on Machine Learning*, pages 15330–15355. PMLR, 2023.
- [50] Andrei Leman and Boris Weisfeiler. A reduction of a graph to a canonical form and an algebra arising during this reduction. *Nauchno-Technicheskaya Informatsiya*, 2(9): 12–16, 1968.
- [51] Oliver T Unke and Markus Meuwly. PhysNet: A neural network for predicting energies, forces, dipole moments, and partial charges. *J. Chem. Theory Comput.*, 15(6): 3678–3693, 2019.
- [52] Kristof Schütt, Oliver Unke, and Michael Gastegger. Equivariant message passing for the prediction of tensorial properties and molecular spectra. In *International Conference on Machine Learning*, pages 9377–9388. PMLR, 2021.
- [53] Uri Alon and Eran Yahav. On the bottleneck of graph neural networks and its practical implications. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=i800PhOCVH2>.
- [54] Stefan Chmiela, Huziel E Sauceda, Igor Poltavsky, Klaus-Robert Müller, and Alexandre Tkatchenko. sgdml: Constructing accurate and data efficient molecular force fields

- using machine learning. *Comput. Phys. Commun.*, 240:38–45, 2019.
- [55] Alexander G Donchev, Andrew G Taube, Elizabeth Decolvenaere, Cory Hargus, Robert T McGibbon, Ka-Hei Law, Brent A Gregersen, Je-Luen Li, Kim Palmo, Karthik Siva, et al. Quantum chemical benchmark databases of gold-standard dimer interaction energies. *Scientific data*, 8(1):55, 2021.
- [56] Heikki Muhli, Xi Chen, Albert P Bartók, Patricia Hernández-León, Gábor Csányi, Tapio Ala-Nissila, and Miguel A Caro. Machine learning force fields based on local parametrization of dispersion interactions: Application to the phase diagram of c 60. *Physical Review B*, 104(5):054106, 2021.
- [57] Julia Westermayr, Shayantan Chaudhuri, Andreas Jeindl, Oliver T Hofmann, and Reinhard J Maurer. Long-range dispersion-inclusive machine learning potentials for structure search and optimization of hybrid organic–inorganic interfaces. *Digital Discovery*, 1(4):463–475, 2022.
- [58] Nongnuch Artrith, Tobias Morawietz, and Jörg Behler. High-dimensional neural-network potentials for multicomponent systems: Applications to zinc oxide. *Physical Review B*, 83(15):153101, 2011.
- [59] Tobias Morawietz, Vikas Sharma, and Jörg Behler. A neural network potential-energy surface for the water dimer based on environment-dependent atomic energies and charges. *The Journal of chemical physics*, 136(6), 2012.
- [60] Tsz Wai Ko, Jonas A Finkler, Stefan Goedecker, and Jörg Behler. A fourth-generation high-dimensional neural network potential with accurate electrostatics including non-local charge transfer. *Nat. Commun.*, 12(1):1–11, 2021.
- [61] Andrea Grisafi and Michele Ceriotti. Incorporating long-range physics in atomic-scale machine learning. *The Journal of chemical physics*, 151(20), 2019.
- [62] Joshua Pagotto, Junji Zhang, and Timothy Duignan. Predicting the properties of salt water using neural network potentials and continuum solvent theory. *ChemRxiv*, 2022. doi: 10.26434/chemrxiv-2022-jndlx.
- [63] Yunyang Li, Yusong Wang, Lin Huang, Han Yang, Xinran Wei, Jia Zhang, Tong Wang, Zun Wang, Bin Shao, and Tie-Yan Liu. Long-short-range message-passing: A physics-informed framework to capture non-local interaction for scalable molecular dynamics simulation. *arXiv preprint arXiv:2304.13542*, 2023.
- [64] Philip Loche, Kevin K Huguenin-Dumittan, Melika Honarmand, Qianjun Xu, Egor Rumiantssev, Wei Bin How, Marcel F Langer, and Michele Ceriotti. Fast and flexible range-separated models for atomistic machine learning. *arXiv preprint arXiv:2412.03281*, 2024.
- [65] Matteo Gori, Philip Kurian, and Alexandre Tkatchenko. Second quantization of many-body dispersion interactions for chemical and biological systems. *Nature Communications*, 14(1):8218, 2023.
- [66] Ilyes Batatia, Lars L Schaaf, Huajie Chen, Gábor Csányi, Christoph Ortner, and Felix A Faber. Equivariant matrix function neural networks. *arXiv preprint arXiv:2310.10434*, 2023.
- [67] Arthur Kosmala, Johannes Gasteiger, Nicholas Gao, and Stephan Günnemann. Ewald-based long-range message passing for molecular graphs. In *International Conference on Machine Learning*, pages 17544–17563. PMLR, 2023.
- [68] Yusong Wang, Tong Wang, Shaoning Li, Xinheng He, Mingyu Li, Zun Wang, Nanning Zheng, Bin Shao, and Tie-Yan Liu. Enhancing geometric representations for molecules with equivariant vector-scalar interactive message passing. *Nature Communications*, 15(1):313, 2024.
- [69] Vyacheslav Ivanovich Lebedev. Quadratures on a sphere. *USSR Computational Mathe-*

- matics and Mathematical Physics*, 16(2):10–24, 1976.
- [70] Oliver Unke, Mihail Bogojeski, Michael Gastegger, Mario Geiger, Tess Smidt, and Klaus-Robert Müller. Se(3)-equivariant prediction of molecular wavefunctions and electronic densities. *Advances in Neural Information Processing Systems*, 34, 2021.
- [71] Maurice Weiler, Mario Geiger, Max Welling, Wouter Boomsma, and Taco S Cohen. 3d steerable CNNs: Learning rotationally equivariant features in volumetric data. *Advances in Neural Information Processing Systems*, 31, 2018.
- [72] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (GELUs). *arXiv preprint arXiv:1606.08415*, 2016.
- [73] Nathaniel Thomas, Tess Smidt, Steven Kearnes, Lusann Yang, Li Li, Kai Kohlhoff, and Patrick Riley. Tensor field networks: rotation-and translation-equivariant neural networks for 3D point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- [74] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- [75] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. URL <http://github.com/google/flax>.
- [76] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107: 3–11, 2018.
- [77] Diederik P Kingma and Jimmy Ba. Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [78] Matteo Hessel, David Budden, Fabio Viola, Mihaela Rosca, Eren Sezener, and Tom Hennigan. Optax: composable gradient transformation and optimisation, in jax!, 2020. URL <http://github.com/deepmind/optax>.
- [79] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in neural information processing systems*, 33:7537–7547, 2020.
- [80] Peter Eastman, Pavan Kumar Behara, David L Dotson, Raimondas Galvelis, John E Herr, Josh T Horton, Yuezhi Mao, John D Chodera, Benjamin P Pritchard, Yuanqing Wang, et al. Spice, a dataset of drug-like molecules and peptides for training machine learning potentials. *Scientific Data*, 10(1):11, 2023.

Supplementary Information

Euclidean Fast Attention: Machine Learning Global Atomic Representations at Linear Cost

Numerical Verification of Rotational Symmetry

Rotational invariance of the predicted energy is verified numerically for 400 example structures of S_N2 with varying spatial extent (increasing separation between the iodine atom (purple) and the central carbon hydrogen complex). As expected from our analysis, the predictions are invariant (up to numerical noise) when the Lebedev grid is sufficiently large, see [Fig. S1](#).

Comparison to Other Models

Pairwise Potential For the pairwise potential we train on a two atom system as the one investigated in the main text with a maximal separation of 10\AA . We train an MP+EFA (with the same settings as in the main text for the two atom system) and a SpookyNet model with the default settings from the original paper. The SpookyNet model includes non-local corrections via linear scaling self-attention on the atomic representations. However, this neglects the relative positioning $\vec{r}_m - \vec{r}_n$ of the atoms w. r. t. each other such that it can not describe the potential beyond the local message passing cutoff ([Fig. S5A](#)).

Cumulene For comparison we choose a global ($r_{\text{cut}} = 12\text{\AA}$) SchNet model, a global (per construction) sGDML⁵⁴ kernel, a SpookyNet model and a NequIP³¹ model with an effective cutoff of 9\AA ([Fig. S5B](#)). Both, the SchNet and sGDML model rely on invariant descriptors that are based on pairwise distances only. Due to the large separation of the hydrogen rotors, changes in the dihedral angle correspond to variations in the pairwise distances which are too small to be resolved. The linear self-attention mechanism in SpookyNet fails since it is not able to discriminate interaction patterns which depend on Euclidean information. The equivariant NequIP model is capable of solving the problem by using a sufficient number of MP layers⁴⁵ but as soon as the effective cutoff becomes too small it can not solve the regression task. Trivially one can always solve such tasks by using a sufficient number of MP layers but this solution is practically infeasible due to the increasing computational cost and information blur (over-squashing).

Analytic Solution of the Surface Integral

In the following we give a step-by-step solution of the surface integral (see [Eq. 10](#))

$$I = \frac{1}{4\pi} \int_{S^2} e^{i\omega\vec{u}\cdot\vec{r}_{mn}} d\vec{u}.$$

We start by re-writing the integral in spherical coordinates

$$I = \frac{1}{4\pi} \int_0^{2\pi} \int_0^\pi e^{i\omega r_{mn} \cos(\theta)} \sin(\theta) d\theta d\phi,$$

where we have substituted $\vec{u} = [\sin(\theta) \cos(\phi) \sin(\theta) \sin(\phi) \cos(\theta)]^\top$ and assumed (without loss of generality) that \vec{r}_{mn} is aligned with the z -axis of our coordinate system, such that $\vec{u} \cdot \vec{r}_{mn} = r_{mn} \cos(\theta)$

(with $r_{mn} = \|\vec{r}_{mn}\|$). Next, we substitute $x := -\cos(\theta)$ such that $\sin(\theta) d\theta = dx$ and $\int_0^\pi \rightarrow \int_{-1}^1$, giving

$$I = \frac{1}{4\pi} \int_0^{2\pi} \int_{-1}^1 e^{-i\omega r_{mn} x} dx d\phi.$$

Using the analytic solution for

$$\int_{-1}^1 e^{-iax} dx = 2 \frac{\sin(a)}{a} = 2 \operatorname{sinc}(a),$$

we arrive at the solution

$$I = \frac{1}{2\pi} \int_0^{2\pi} \operatorname{sinc}(\omega r_{mn}) d\phi = \operatorname{sinc}(\omega r_{mn}).$$

For the general case with $\ell \geq 0$, recall the encoding map used in ERoPE (see [Methods](#)), which is given as

$$\varphi_{km\vec{u}} := c_{km} e^{i\omega_k \vec{u} \cdot \vec{r}_m}. \quad (\text{S1})$$

The dot product for encoding maps at atomic positions \vec{r}_m and \vec{r}_n is given as

$$\langle \varphi_{km\vec{u}}, \varphi_{kn\vec{u}} \rangle = c_{kmn} e^{i\omega_k \vec{u} \cdot \vec{r}_{mn}}, \quad (\text{S2})$$

where $c_{kmn} := c_{km} \bar{c}_{kn} \in \mathbb{C}$. The complex exponential can be written in terms of spherical harmonic functions using plane wave expansion

$$e^{i\omega_k \vec{u} \cdot \vec{r}_{mn}} = 4\pi \sum_{\ell=0}^{\infty} \sum_{M=-\ell}^{\ell} i^\ell j_\ell(\omega_k r_{mn}) \bar{Y}_\ell^M(\vec{u}) Y_\ell^M(\hat{r}_{mn}), \quad (\text{S3})$$

where j_ℓ is the Bessel function and Y_ℓ^M is a spherical harmonic of degree ℓ and order M . The integral for a single component of order M and degree ℓ in the spherical harmonics vector \mathbf{Y} in [Eq. 13](#) is given as

$$I = \frac{1}{4\pi} \int_{S^2} d\vec{u} \langle \varphi_{km\vec{u}}, \varphi_{kn\vec{u}} \rangle Y_{\vec{\ell}}^{\bar{M}}(\vec{u}). \quad (\text{S4})$$

Using [Eq. S3](#) one can obtain a compact expression for the integral

$$I_{\vec{\ell}}^{\bar{M}} = c_{kmn} \sum_{\ell=0}^{\infty} \sum_{M=-\ell}^{\ell} i^\ell \int_{S^2} d\vec{u} j_\ell(\omega_k r_{mn}) \bar{Y}_\ell^M(\vec{u}) Y_\ell^M(\hat{r}_{mn}) Y_{\vec{\ell}}^{\bar{M}}(\vec{u}) \quad (\text{S5})$$

$$= c_{kmn} \sum_{\ell=0}^{\infty} \sum_{M=-\ell}^{\ell} i^\ell j_\ell(\omega_k r_{mn}) Y_\ell^M(\hat{r}_{mn}) \delta_{M\bar{M}, \ell\vec{\ell}} \quad (\text{S6})$$

$$= i^{\vec{\ell}} c_{kmn} j_{\vec{\ell}}(\omega_k r_{mn}) Y_{\vec{\ell}}^{\bar{M}}(\hat{r}_{mn}), \quad (\text{S7})$$

where we used the orthonormality relation of the spherical harmonics $\int_{S^2} d\vec{u} \bar{Y}_\ell^{\bar{M}} Y_\ell^M = \delta_{M\bar{M}, \ell\vec{\ell}}$. Evaluation of the integral for all orders $M = -\ell, \dots, +\ell$ for a given degree ℓ in \mathbf{Y} and subsequent concatenation can be written using the tensor product

$$I_{\vec{\ell}} = \frac{1}{4\pi} \int_{S^2} d\vec{u} \langle \varphi_{km\vec{u}}, \varphi_{kn\vec{u}} \rangle \otimes \mathbf{Y}_{\vec{\ell}} \quad (\text{S8})$$

$$= i^{\vec{\ell}} c_{kmn} j_{\vec{\ell}}(\omega_k r_{mn}) \otimes \mathbf{Y}_{\vec{\ell}}(\hat{r}_{mn}), \quad (\text{S9})$$

where Y_ℓ is a vector containing all spherical harmonics with degree ℓ in Y . Repeating this for all different $\ell = 0, \dots, \ell_{\max}$ in Y recovers Eq. 14 from the main text, which puts equivariant SO(3) convolutions into relation to the equivariant EFA update. It should be noted, that the above derivation assumed \mathbf{x} to be an invariant representation ($\ell = 0$), for the convenience of notation and to highlight the fundamental relation between the dot product and the spherical harmonics. In the setting of invariant features, the output of EFA equals equivariant filters (e.g. used in SpookyNet³⁰). However, EFA similarly works (and is actually implemented) for equivariant features \mathbf{x} , where technical implementation details are explained in the [Methods](#) section. This corresponds to full SO(3) convolutions as e.g. found in NequIP³¹.

For the invariant case $\ell = 0$, the spherical harmonics vector Y_ℓ is simply a scalar 1 (when using Racah normalization) and Eq. S9 simplifies to

$$I_0 = c_{kmn} j_0(\omega_k r_{mn}). \quad (\text{S10})$$

As the zeroth Bessel function j_0 is identical to the sinc function, the general solution correctly recovers the integral solution for the invariant case from above.

We verify numerically, that the output of the Euclidean fast attention implementation recovers the functional form of the first three Bessel functions (see Fig. S2A). As long as the integration grid is sufficiently large, the output is equivalent to the expected functional behavior up to numerical precision (see Fig. S2B). The larger the degree ℓ in Y_ℓ , the sooner the output starts to deviate from the exact solution. However, the effect of increasing ℓ only marginally reduces the threshold distance when numerical precision (here 10^{-5}) is violated. The driving source for deviation from the exact solution remains the maximum distance in the data.

Other Symmetrization Operations

As argued in the main text, evaluating the product of encoding maps

$$\langle \varphi_{km\vec{u}}, \varphi_{kn\vec{u}} \rangle \sim e^{i\omega_k \vec{u} \cdot \vec{r}_{mn}} = e^{i\omega_k r_{mn} \cos(\theta)}, \quad (\text{S11})$$

depends on the angle θ between vector \vec{u} and displacement vector $\vec{r}_{mn} = \vec{r}_m - \vec{r}_n$. This dependency is equivalent to a sensitivity to the choice of \vec{u} , which effectively defines a reference frame for each pairwise atomic interaction. Thus, the output is sensitive to the orientation θ between positions and \vec{u} , which gives atomic representations that are neither invariant nor equivariant under global rotation of the input positions.

Instead of integration over the unit sphere, invariance w. r. t. global rotations can be achieved via other symmetrization operations which are described in the following.

Lattice Vectors In the setting of periodic systems, lattice vectors can be used to define a reference frame as e.g. done in Ewald message passing (MP).⁶⁷ In the setting of present lattice vectors and of maximal degree $L_Y = 0$ for the spherical harmonics vector, the symmetrization of the EFA update can be written as

$$\text{EFA}_{\text{lattice}}(\mathcal{X}) = \sum_{p=1}^3 \mathbf{B}_{\vec{u}_p}, \quad (\text{S12})$$

where $\vec{u}_p = \vec{l}_p$ and $\vec{l}_p \in \{\vec{l}_1, \vec{l}_2, \vec{l}_3 \mid \vec{l}_p \in \mathbb{R}^3\}$ are the lattice vectors and $\mathbf{B}_{\vec{u}_p}$ is defined as in Eq. 22. Although this ensures invariance w. r. t. global rotations it does induce a dependence w. r. t. the lattice

vectors. As such, generalization to the very same structures simply oriented differently w. r. t. the cell (as in the training data) and other super cells is likely to be scrutinized. Even for a single system, as soon as entering realistic simulation settings (e.g. constant pressure simulations) lattice vectors undergo non-trivial changes in length and orientation, questioning the applicability of models trained on fixed lattice vectors.

Canonicalization of the Input Geometry In the absence of lattice vectors, another option is to define the vectors $\vec{u}_p = \vec{b}_p$ as the right singular vectors $\{\vec{b}_1, \vec{b}_2, \vec{b}_3 \mid \vec{b}_p \in \mathbb{R}^3\}$ of the singular value decomposition (SVD) of the input geometry $\mathbf{R} \in \mathbb{R}^{N \times 3}$. However, SVD is only defined up to a sign and basis vectors are non-unique for degenerate eigenvalues, which can appear in symmetric systems. As for lattice vectors, this results in learned interactions which can be expected to be hard to apply in generalization or realistic simulation settings.

Degree Re-Scaling

We want to highlight, that all experiments reported in the main body of the text were performed without degree re-scaling. For very small values of ω_{\max} , the numerical difference between different distance values becomes small. Further the absolute values of the Bessel functions for $\ell > 0$ (Fig. S3) are close to zero. We empirically find that this can slow down training due to a vanishing signal at initialization, also capable of affecting overall performance for a finite number of gradient steps. This problem can be mitigated by re-scaling the entries in the integration degree-wise. A simple strategy is to divide each degree in the spherical harmonics vectors \mathbf{Y} by the standard deviation of the radial Bessel function over the interval from 0 to ω_{\max} (for $\ell > 0$) and subtract the mean (computed over the same interval) and divide by the standard deviation for $\ell = 0$. Subtraction by the mean can only be performed for the invariant part since it would destroy equivariance for $\ell > 0$.

The benefit for the training dynamics is highlighted in Fig. S9, where we compare EFA augmented models for $L_{\text{EFA}} = 0$ and $L_Y = 2$ with and without degree scaling. It should be noted, that the EFA-augmented models for cumulene in the main body of the text used different maximal degrees of $L_{\text{EFA}} = 2$ and $L_Y = 0$, but here we choose degrees with the aim of highlighting the potential benefit of degree re-scaling. Other model hyperparameters are the same as for the EFA augmented cumulene models used in the main body of the text.

S_N2 Reactions - SchNet

For the SchNet model we use the same feature dimensions and follow the implementation from the original publication¹², which results in an equal parameter number of 255k for SchNet and SchNet+EFA. The results are shown in Tab. S4. It highlights the ability of EFA to improve the performance across MPNN backbones.

	MP _{5Å}	MP _{10Å}	MP _{5Å} +EFA
MAE _{Energy}	72.1 (±0.3)	38.6 (±1.5)	2.1 (±0.3)
MAE _{Forces}	13.5 (±0.1)	3.2 (±0.0)	1.6 (±0.1)

Table S1 | **Performance on the S_N2 dataset.** Mean absolute errors (MAEs) for energy and forces in meV and meV/Å. MAEs are reported with local cutoffs of 5 Å (with and without EFA augmentation) and 10 Å. The term in brackets denotes the standard deviation over three different training runs (with different random seeds). Best model shown in bold.

Data set	r_{\max}	N_{tot}	N_{train}	N_{valid}	B	N_{epochs}	λ_E	λ_F
k -chains _{$k=2$}	15 Å	2	2	0	2	1000	-	-
k -chains _{$k=4$}	25 Å	2	2	0	2	1000	-	-
k -chains _{$k=6$}	35 Å	2	2	0	2	1000	-	-
k -chains _{$k=8$}	45 Å	2	2	0	2	1000	-	-
Pair	30 Å	10k	3500	500	10	1000	0.01	0.99
Cluster _{$N=16$}	5 Å	2500	2000	500	32	5000	0.01	0.99
Cluster _{$N=32$}	10 Å	2500	2000	500	32	5000	0.01	0.99
Charge-Dipole	10 Å	10k	2500	500	10	3000	0.01	0.99
S _N 2	20 Å	452k	405k	5000	32	500	0.01	0.99
Cumulene	15 Å	4973	1500	500	5	2000	0.01	0.99
Dimer	15 Å	4612	4500	250	16	6000	0.50	0.50

Table S2 | **Data set information and training hyperparameters.** The columns denote the maximal separation in the data set between atoms r_{\max} rounded to the next highest multiple of 5 (this value is used in Eq. 24 to determine the maximal frequency), the total number of data points N_{tot} in the data set, the total number of points used for training N_{train} including the number of points N_{valid} used for validation to select the best model during training, the batch size B , the number of epochs N_{epochs} , and the trade-off parameters for energy λ_E and forces λ_F in the loss function. As k -chains is a classification task we train it via softmax cross-entropy loss.

N_Ω	50	86	110	146	194	230	266	302	350	434	590	770	974	6000
b_{\max}	π	2π	2.5π	3π	4π	4.5π	5π	5.5π	6.5π	7.5π	9π	11π	12.5π	35π

Table S3 | **Lebedev number vs. maximal value for b .** Number of Lebedev points and the maximal possible value b_{\max} according to Eq. 24 from the Methods section in the main text. Values have been determined via comparison to the analytic solution for maximal degree of the spherical harmonics vector $L_Y = 0$. Note that a larger maximal degree $L_Y > 0$, slightly decreases the maximal value for b as also shown in Fig. S2.

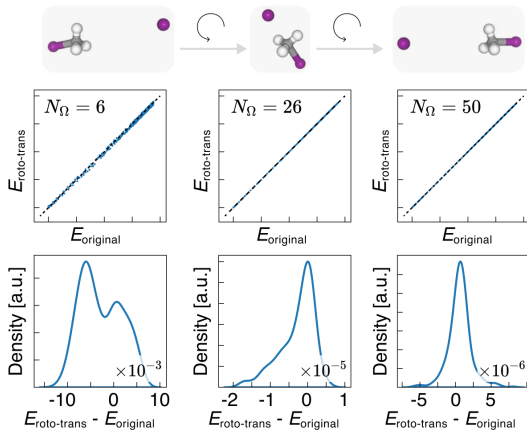


Figure S1 | **Verification of a Model's Rotational and Translation Invariance.** Analysis of the invariance of an MP+EFA model after training, which uses the EFA block for learning long-range dependencies from the data. In total 400 geometries with different spatial extend (varying separation between the Iodine atom (purple) and the central Hydrogen-Carbon complex) are randomly rotated and translated and compared against the model prediction for the non-rotated and non-translated structure. For a Lebedev grid sufficiently large, one finds perfect agreement (roto-translational invariance) up to numerical noise.

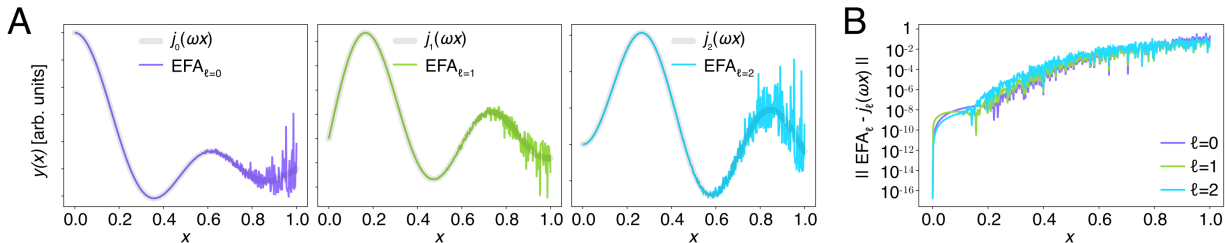


Figure S2 | **Comparison of Euclidean fast attention (EFA) with the analytic integral solution.** (A) Comparison of the output of EFA and the first three Bessel functions j_ℓ for increasing degree ℓ which are the theoretical solution of the surface integral (Eq. S9). (B) Deviation between the Bessel function and the output of EFA as function of x . For the frequency we chose $\omega = 4\pi$ and the number of Lebedev grid points was chosen to be $N_\Omega = 50$.

	SchNet _{5Å}	SchNet _{10Å}	SchNet _{5Å} + EFA
MAE _{Energy}	76.3 (± 0.0)	38.7 (± 0.7)	2.0 (± 0.2)
MAE _{Forces}	18.1 (± 0.1)	3.7 (± 0.1)	1.8 (± 0.1)

Table S4 | **Performance on the S_N2 dataset for SchNet.** Mean absolute errors (MAEs) for energy and forces in meV and meV/Å. MAEs are reported for SchNet with local cutoffs of 5 Å (with and without EFA augmentation) and 10 Å (without EFA augmentation). The term in brackets denotes the standard deviation over three different training runs (with different random seeds). Best model shown in bold.

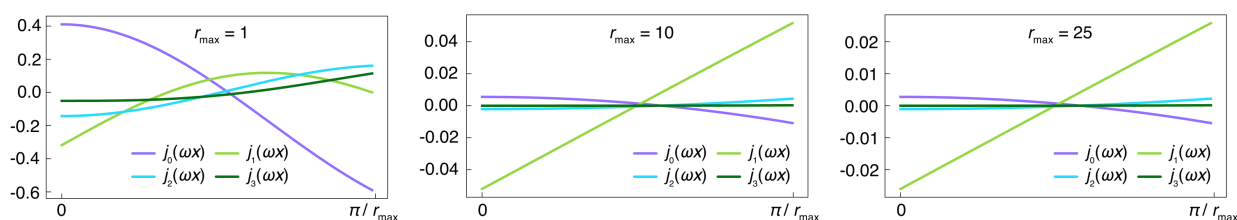


Figure S3 | **Radial Bessel function for varying maximal distance.** Shape of the Bessel function of different degree ℓ and varying value of r_{\max} , which changes the maximal frequency value ω_{\max} as described in the [Methods](#) section (Eq. 24).

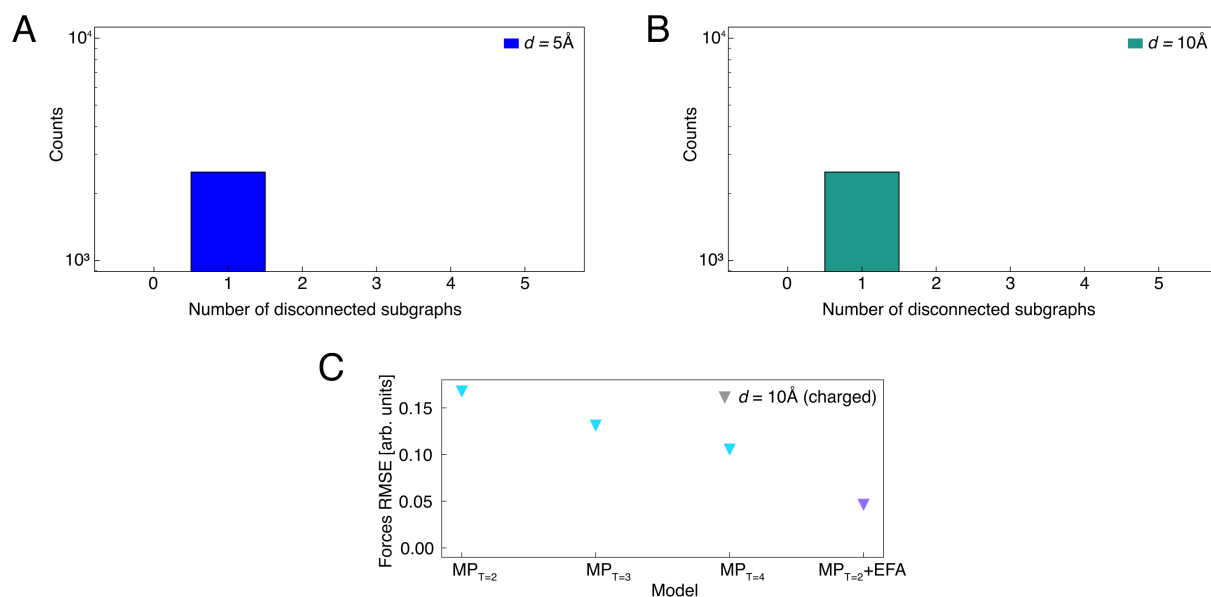


Figure S4 | **Analysis of 3D cluster experiments.** (A) Number of disconnected sub-graphs in the $d = 5 \text{ \AA}$ system. A single disconnected sub-graphs means all nodes are connected via some path. Threshold for two nodes to share an edge is $r_{\text{cut}} = 5 \text{ \AA}$. (B) Number of disconnected sub-graphs in the $d = 10 \text{ \AA}$ system. (C) Results for $N = 32$ systems with partial charges q_i located on the atoms.

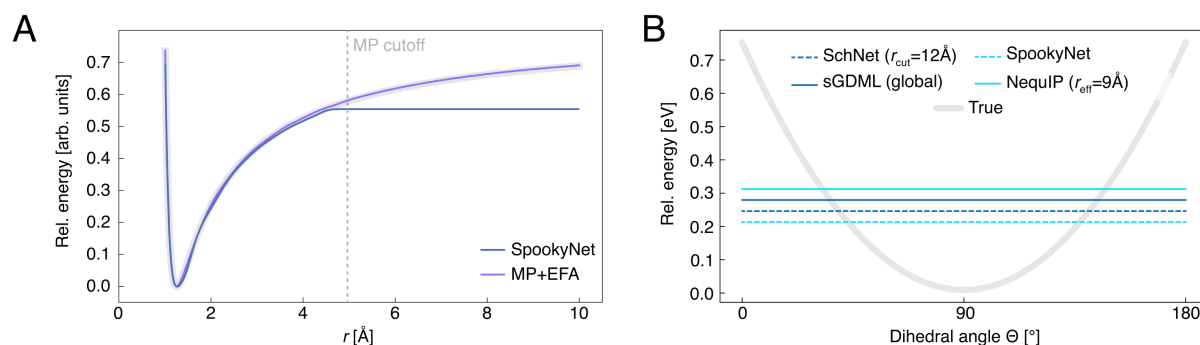


Figure S5 | **Comparison to other models.** (A) Pairwise potential for $N = 2$ atoms system using SpookyNet with linear scaling attention and MP+EFA. (B) Learned energy profile on the cumulene structure using different approaches from the literature.

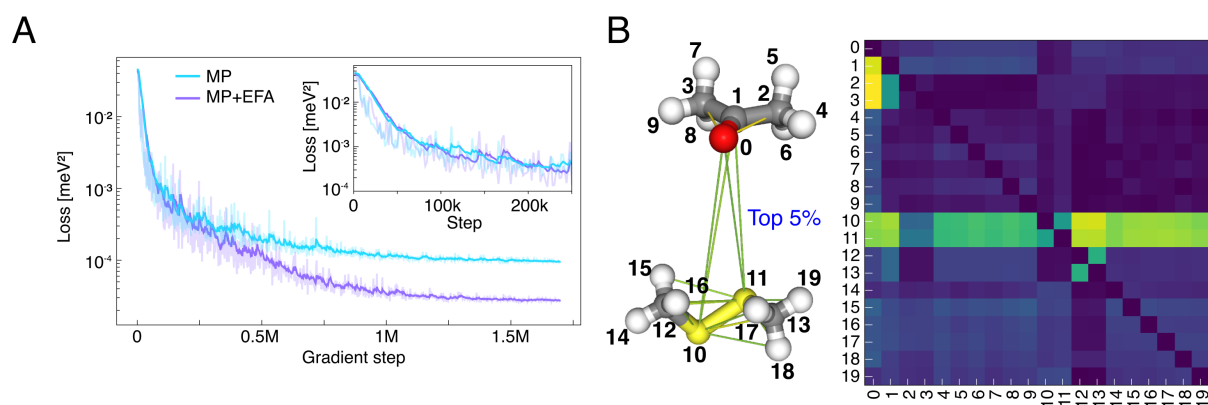


Figure S6 | **Training dynamics and attention analysis.** (A) Loss as a function of gradient step for MP and MP+EFA on the dimer data set. Inset shows the loss over the first 250k steps. (B) Visualization of the learned attention map for a randomly selected dimer. In the 3D view, the pairwise attention values which belong to the largest 5% in magnitude are shown.

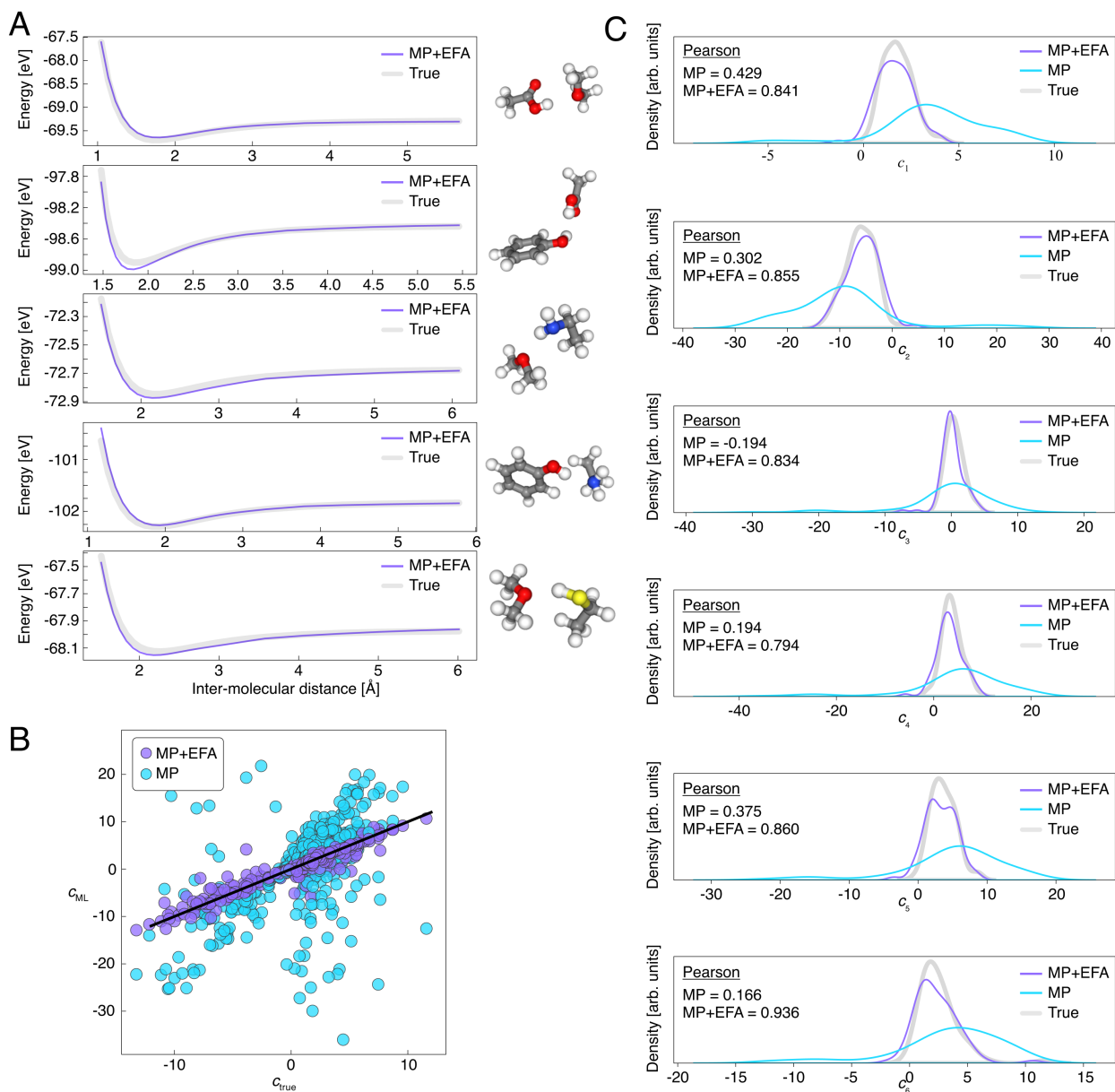


Figure S7 | **Analysis of dimer experiments.** (A) Energy profile for five completely unknown dimers that have not been part of the training data. (B) Scatter plot of the coefficients fitted to the long range tail of the true energy profiles vs. the coefficients fitted to the prediction of the message passing (MP) model and the MP + Euclidean fast attention (EFA) model. (C) Individual distribution for each fitted coefficient c_1, \dots, c_6 . Additionally the Pearson correlation coefficient between the coefficients fitted to the true energy profile and the MP and MP+EFA model is reported as inset.

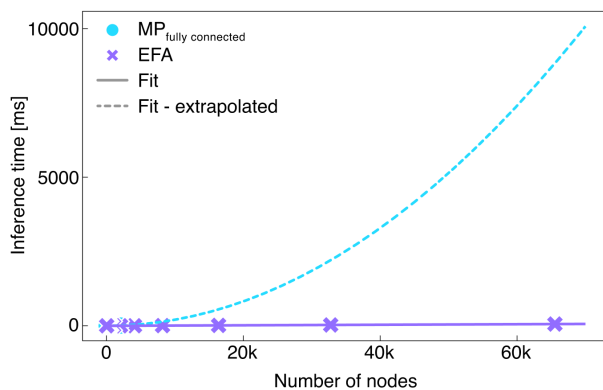


Figure S8 | **Scaling analysis.** Scaling of Euclidean fast attention (EFA) and message passing (MP) on a fully connected graph. For the fully connected MP, the computation and memory complexity scales quadratic leading to out of memory for graphs with more than 2048 nodes. EFA scales linear in both, time and memory which allows scaling to tenth of thousands of nodes.

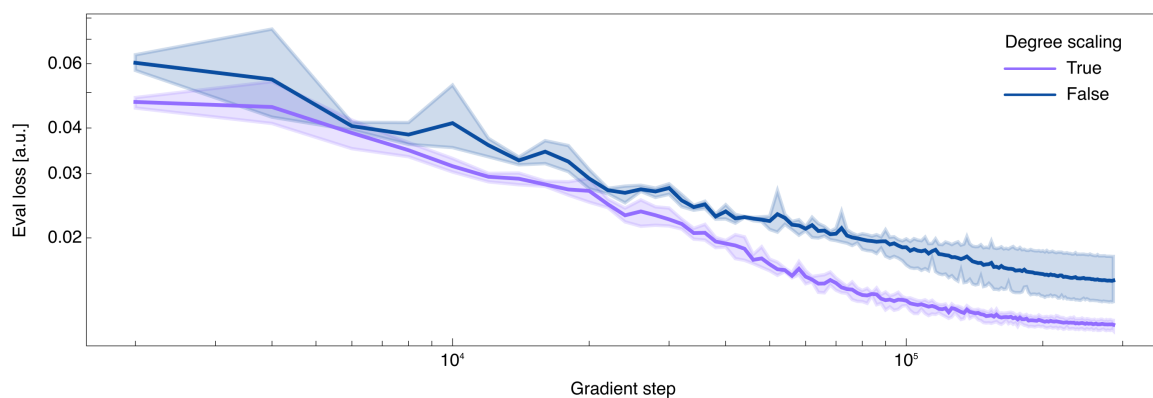


Figure S9 | **Degree Re-Scaling in Euclidean Fast Attention (EFA).** Evaluation loss as a function of the gradient step for EFA augmented models with and without degree re-scaling as described in section [Degree Re-Scaling](#).