



Cite this: *J. Mater. Chem. A*, 2026, **14**, 7628

# Accelerating discovery through integration: a DFT validated machine learning framework for screening MOF photocatalysts

Marco Anselmi,<sup>a</sup> Gregory Slabaugh,<sup>\*b</sup> Rachel Crespo-Otero <sup>\*c</sup> and Devis Di Tommaso <sup>\*a</sup>

The discovery of Metal–Organic Framework (MOF) photocatalysts for CO<sub>2</sub> reduction is hindered by the computational cost of quantum chemical screenings. To overcome this barrier, we introduce a Machine Learning (ML)-accelerated workflow that integrates the speed of ML with the accuracy of Density Functional Theory (DFT). While a DFT-based screening of over 20 000 MOFs identified 105 promising candidates in nearly a month, a ML-driven approach using the Molecular Graph Transformer (MGT) required only 4.5 hours. Here, we present a quantitative assessment of ML performance compared with hybrid DFT for MOF electronic screening, showing that prediction errors are related to the chemistry of the MOFs. We therefore derive an error-aware ML candidate selection strategy that raises DFT candidate recovery from 20% to 70% while keeping a sensible selection set. Building on this, we propose a practical ML to DFT screening workflow in which ML serves as a fast pre-filter to define a small subset for hybrid DFT evaluation, enabling efficient discovery of promising MOFs.

Received 3rd October 2025  
Accepted 2nd January 2026

DOI: 10.1039/d5ta08107f

rsc.li/materials-a

## 1. Introduction

Since the industrial revolution, atmospheric carbon dioxide (CO<sub>2</sub>) levels have been rising consistently due to the burning of fossil fuels, deforestation, the production of cement and many other practices.<sup>1</sup> The rising concentration of atmospheric CO<sub>2</sub> and its impact on Earth's climate have prompted researchers to take action to reduce emissions. Carbon capture and utilization (CCU) techniques are important for bridging the gap between our high-carbon industries and a sustainable, low-carbon future. In CCU processes, CO<sub>2</sub> is upgraded into value added chemicals and materials,<sup>2–7</sup> rather than accumulating in the atmosphere, thereby enabling a circular economy.<sup>8</sup> Photochemical CO<sub>2</sub> reduction (PCO<sub>2</sub>R) is a promising technique within CCU.<sup>9,10</sup> The process mimics natural photosynthesis, where sunlight is used to drive chemical reactions. In PCO<sub>2</sub>R, photocatalysts absorb sunlight and use this energy to power the thermodynamically uphill CO<sub>2</sub> reduction reaction (CO<sub>2</sub>RR) into valuable chemicals such as formic acid (HCOOH), methanol (CH<sub>3</sub>OH), and methane (CH<sub>4</sub>).<sup>3–5,11,12</sup> Given that, in one hour, the sun provides Earth with the same amount of energy that

humanity consumes in one year, it makes the sun one of the best sources of clean energy.<sup>13</sup>

A photocatalytic reaction such as PCO<sub>2</sub>R generally occurs in three steps.<sup>14–16</sup> (i) Photon absorption: in the first step, photons with energy equal to or greater than the optical gap energy are absorbed by the photocatalyst. This absorption excites electrons from the occupied orbitals to the unoccupied orbitals, creating holes in the highest occupied molecular orbital. (ii) Electron–hole dynamics: in the second step, the excited electrons and the created holes can either recombine within the catalyst or migrate to the surface. (iii) Surface reactions: in the third step, the electrons and holes that reach the surface participate in oxidation and reduction reactions.

The electronic structure of a photocatalyst is a key factor in determining its suitability for PCO<sub>2</sub>R, particularly the positions of the Highest Occupied Molecular Orbital (HOMO) and Lowest Unoccupied Molecular Orbital (LUMO), as well as the bandgap between them. To enable direct comparison with the redox potentials of the Oxygen Evolution Reaction (OER) and the CO<sub>2</sub> Reduction Reaction (CO<sub>2</sub>RR) for a specific product, the HOMO and LUMO energy levels are converted to Ionization Potential (IP) and Electron Affinity (EA), respectively, by alignment with the vacuum energy level. This comparison assesses whether the photocatalyst has the necessary energetic properties for the CO<sub>2</sub>RR to take place.<sup>17,18</sup>

Metal–organic frameworks (MOFs) are a class of crystalline porous materials constructed from organic linkers and metal ions/clusters, known as secondary building units (SBUs). The various combinations of these components can lead to an almost infinite

<sup>a</sup>Department of Chemistry, Queen Mary University of London, Mile End Rd, Bethnal Green, London, E1 4NS, UK. E-mail: d.ditommaso@qmul.ac.uk

<sup>b</sup>Digital Environment Research Institute, Queen Mary University of London, Empire House, 67–75 New Road, London, E1 1HH, UK. E-mail: g.slabaugh@qmul.ac.uk

<sup>c</sup>Department of Chemistry, University of College London, 20 Gordon Street, London, WC1H 0AJ, UK. E-mail: r.crespo-otero@qmul.ac.uk



number of structures with diverse functionalities. Through fine-tuning, it is possible to create and screen materials with easily accessible catalytic sites and desirable electronic structures.<sup>19–21</sup> Due to their unique structural properties, such as high surface area, high porosity, tuneable morphology, and adjustable chemical composition, MOFs are highly suitable for various CCU applications,<sup>22</sup> including photocatalytic CO<sub>2</sub> reduction.<sup>5–7,23,24</sup>

Nevertheless, the practical deployment of MOF photocatalysts is currently hindered by significant experimental challenges. Foremost among these are the poor stability of many frameworks in the aqueous environments required for CO<sub>2</sub> reduction and the complexity of synthesis, which limits scalability.<sup>25</sup> From an electronic perspective, efficiency is often stifled by rapid electron–hole recombination and light absorption spectra that are typically confined to the ultraviolet range, a limitation shared with traditional photocatalysts such as TiO<sub>2</sub>.<sup>26</sup> To address these limitations, this study utilizes a computational screening approach centred on band-edge engineering. Unlike previous screenings that utilize broad energetic windows, we use a visible-light based bandgap window of 1.9–2.5 eV. By prioritizing this range, we aim to identify candidates that maximize solar energy utilization efficiency while maintaining the energy requirements to drive the thermodynamically challenging CO<sub>2</sub> reduction.

Quantum chemical methods, typically within the framework of Density Functional Theory (DFT), can be used to determine the electronic structure of MOFs, including the energy positions of the HOMO, LUMO, and band gap. This information is crucial for assessing their potential as PCO<sub>2</sub>R catalysts. However, DFT calculations of band gaps generally require the use of hybrid DFT functionals such as HSE06 (ref. 27) and B3LYP3,<sup>28</sup> rather than standard GGA functionals such as PBE,<sup>29</sup> which tend to underestimate band gaps.<sup>30</sup> The large unit cells required to model MOFs mean that a significant number of atoms must be considered, making the application of DFT, especially hybrid DFT methods, computationally very demanding. These calculations can take weeks to months to complete when applied to large datasets with tens or hundreds of thousands of structures, and because of its complexity, the explicit calculation of the excited states in MOFs hasn't been extensively explored in the literature.<sup>31</sup>

Recently, Machine Learning (ML) models have been gaining popularity as a methodology to accelerate the discovery of new materials. ML models can be used to predict properties of structures, such as formation energy, potential energy surfaces and electronic properties, like the bandgap, of various materials.<sup>32–38</sup> Nevertheless, MOFs have seen less attention in the literature due to the lack of substantial electronic structure–property datasets, with only three contributions found for this study.<sup>39–41</sup>

In a major contribution to the field of ML prediction of MOF bandgaps, Rosen *et al.*<sup>40</sup> introduced a new dataset called the Quantum MOF (QMOF) database and an analysis of the performance of five different ML representations on this dataset, namely the Sine Coulomb matrix,<sup>42</sup> the stoichiometric, the orbital field matrix, the Smooth Overlap of Atomic Positions<sup>43</sup> (SOAP) and the Crystal Graph Convolutional Neural Network<sup>33</sup> (CGCNN) representations. The QMOF database introduced by Rosen *et al.*<sup>40</sup> contains the computed properties for 15 713 MOFs after structure relaxation *via* DFT, which they collected from the Cambridge

Structural Database MOF subset<sup>44</sup> and the 2019 Computationally Ready MOF database,<sup>45</sup> with the aim of creating a database for the development of novel ML algorithms for MOFs. The properties reported included, but were not limited to, relaxed geometries, bandgaps, HOMOs, LUMOs, charge densities, density of states, partial charges, spin densities and bond orders. Since its inception, the database has been expanded to contain the properties for 20 375 MOFs.

Unlike earlier MOF ML models such as CGCNN,<sup>33</sup> which primarily capture local atomic connectivity, our approach employs the Molecular Graph Transformer<sup>46</sup> (MGT) architecture to address the unique challenges posed by MOFs. MGT uses a multi-graph representation that encodes local bonded interactions, many-body effects, and long-range electrostatic interactions through a global graph derived from Coulomb matrix features. This design is particularly suited to MOFs with large unit cells and extended pore networks, where non-bonded interactions strongly influence band edges and band gaps. By combining local attention with message passing across these graphs, MGT achieves improved accuracy (average error 0.34 eV for bandgap, EA, and IP) and enables an error-aware ML to DFT workflow that accelerates screening while maintaining chemical interpretability.

This study involves a high-throughput screening of the QMOF database to identify potential photocatalysts for CO<sub>2</sub> reduction using vacuum-aligned HOMOs and LUMOs (band edges), together with the bandgap, and employs both hybrid DFT and ML techniques throughout. Initially, 105 promising photocatalyst candidates were identified *via* hybrid DFT calculations. The electronic-structure data resulting from these calculations were then used to construct a training dataset for the ML model. A subsequent ML based screening was then performed to predict vacuum-aligned band edges, yielding a set of 45 candidates. The outcomes of DFT and ML high-throughput screening were then compared to assess the feasibility of ML as a stand-alone approach. In doing so, the work provides a systematic evaluation of ML screening *versus* DFT for MOF electronic structure prediction, including an analysis of prediction error with respect to the chemistry of MOFs. Using these insights, we propose an error-informed ML-to-DFT workflow, in which ML serves as a fast pre-filter for the selection of a small subset of candidates for hybrid-DFT evaluation, thereby substantially accelerating the screening speeds of MOF photocatalyst candidates.

## 2. Methods

### 2.1. High-throughput DFT calculations

Fig. 2 shows the process of the screening methodology presented in this study. For this study, the QMOF database was used as the candidate dataset, and the properties reported in it were used for the following steps. The results of the DFT calculations recorded in the QMOF database had been obtained using the GGA PBE functional for DFT and computed using VASP. Further DFT calculations, which were conducted for this study, were also performed using the VASP code (version 6.4.3) with the Perdew–Burke–Ernzerhof<sup>29</sup> (PBE) functional or the Heyd–Scuseria–Ernzerhof<sup>27</sup> (HSE) exchange–correlation functionals and the Projector-Augmented-Wave<sup>47</sup> (PAW) method. The kinetic energy cutoff was



520 eV, and the  $k$ -point mesh used to sample the Brillouin zone of the simulation supercell was determined using the VASPKIT<sup>48</sup> package for each structure.

The VASP calculations were run on a mixture of GPU and CPU nodes on the High Performance Computing (HPC) Hub in Materials and Molecular Modelling Young (MMM Hub Young), on Queen Mary's Apocrita HPC facilities supported by QMUL Research-IT, and on the Sulis Tier 2 HPC platform. The alignment calculations, on the other hand, were run using 8 CPU cores and 16 GB memory on the Apocrita HPC nodes.

**2.1.1. Candidate dataset.** The first step in the process consists of generating a dataset of structurally optimized MOFs. This can be either achieved by compiling structures from external databases, such as the QMOF,<sup>40</sup> the CoreMOF,<sup>45</sup> or the CSD MOF Collection,<sup>44</sup> or it can also be created by generating new structures.

In this step, the collection of DFT computed bandgaps is essential. To accelerate the process, utilizing the generalized gradient approximation (GGA) PBE functional is recommended. Although the PBE functional is less accurate than hybrid functionals like HSE, its computational efficiency makes it preferable for large datasets. Notably, since the only property needed for this step is the bandgap, and a linear relationship exists between the bandgaps obtained using the GGA PBE functional and those obtained using the HSE06 Hybrid functional,<sup>49</sup> this relationship can be expressed as:

$$E_{g,HSE} = 1.09 E_{g,PBE} + 1.04 \text{ eV} \quad (1)$$

Similar linear equations can also be derived for other hybrid functionals, reinforcing the preference for using the PBE functional in this step.

**2.1.2. Bandgap screening.** In the context of single semiconductor photocatalysis, the materials must exhibit specific characteristics that enable efficient light absorption while satisfying the energy requirements for the desired reaction. Specifically, for PCO<sub>2</sub>R to occur, as shown in Fig. 1, the bandgap of the catalyst must exceed the difference between the redox potentials for the OER and the CO<sub>2</sub>RR to specific products.

In this regard, a techno-economic evaluation of PCO<sub>2</sub>R revealed that the production of C1 products is more economically competitive compared to the production of C2 products, which entails higher energy requirements, longer reaction times, and more complex processes that contribute to increased production costs.<sup>50</sup> Consequently, this study focuses on three key C1 products: methane (CH<sub>4</sub>), methanol (CH<sub>3</sub>OH) and formic acid (HCOOH). Among these, the smallest energy gap between the redox potentials for the OER and the reduction to a specific product is 1.88 eV, while the largest gap is 2.25 eV. However, practical considerations necessitate adjusting the required bandgap values. Overpotentials and losses in current efficiencies must be accounted for, leading to higher bandgap requirements for catalyst selection.<sup>17</sup>

Furthermore, the ability of a material to absorb light depends on the energy of incident photons, which should match or be higher than the bandgap of the material. For instance, TiO<sub>2</sub>, one of the most studied photocatalysts, has a bandgap of 3.2 eV, limiting its light absorption to the ultraviolet range, which constitutes less

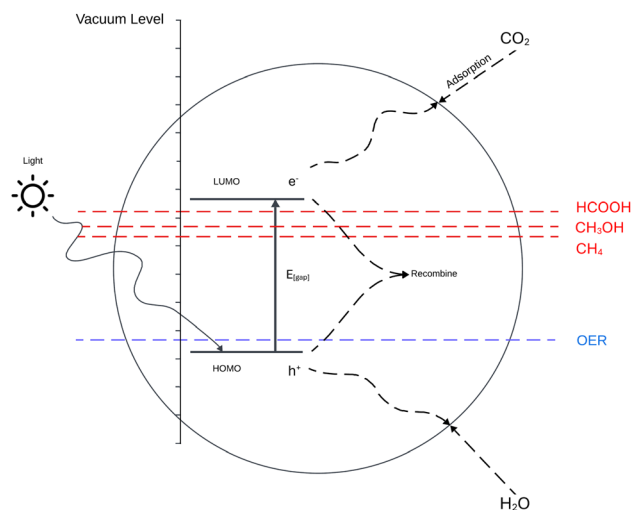


Fig. 1 Schematic showing the steps of the photocatalytic CO<sub>2</sub> reduction on a single catalyst. Red and blue dotted lines show the redox potentials for the CO<sub>2</sub> reduction and water oxidation, respectively. Black dotted lines show the possible paths that the excited electrons and holes can take, and the circle represents the surface of the catalyst.

than 5% of the solar spectrum.<sup>51</sup> Thus, to maximize light absorption, it is essential to tune the bandgap.

Considering all these criteria, an ideal bandgap for efficient photocatalysis would fall within the range of 1.9 eV and 2.5 eV. To account for bandgap underestimation errors caused by the PBE functional used in the calculation of bandgaps in the QMOF database, we adjust the bandgap selection criteria to 0.8 eV and 1.34 eV using eqn (1).

**2.1.3. Higher level DFT.** In the subsequent band alignment step, precise HOMO and LUMO values play a crucial role. If the functional employed in generating the candidate dataset for bandgap computation operates at the GGA level of theory, it is highly advisable to proceed with this step. Although GGA functionals offer computational efficiency and reasonable accuracy,

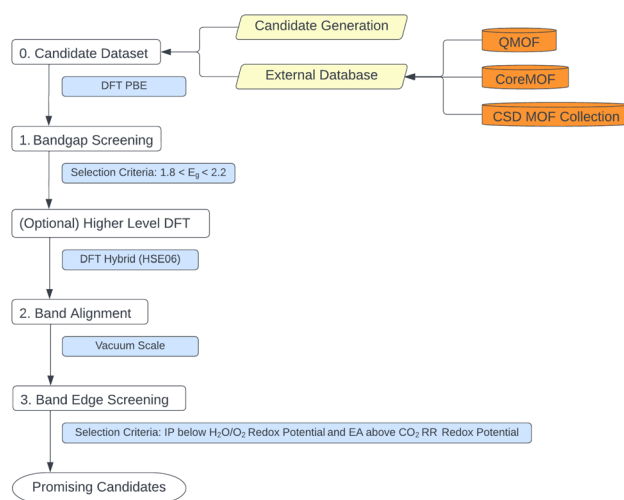


Fig. 2 Screening approach for MOF CO<sub>2</sub>RR photocatalysts.



they tend to underestimate experimental bandgaps,<sup>19</sup> leading to inaccuracies in the derived HOMO and LUMO levels.

To enhance accuracy, a hybrid functional can be created by incorporating a percentage of exact exchange computed using the Hartree-Fock (HF) exchange functional.<sup>52</sup> In the study, single-point calculations were performed for bandgaps, band edges, and local potential using the HSE06 functional. The input structures were obtained from the QMOF database, and the HSE06 setup involved a 25% addition of HF exchange and an inter-electronic range of  $0.2 \text{ \AA}^{-1}$  for applying the HF exchange.

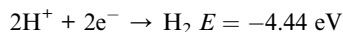
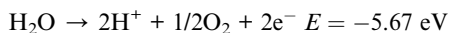
**2.1.4. Band alignment.** For the selection of suitable structures for PCO<sub>2</sub>R, it is crucial to consider the alignment of electronic energy levels. Specifically, the HOMO and LUMO values must straddle the redox potentials associated with both the OER and the CO<sub>2</sub> reduction processes (Fig. 1). However, the use of pseudo-potentials in DFT plane wave calculations causes the HOMO, LUMO, and redox potentials to be defined on different energy scales, necessitating alignment for meaningful comparison.

The methodology outlined in this study addresses this issue by adjusting the bandgap and band edges to the energy level of the vacuum. Traditionally, band alignment involves referencing the energy of the vacuum above the surface of the material. However, for MOFs, direct calculations of the surface electronic structure can be computationally demanding due to their large unit cells containing hundreds of atoms. To overcome this challenge, the approach proposed by Butler *et al.*<sup>53</sup> uses the electrostatic potential at the centre of MOF pores as an approximation for the vacuum level. By using this reference point, electronic energy levels of MOFs can be placed on a common energy scale.

To determine the pore centre required for the vacuum energy approximation, the Pore Size Distribution (PSD) method introduced by Trepte and Schwalbe<sup>54</sup> is employed. This approach utilizes a Monte Carlo procedure, generating random points within the MOF unit cell. At each step, these points are adjusted to maximize their distance from the atoms in the MOF, ultimately identifying the pore center and diameter.

The combination of these two techniques allows consistent alignment of electronic energy levels across MOF structures, facilitating informed material selection for efficient CO<sub>2</sub>RR.

**2.1.5. Band edge screening.** Once the band alignment has been completed, the last step is to select the MOF catalysts. The selection is carried out by comparison of the aligned HOMO and LUMO values with respect to the reduction potentials for the reactions of interest. For example, for water splitting, the conduction and valence band edges (the electron affinity (EA) and the ionization potential (IP)) must straddle the redox potentials of the two reactions involved: the OER and the hydrogen evolution reaction (HER). Specifically, the valence band must be below the OER ( $E = -5.67 \text{ eV}$ ) and the conduction band above the HER ( $E = -4.44 \text{ eV}$ ).<sup>17,18</sup>



In the case of water splitting photocatalytic CO<sub>2</sub>RR, as the redox potentials for the conversion of CO<sub>2</sub> into added value chemicals are higher than those of both the OER and HER, the required bandgap of the MOF will need to be larger than that required for water splitting, but not too wide, in order to allow for the absorption of photons. The bandgap edges will then need to straddle the OER and the redox potential for the CO<sub>2</sub>RR to the desired product, with the valence band below the OER and the conduction band above the CO<sub>2</sub>RR redox potential.



## 2.2. Machine learning

The Molecular Graph Transformer (MGT)<sup>46</sup> is a recent Graph Neural Network (GNN) architecture designed to enhance material property prediction by incorporating long-range interactions alongside local atomic interactions. The model relies on a Molecular Graph Representation (MGR) that encodes the periodic MOF structure into three distinct graphs: the Local Graph ( $G_{\text{local}}$ ), which captures pairwise bond interactions by representing atoms as nodes, with edges formed between atoms within a cutoff radius of  $8 \text{ \AA}$ ; the Line Graph ( $G_{\text{line}}$ ), which captures many-body interactions (bond angles) by using the edges of the local graph as nodes and representing the angles between connected bonds as edges; the Global Graph ( $G_{\text{global}}$ ), which captures non-bonded, long-range electrostatic interactions. This is a fully connected graph in which edges contain features derived from the Coulomb matrix, representing the interaction between atom pairs regardless of distance. The architecture of the model involves encoding nodes and edges from the MGR into feature vectors through encoding layers, followed by passing them through multiple MGT Encoder layers to capture structural information effectively. Each MGT Encoder combines local attention mechanisms with message passing on bond graphs, their line graphs and long range electrostatic graphs, allowing it to explicitly capture long-range electrostatic forces that conventional GNN models often neglect.

**2.2.1. Training details.** To evaluate the performance of the MGT, the results obtained through the DFT methodology described in Section 2.1 were used. Data preprocessing was minimal; fewer than 50 structures were excluded solely due to failed DFT convergence. Statistical outliers were explicitly retained in the dataset to rigorously assess the model's robustness across chemically complex environments. The resulting dataset contains the structures and bandgap, HOMO, LUMO, EA and IP values for the 2169 structures on which Hybrid DFT calculations and band alignments were performed. The model was trained for 200 epochs using the Adam optimizer with a learning rate of  $10^{-4}$ . Due to the memory intensity of the global graph attention, a batch size of 2 was employed with gradient accumulation steps of 8. The network architecture adopts the optimized configuration described by Anselmi *et al.*,<sup>46</sup> comprising a single MGT Encoder layer with three angle convolution layers, three bond convolution



layers, and one electrostatic convolution layer. Full training details are provided in SI Section 1.2. All of the ML training and testing were run on Queen Mary's Apocrita HPC facilities. The training was performed on 2 A100 40 GB GPUs, while the testing and inference were performed using 1 A100 40 GB GPU.

**2.2.1.1. Training and testing.** In this work, the model is trained to predict the bandgap, HOMO, LUMO, EA and IP of a MOF. This is achieved by training the MGT using the labelled data from the screening process described in section 2.1. From the 2169 structures, on which hybrid DFT calculations and band alignments were performed, 1518 have been used for training while the remaining 651 have been used for testing. This 70–30% training-test split was performed randomly to ensure an unbiased evaluation of the model's performance. Furthermore, the 1518 structures for training were further split with a 80–20% ratio for training and validation, respectively. Moreover, to enable comparison of the ML screening with the DFT screening, the structures identified as potential CO<sub>2</sub> reduction catalysts have been incorporated into the dataset used to evaluate the model.

## 3. Results and discussion

### 3.1. DFT

Following the initial bandgap screening detailed in Section 2.1, 2169 potential MOF structures were selected from the 20 375 data points in the QMOF database. Hybrid DFT calculations were then performed on these structures to obtain more accurate bandgaps and band-edges (Table S3). Subsequently, the band alignment step was carried out. The DFT calculations for all 2169 structures required a total of 24 days of CPU time, as determined from the VASP output files.

From these 2169 MOFs, 1132 structures exhibited an EA higher than the CO<sub>2</sub>R redox potential (for the CO<sub>2</sub>RR with the

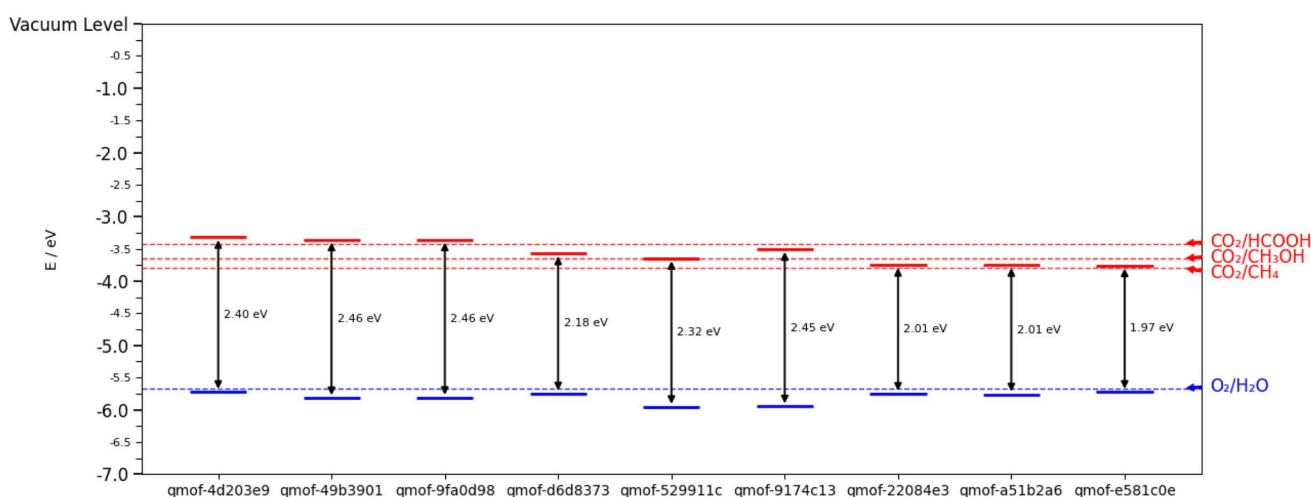
highest gap requirement) and an IP lower than the water OER. However, by applying the ideal bandgap range of 1.9 eV to 2.5 eV (refer to Section 2.1.2), the screening was further refined to identify 105 structures that meet the criteria to be suitable as catalysts for at least one of the CO<sub>2</sub>R reactions, producing either CH<sub>4</sub>, CH<sub>3</sub>OH or HCOOH.

Out of these 105 structures, 61 have been synthesised, with their structures available on the CSD dataset. Among these, qmof-e100550 (ref. 55) presented a experimentally calculated bandgap of 2.23 eV, qmof-d6d8373 (ref. 56) has a reported bandgap of 3.1 eV, lastly for qmof-306fb3e,<sup>57</sup> its bandgap was measured to be 2.49 eV. Meanwhile, qmof-3eec122 (ref. 58) doesn't have a reported bandgap; it has, however, been tested for the electrocatalytic oxygen reduction reaction and for CO<sub>2</sub> trapping. The remaining 44 structures are computationally theorized compounds sourced from various databases.<sup>59,60</sup>

Fig. 3 shows the bandgaps and band edges positions for the structures that are suitable for all three CO<sub>2</sub>R reactions. Among these MOFs, qmof-e581c0e exhibits the lowest bandgap of approximately 1.97 eV, allowing for greater light absorption. Among these, 13 are suitable for the reduction to formic acids, 42 structures are suitable for the reduction to methanol, while the remaining 43 are suitable for the reduction to methane.

### 3.2. ML

The training and testing took around 4.5 hours. Evaluated on a test set of 651 structures, as shown in Fig. 4 and Table S1, the predictions of the model are in general agreement with the DFT-calculated values, demonstrating a general high accuracy of the ML approach, with an average error of 0.34 eV across all three properties. Nevertheless, despite the low average error, a visual analysis of its performance (see Fig. 4) reveals significant outliers. This indicates that the model can be substantially inaccurate for certain structures, which can introduce the risk



**Fig. 3** Band-edges and bandgaps obtained using the HSE06 functional. HOMO levels are indicated by the blue lines and LUMO levels by the red lines. The dashed horizontal lines denote the redox potentials for water splitting (O<sub>2</sub>/H<sub>2</sub>O) and CO<sub>2</sub> reduction to formic acid (CO<sub>2</sub>/HCOOH), methanol (CO<sub>2</sub>/CH<sub>3</sub>OH), and methane (CO<sub>2</sub>/CH<sub>4</sub>). The displayed MOFs are illustrative examples randomly selected for meeting selection criteria for specific CO<sub>2</sub> reduction reactions: qmof-4d203 × 10<sup>9</sup>, qmof-49b3901, and qmof-9fa0d98 are suitable for HCOOH production; qmof-06d8373, qmof-529911c, and qmof-9174c13 for CH<sub>3</sub>OH production; and qmof-22084 × 10<sup>3</sup>, qmof-a51b2a6, and qmof-e581c0e for CH<sub>4</sub> production.



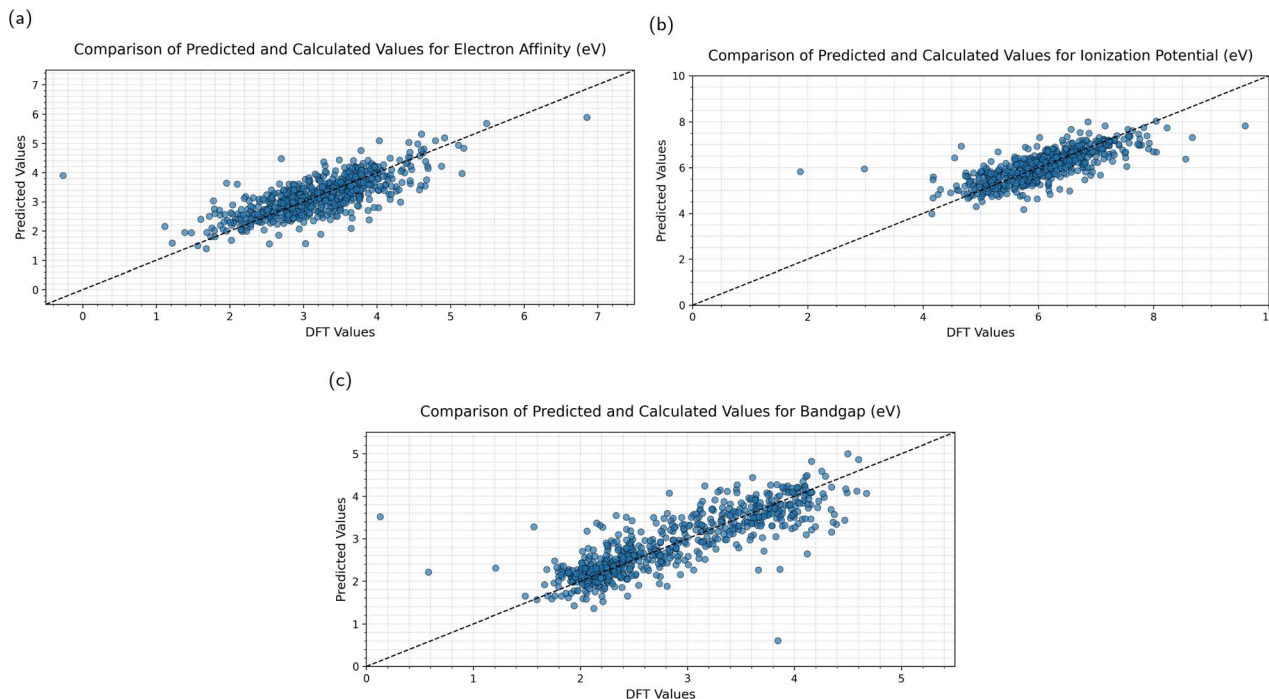


Fig. 4 Parity plots showing the performance of the MGT model on the test set. The plots compare the DFT-calculated ground truth values against the ML-predicted values for (a) electron affinity (EA), (b) ionization potential (IP) and (c) bandgap. The black dashed line in each plot represents perfect agreement.

of both false positives and false negatives in a high-throughput screening workflow.

By further analysing the chemical nature of the structures that exhibit high-error predictions, it is revealed that the cases in which the model fails are not random but related to chemical complexity. This analysis identifies three specific structural features in MOFs that consistently lead to large prediction errors for key electronic properties including IP, EA, and bandgap: (i) the presence of multi-metal centres (*e.g.* qmof-6e89a67 with CuNa, qmof-cfee5a8 with Cu<sub>2</sub>Yb<sub>2</sub>, and qmof-8281b85 with Cd<sub>2</sub>Ni<sub>2</sub>); (ii) the inclusion of heavy d-block and f-block elements; and (iii) the incorporation of halides in the ligands.

A detailed breakdown of error distribution by metal type (Fig. 5 and Table S2) further demonstrates the correlation between the inaccuracies of the model and the chemical environment. Alkali metals show a trend where the average error on unseen data (0.33 eV) is lower than on seen data (0.41 eV). This is likely due to the fact that these metals are very similar to each other and all show stable oxidation states of either 0 or +1. Meanwhile, transition metals, which comprise the majority of the dataset, consistently show a mediocre performance (0.42 eV) for both seen and unseen data, suggesting that the model struggles to capture the variable oxidation states and diverse coordination geometries of d-block metals.

Post-transition metals are similar in performance to transition metals, but they also show a clear trend of being more reliable for lighter metals (*e.g.* Al with 0.28 eV error) than for heavy ones (*e.g.* Ga with 0.74 eV error). Lastly, lanthanides

display the largest disparity: they have the lowest error of any group during training (0.18 eV) but the highest error when unseen (1.09 eV). This is likely due to the fact that lanthanide complexes have a wide range of coordination geometries, which indicates that the model is likely overfitting to the specific coordination environments of seen lanthanide clusters. Thus, the model is effectively memorizing the geometry rather than learning the underlying physics of f-orbitals, leading to significant failures when encountering new elements such as ytterbium. Therefore, the significant outliers in the parity plots (Fig. 4) are primarily associated with these complex electronic phenomena, ranging from the challenge of d-orbital directionality to the poor generalisation of f-electron correlation.

From the results obtained using the MGT model on the test set for this study, 43 MOFs were identified as suitable catalysts for the CO<sub>2</sub>R reactions producing either CH<sub>4</sub>, CH<sub>3</sub>OH, or HCOOH. Out of the 43 selected structures, 20 have been synthesized and studied in experimental papers. The remaining 23 are computationally theorized compounds. For all 43 structures, their bandgaps and band edges have been plotted in relation to the redox potentials for the OER and CO<sub>2</sub>R reactions (Fig. 6). Among these MOFs, qmof-8b03127 exhibits the lowest bandgap of approximately 2.06 eV, allowing for greater light absorption. Among all the structures, 5 are suitable for the reduction to formic acids, 22 are suitable for the reduction to methanol, while the remaining 16 are suitable for the reduction to methane.



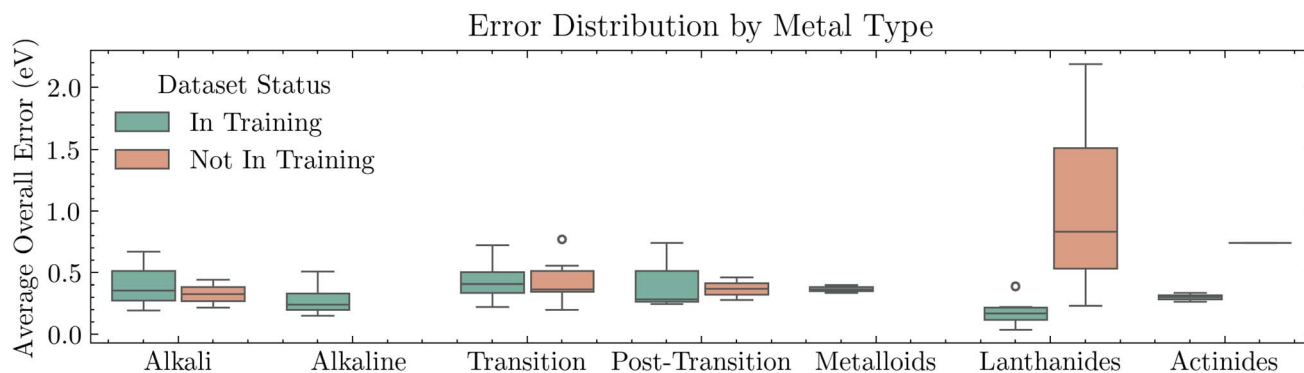


Fig. 5 Distribution of average overall prediction error (eV) by metal type. The boxplots compare the model's performance on metals present in the training set (green) versus those present only in the test set (orange). Alkali metals show better performance on unseen data, suggesting high generalisability of simple ionic interactions. In contrast, lanthanides exhibit overfitting behaviour, achieving the lowest error in training (0.18 eV) but the highest error on unseen species (1.09 eV), indicating that the model memorises specific coordination geometries rather than learning f-orbital physics.

### 3.3. Comparison of methodologies

When performing a screening with a ML model, however, it is important to note that any machine learning model will have a certain amount of error with respect to the results obtained using DFT. From the results reported in the study of the MGT, and the testing results obtained in this study, the average error for bandgap, EA and IP predictions of the MGT is 0.34 eV. Given these prediction errors, a practical workflow must account for them to avoid discarding viable candidates. We find that using the raw ML predictions is a poor strategy, as it recovers only 20% of the promising MOFs identified by DFT. A more effective approach is to apply an error margin ( $\delta$ ) around the selection criteria, creating a widened window for candidate selection. For the bandgap ( $E_g$ ) screening, the bandgap based on DFT results should satisfy  $E_{g,\min} \leq E_g \leq E_{g,\max}$ , and the ML screening would accept candidates satisfying:

$$E_{g,\text{pred}} > E_{g,\min} - \delta \text{ and } E_{g,\text{pred}} < E_{g,\max} + \delta \quad (2)$$

Similarly, for the band edges, where the IP must be lower than the water OER ( $V_{\text{OER}}$ ) and the EA higher than the  $\text{CO}_2$ R redox potential ( $V_{\text{CO}_2\text{R}}$ ), the ML screening criteria can be expanded to accept candidates satisfying:

$$\text{IP}_{\text{pred}} < V_{\text{OER}} + \delta \text{ and } \text{EA}_{\text{pred}} > V_{\text{CO}_2\text{R}} - \delta \quad (3)$$

By accounting for the model's uncertainty in this manner, the screening becomes more in line with the DFT-based screening. As shown in Table 1, adjusting the selection criteria by the error of the model ( $\delta = 0.34$  eV), successfully recovers 69% of the DFT-selected structures. This highlights a clear trade-off between computational cost and scientific accuracy. Regarding the types of DFT hits missed by ML, we find that DFT selected candidates are most often missed by ML in

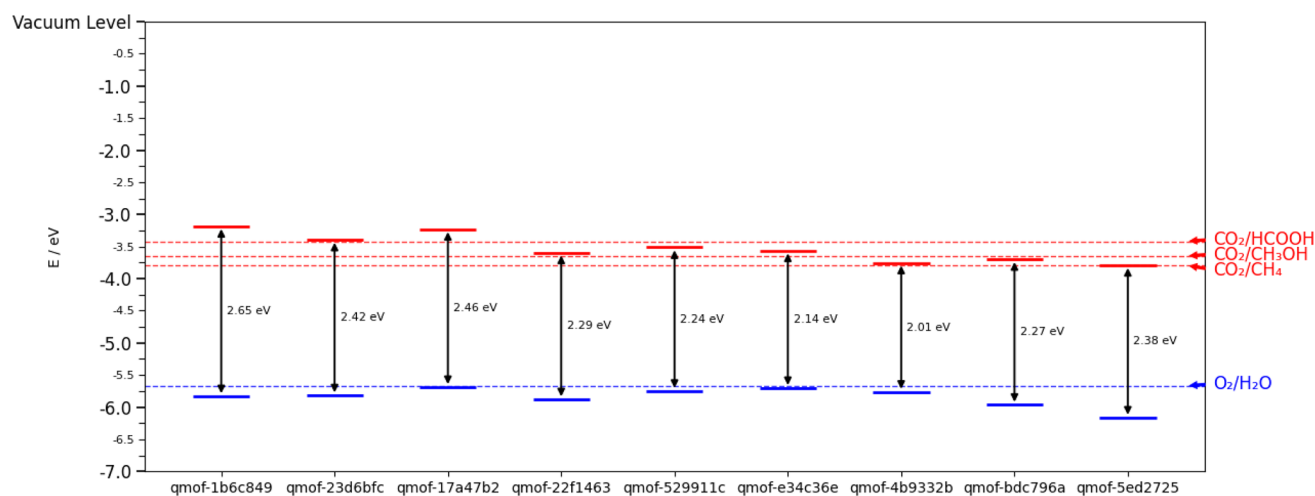


Fig. 6 Band edges and bandgaps obtained using the MGT model, for representative MOFs meeting selection criteria for specific  $\text{CO}_2$  reduction reactions. HOMO levels are indicated by the blue lines and LUMO levels by the red lines. The dashed horizontal lines denote the redox potentials for water splitting ( $\text{O}_2/\text{H}_2\text{O}$ ) and  $\text{CO}_2$  reduction to formic acid ( $\text{CO}_2/\text{HCOOH}$ ), methanol ( $\text{CO}_2/\text{CH}_3\text{OH}$ ), and methane ( $\text{CO}_2/\text{CH}_4$ ). The displayed MOFs are illustrative examples: qmof-1b6c849, qmof-23d6bfc, and qmof-17a47b2 are suitable for  $\text{HCOOH}$  production; qmof-22f1463, qmof-529911c, and qmof-e34c36e for  $\text{CH}_3\text{OH}$  production; and qmof-4b9332b, qmof-bdc796a, and qmof-5ed2725 for  $\text{CH}_4$  production.

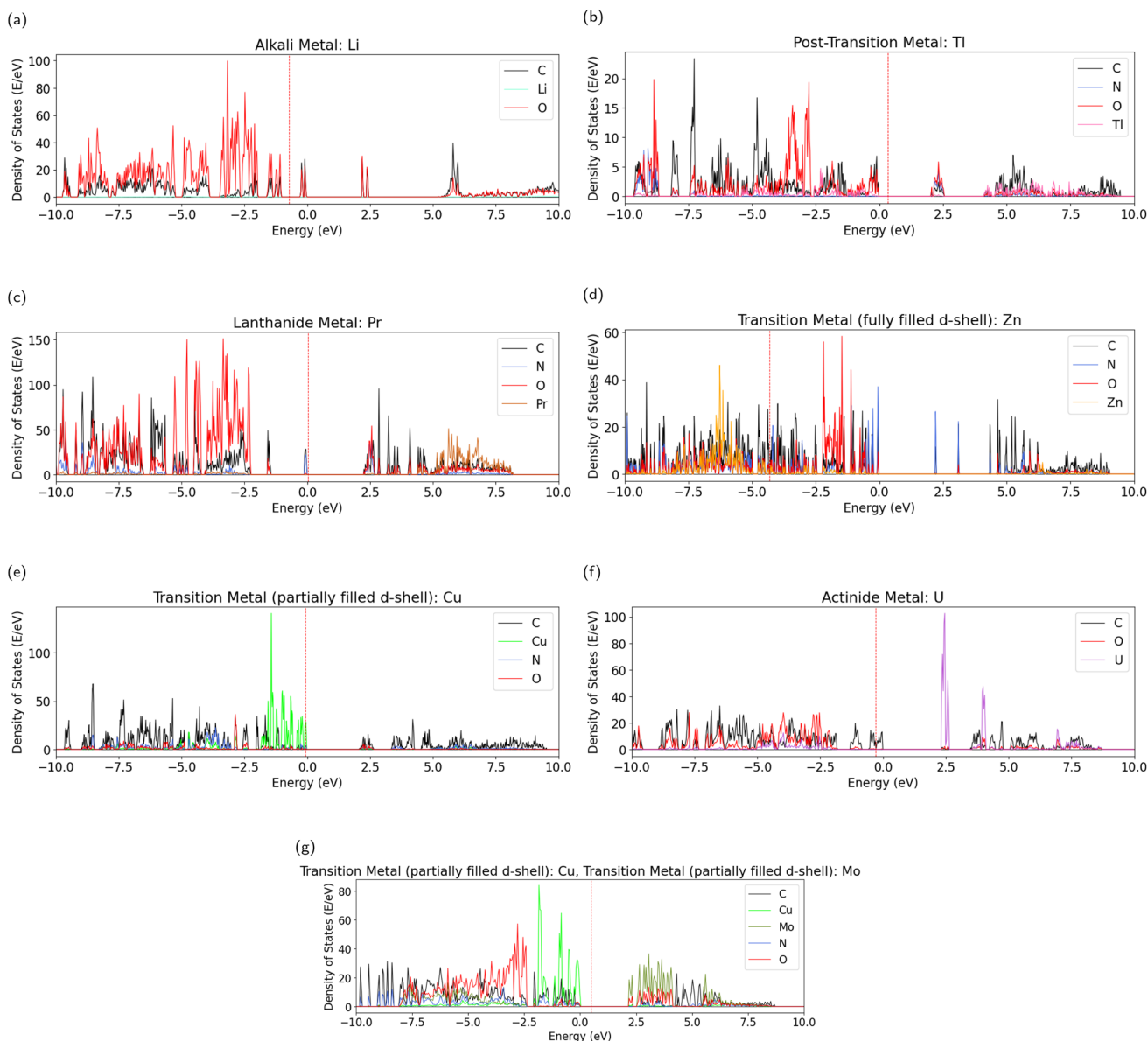


**Table 1** Comparison of screening strategies, highlighting the trade-off between the number of selected candidates and the recovery rate ("hit percentage") of structures selected using the DFT-based screening. Using a wider error margin significantly increases the likelihood of finding promising candidates at the cost of a larger initial pool for follow-up calculations

Screening strategy	Selection criteria	Candidates selected (% of test set)	DFT candidate recovery ("hit percentage")
Raw ML predictions	Use unchanged selection criteria	43 (7%)	21 (20%)
ML error margin	Widen window by the error of the model	204 (31%)	73 (69%)
Conservative margin	Widen window by 0.5 eV	293 (45%)	94 (89%)
High-recall margin	Widen window by 1.0 eV	470 (72%)	103 (98%)

two regimes: (i) near threshold band edge cases, where vacuum aligned IP/EA lie within the model's MAE of the redox or band gap limits, leading to classification flips unless thresholds are

widened by  $\delta \approx 0.34$  eV (see eqn (2) and (3) and Table 1); and (ii) chemically complex secondary building units (SBUs), notably transition metals with partially filled d shells, lanthanides/



**Fig. 7** Projected density of states (pDOS) for selected MOFs incorporating different metal centers: (a) alkali metal lithium (Li), (b) post-transition metal thallium (Tl), (c) lanthanide metal praseodymium (Pr), (d) group 12 metal zinc (Zn), (e) transition metal copper (Cu), (f) actinide metal uranium (U) and lastly (g) double transition metals copper (Cu) and molybdenum (Mo). In MOFs with non-transition metals such as (a)–(c), the DOS at the band edges is mostly composed of orbitals from the organic linkers (carbon, nitrogen, and oxygen). In contrast, in MOFs containing transition metals (d)–(g) or actinides (f), a significant contribution from metal-derived orbitals to the DOS at or near the band edges is observed.



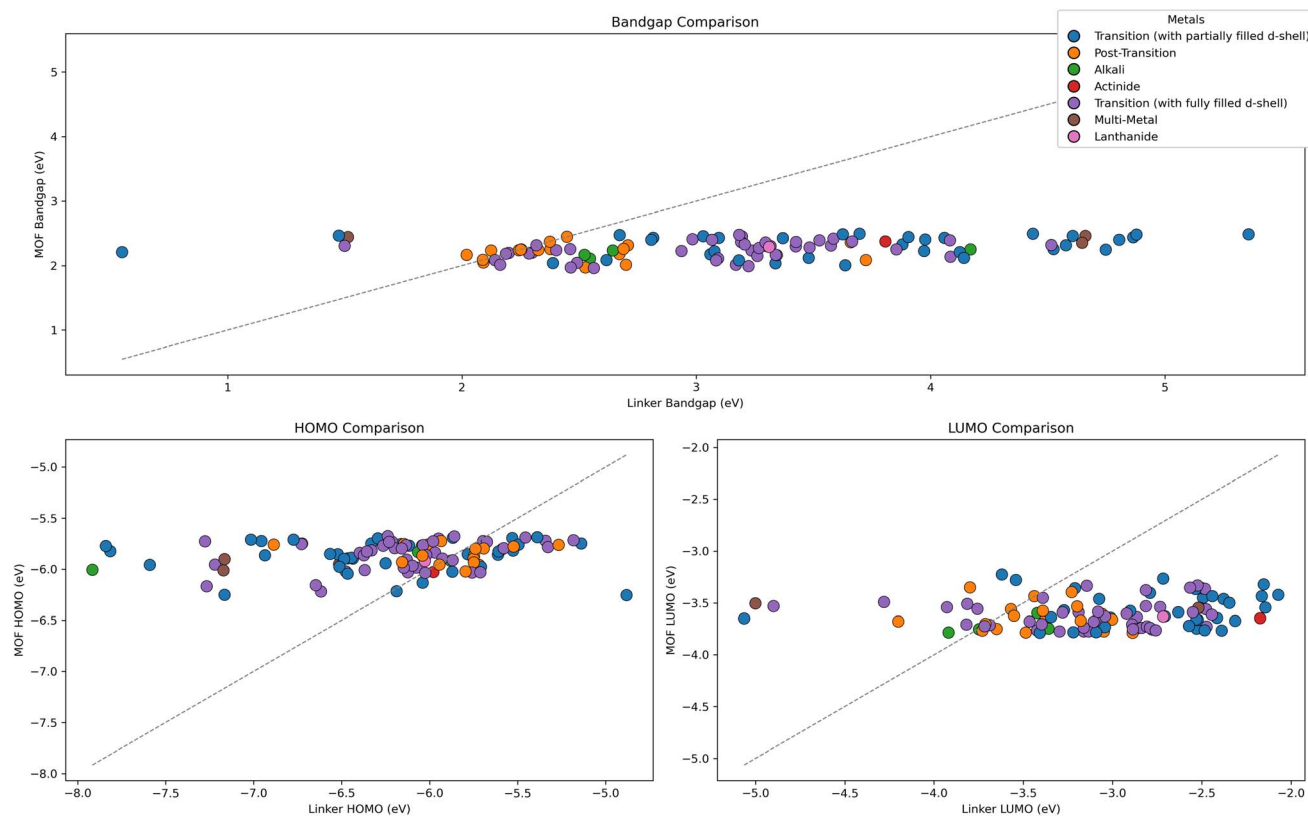


Fig. 8 Comparison of the bandgap, HOMO, and LUMO levels of the MOFs versus their isolated linkers. The lack of a strong linear correlation highlights that the final electronic structure is not dictated solely by the linker. The type of metal cation causes distinct deviations from the baseline linker properties, with transition metals, actinides, and lanthanides showing the most significant influence.

actinides, multi metal centres, and halide bearing ligands, where pDOS indicates metal dominated edges and error diagnostics show poor generalisation. Accordingly, ML only screening is generally reliable for alkali/light post transition SBU, single metal frameworks and for candidates whose predicted edges sit well inside the selection window, whereas hybrid DFT validation is advisable for f block/actinide systems, multi metal nodes, TM dominated edges, and borderline alignment cases.

Beyond the selection of structures, a more fundamental methodological difference lies in the accuracy and internal consistency of the electronic property values obtained from DFT versus ML. In DFT calculations, for instance, the bandgap is derived directly from the computed HOMO and LUMO levels, and the alignment to the vacuum level also utilizes these HOMO and LUMO values. This interconnectedness ensures that all calculated electronic properties are consistent; a shift in the HOMO or LUMO energy propagates through to the bandgap, IP, and EA. In contrast, the machine learning approach predicts each property (HOMO, LUMO, bandgap, IP, and EA) independently. As a result, each prediction has its own distinct error relative to the reference DFT values. This independence can lead to discrepancies: the bandgap predicted directly by the ML model may differ from a bandgap calculated as the difference between the predicted HOMO and LUMO levels. Similarly, these values might also diverge from a bandgap inferred from the

predicted EA and IP. Due to the nature of machine learning algorithms, these inconsistencies cannot be entirely eliminated, although they can be reduced by developing more accurate models.

These inconsistencies, however, lead to another problem. When evaluating candidate photocatalysts, it is often necessary to consider more than just the bandgap, HOMO, and LUMO levels. For instance, our current study reveals that the Density of States (DOS) at the HOMO and LUMO levels is predominantly localized on the organic linkers when the metals are alkali (Fig. 7a), post-transition metals (Fig. 7b), lanthanides (Fig. 7c), or Group 12 metals (Fig. 7d). Conversely, in MOFs with transition metals featuring partially filled d-shells (Fig. 7e), the DOS at the HOMO can shift significantly onto the metal center (Fig. 7). MOFs with Actinide metal centers (Fig. 7f) appear to localize the LUMO DOS onto the metal. Notably, when multiple transition metals with partially filled d-shells are present (Fig. 7g), the DOS at both the HOMO and LUMO levels shifts onto the metal centers.

This observation highlights the critical role of the organic linker in defining the foundational electronic structure of a MOF, a concept explored by Grau-Crespo *et al.* for zeolitic imidazolate frameworks (ZIFs).<sup>61</sup> Our findings align with their conclusion that the positions of the band edges are largely determined by the linkers. However, while Grau-Crespo *et al.* also reported a direct correlation between the bandgaps of the



isolated linkers and the energy difference between the Highest Occupied Crystal Orbital (HOCO) and the Lowest Unoccupied Crystal Orbital (LUCO) of the corresponding ZIFs, this specific relationship is not apparent in our broader dataset (Fig. 8). Nevertheless, the principle that metals modify these linker-defined band edges to achieve a desired electronic structure holds true. The comparison plots (Fig. 8) clearly show that transition metals and multi-metal systems cause the largest deviations from the linker-only properties, affecting the bandgap, HOMO, and LUMO. Actinides and lanthanides also show a strong influence, particularly on the LUMO energy, while having less effect on the HOMO. In contrast, MOFs with alkali and post-transition metals stick closest to the electronic properties of their linkers, generally showing the smallest influence. As also suggested by Grau-Crespo *et al.*, transition metals can play a crucial role by providing alternative sites for excited electrons or holes, thereby promoting electron–hole separation and increasing the lifetime of the excited state by reducing recombination rates. Grau-Crespo *et al.* also discussed that charge separation, and thus potentially slower recombination, can also be achieved in ZIFs by using multiple linker types, where one type primarily contributes to the HOMO and another to the LUMO.

Performing a comparable detailed analysis of linker contributions and metal–ligand interactions using machine learning presents significant challenges. It would require the prediction of additional properties, each introducing its own error margin. Furthermore, obtaining detailed electronic data for isolated ligands would likely require a separate ML model trained specifically for that purpose. The accumulation of errors from these multiple steps could render the overall comparison too uncertain to yield any meaningful conclusions.

Design trends across the 105 DFT-selected and 43 ML-selected MOFs reveal qualitatively a linker first, metal tune paradigm: organic linkers broadly set the baseline band edge positions, while the choice of metal centres (especially d and f block elements and multi metal nodes) introduces controllable deviations that tune the HOMO/LUMO alignment and the band gap within the visible light window. Alkali and light post transition metals preserve linker dominated edges and simplify alignment to OER/CO<sub>2</sub>RR potentials, whereas transition metals and actinides shift edge densities onto the metal, enabling band edge tuning and charge separation at the cost of greater variability that requires hybrid DFT validation. Halide-bearing ligands and multi-metal centres correlate with higher ML uncertainty, so error aware selection windows are essential when ML is used as a pre filter.

## 4. Conclusions

Given the promise of MOFs as tuneable materials for photocatalytic CO<sub>2</sub> reduction, it is essential to develop efficient screening strategies to identify candidates suitable for CCU technologies. This study critically compared high-throughput DFT calculations with ML techniques for identifying novel MOF photocatalysts. A crucial aspect for both methodologies is the accurate alignment of electronic energy levels (HOMO,

LUMO, IP, and EA) to a common vacuum reference, essential for comparing against the OER and CO<sub>2</sub> redox potentials. The DFT screening, starting from the QMOF database and followed by more accurate hybrid HSE06 (ref. 27) calculations for 2169 candidates, identified 105 promising MOFs. However, while providing high confidence and consistency, this DFT approach was resource-intensive, taking 24 days of computation time using 3 to 8 GPUs per node. On the other hand, the ML approach using the MGT<sup>46</sup> model trained with data from 1518 DFT-calculated structures identified 43 distinct MOF candidates in approximately 4.5 hours. However, this study reveals that the speed of ML comes with critical trade-offs. The largest errors of the model are not random but systematic, concentrating in MOFs with complex chemistries, such as those containing multiple metal centers or heavy elements. Furthermore, because the ML model predicts each property independently of the other, potential internal inconsistencies can arise (*e.g.*, the predicted bandgap not matching the difference between the predicted LUMO and HOMO).

Therefore, directly replacing DFT with ML is not yet feasible. Instead, an integrated workflow is essential. ML should serve as a high-speed engine for initial screening, the effectiveness of which hinges on the screening strategy. We demonstrate that applying a margin based on the error of the model to the selection criteria is crucial for creating a candidate pool that is both computationally feasible and can include significant percentage of candidates that would be selected with a DFT screening. To build more confidence into the candidate screening process, ML models could also be enhanced with uncertainty quantification. Methods such as Monte Carlo dropout and model ensembling can show how certain a prediction is, allowing the selection of potential candidates with greater certainty. DFT calculation can then be performed on these ML-identified candidates for validation, accurate electronic structure determination, and detailed analysis. This workflow can leverage speed of the MGT and accuracy of hybrid DFT calculations.

Looking forward, the full potential of ML in materials science will be unlocked by developing models capable of preserving physical interdependencies between properties. For instance, using a consistency loss that enforces a relationship between the predictions could potentially reduce inconsistencies and provide better bandgap predictions. Until then, an integrated workflow, which leverages the speed of ML and the accuracy of DFT, offers the most powerful and pragmatic strategy to accelerate the design and discovery of next-generation photocatalysts for a sustainable future.

## Conflicts of interest

There are no conflicts to declare.

## Data availability

The datasets generated and/or analyzed during the current study are available in the supplementary information (SI) and from the corresponding author on reasonable request.



Supplementary information is available. See DOI: <https://doi.org/10.1039/d5ta08107f>.

## Acknowledgements

M. A. acknowledges the mini-centre-for Doctoral-Training in CO<sub>2</sub>-conversion at QMUL for a PhD scholarship. D. D. T. acknowledges the Leverhulme Trust (RPG-2023-239) for resources supporting projects on advancing materials characterization by computational means. This work was supported by the Engineering and Physical Sciences Research Council [grant number EP/Y009800/1], through funding from Responsible AI UK (KP0016). This project has received funding from the AI for Chemistry: Alchemy Hub (EPSRC grant EP/Y028775/1 and EP/Y028759/1). Calculations were performed using the Sulis Tier 2 HPC platform hosted by the Scientific Computing Research Technology Platform at the University of Warwick, on the MMM Hub Young Tier 2 HPC platform, part of the Materials and Molecular Modelling Hub, and on the JADE 2 Tier 2 HPC platform hosted by the Hartree Center at the University of Oxford. Sulis is funded by EPSRC Grant EP/T022108/1 and the HPC Midlands Plus consortium. JADE 2 is funded by EPSRC (EP/T022205/1). MMM Hub Young is partially funded by EPSRC (EP/T022213/1, EP/W032260/1 and EP/P020194/1).

## References

- 1 M. Office, *Mauna Loa Carbon Dioxide Forecast for 2022, 2022*, <https://www.metoffice.gov.uk/research/climate/seasonal-to-decadal/long-range/forecasts/co2-forecast>.
- 2 Q. Wang, C. Cai, M. Dai, J. Fu, X. Zhang, H. Li, H. Zhang, K. Chen, Y. Lin, H. Li, J. Hu, M. Miyauchi and M. Liu, *Small Sci.*, 2021, **1**, 2000028.
- 3 J. Fu, K. Jiang, X. Qiu, J. Yu and M. Liu, *Mater. Today*, 2020, **32**, 222–243.
- 4 M. A. Tekalgne, H. H. Do, A. Hasani, Q. V. Le, H. W. Jang, S. H. Ahn and S. Y. Kim, *Mater. Today Adv.*, 2020, **5**, 100038.
- 5 J. H. Qin, P. Xu, Y. D. Huang, L. Y. Xiao, W. Lu, X. G. Yang, L. F. Ma and S. Q. Zang, *Chem. Commun.*, 2021, **57**, 8468–8471.
- 6 S. S. A. Shah, T. Najam, M. Wen, S. Q. Zang, A. Waseem and H. L. Jiang, *Small Struct.*, 2022, **3**, 2100090.
- 7 Y. Zhao, L. Zheng, D. Jiang, W. Xia, X. Xu, Y. Yamauchi, J. Ge and J. Tang, *Small*, 2021, **17**, 2006590.
- 8 B. M. Tackett, E. Gomez and J. G. Chen, *Nat. Catal.*, 2019, **2**(5), 381–386.
- 9 S. Fang, M. Rahaman, J. Bharti, E. Reisner, M. Robert, G. A. Ozin and Y. H. Hu, *Nat. Rev. Methods Primers*, 2023, **3**, 61.
- 10 S. C. Shit, I. Shown, R. Paul, K. H. Chen, J. Mondal and L. C. Chen, *Nanoscale*, 2020, **12**, 23301–23332.
- 11 T. Billo, I. Shown, A. kumar Anbalagan, T. A. Effendi, A. Sabbah, F. Y. Fu, C. M. Chu, W. Y. Woon, R. S. Chen, C. H. Lee, K. H. Chen and L. C. Chen, *Nano Energy*, 2020, **72**, 104717.
- 12 S. Wu, H. Pang, W. Zhou, B. Yang, X. Meng, X. Qiu, G. Chen, L. Zhang, S. Wang, X. Liu, R. Ma, J. Ye and N. Zhang, *Nanoscale*, 2020, **12**, 8693–8700.
- 13 G. W. Crabtree and N. S. Lewis, *Phys. Today*, 2007, **60**, 37–42.
- 14 C. Li, Y. Xu, W. Tu, G. Chen and R. Xu, *Green Chem.*, 2017, **19**, 882–899.
- 15 R. Cauwenbergh and S. Das, *Green Chem.*, 2021, **23**, 2553–2574.
- 16 M. Ingham, A. Aziz, D. D. Tommaso and R. Crespo-Otero, *Mater. Adv.*, 2023, **4**, 5388–5419.
- 17 M. G. Walter, E. L. Warren, J. R. McKone, S. W. Boettcher, Q. Mi, E. A. Santori and N. S. Lewis, *Chem. Rev.*, 2010, **110**, 6446–6473.
- 18 S. Hamad, N. C. Hernandez, A. Aziz, A. R. Ruiz-Salvador, S. Calero and R. Grau-Crespo, *J. Mater. Chem. A*, 2015, **3**, 23458–23465.
- 19 J. L. Mancuso, A. M. Mroz, K. N. Le and C. H. Hendon, *Chem. Rev.*, 2020, **120**, 8641–8715.
- 20 B. J. Burnett, P. M. Barron, C. Hu and W. Choe, *J. Am. Chem. Soc.*, 2011, **133**, 9984–9987.
- 21 W. Lu, Z. Wei, Z. Y. Gu, T. F. Liu, J. Park, J. Park, J. Tian, M. Zhang, Q. Zhang, T. Gentle, M. Bosch and H. C. Zhou, *Chem. Soc. Rev.*, 2014, **43**, 5561–5593.
- 22 S. E. M. Elhenawy, M. Khraisheh, F. AlMomani, G. Walker, S. E. M. Elhenawy, M. Khraisheh, F. AlMomani and G. Walker, *Catalysts*, 2020, **10**, 1–33.
- 23 X. Li and Q. L. Zhu, *EnergyChem*, 2020, **2**, 100033.
- 24 Y. S. Bae and R. Q. Snurr, *Angew. Chem., Int. Ed.*, 2011, **50**, 11586–11596.
- 25 M. Khan, Z. Akmal, M. Tayyab, S. Mansoor, A. Zeb, Z. Ye, J. Zhang, S. Wu and L. Wang, *Carbon Capture Sci. Technol.*, 2024, **11**, 100191.
- 26 C. I. Ezugwu, S. Liu, C. Li, S. Zhuiykov, S. Roy and F. Verpoort, *Coord. Chem. Rev.*, 2022, **450**, 214245.
- 27 J. Heyd, G. E. Scuseria and M. Ernzerhof, *J. Chem. Phys.*, 2003, **118**, 8207–8215.
- 28 A. D. Becke, *J. Chem. Phys.*, 1993, **98**, 1372–1377.
- 29 J. P. Perdew, K. Burke and M. Ernzerhof, *Phys. Rev. Lett.*, 1996, **77**, 3865.
- 30 P. Borlido, T. Aull, A. W. Huran, F. Tran, M. A. Marques and S. Botti, *J. Chem. Theory Comput.*, 2019, **15**, 5069–5079.
- 31 M. Ingham, M. Brady and R. Crespo-Otero, *J. Chem. Theory Comput.*, 2025, **21**, 7576–7592.
- 32 K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O. A. V. Lilienfeld, K. R. Müller and A. Tkatchenko, *J. Phys. Chem. Lett.*, 2015, **6**, 2326–2331.
- 33 T. Xie and J. C. Grossman, *Phys. Rev. Lett.*, 2018, **120**, 145301.
- 34 K. Choudhary and B. DeCost, *npj Comput. Mater.*, 2021, **7**(1), 185.
- 35 J. Gasteiger, J. Groß and S. Günnemann, *8th International Conference on Learning Representations, ICLR*, 2020.
- 36 K. T. Schütt, H. E. Sauceda, P. J. Kindermans, A. Tkatchenko and K. R. Müller, *J. Chem. Phys.*, 2018, **148**, 241722.
- 37 S. S. Omeel, S. Y. Louis, N. Fu, L. Wei, S. Dey, R. Dong, Q. Li and J. Hu, *Patterns*, 2022, **3**, 100491.
- 38 S. Zhang, Y. Liu and L. Xie, *Sci. Rep.*, 2023, **13**, 19171.



- 39 Y. He, E. D. Cubuk, M. D. Allendorf and E. J. Reed, *J. Phys. Chem. Lett.*, 2018, **9**, 4562–4569.
- 40 A. S. Rosen, S. M. Iyer, D. Ray, Z. Yao, A. Aspuru-Guzik, L. Gagliardi, J. M. Notestein and R. Q. Snurr, *Matter*, 2021, **4**, 1578–1597.
- 41 K. Choudhary, T. Yildirim, D. W. Siderius, A. G. Kusne, A. McDannald and D. L. Ortiz-Montalvo, *Comput. Mater. Sci.*, 2022, **210**, 111388.
- 42 F. Faber, A. Lindmaa, O. A. V. Lilienfeld and R. Armiento, *Int. J. Quantum Chem.*, 2015, **115**, 1094–1101.
- 43 A. P. Bartók, R. Kondor and G. Csányi, *Phys. Rev. B: Condens. Matter Mater. Phys.*, 2013, **87**, 184115.
- 44 P. Z. Moghadam, A. Li, S. B. Wiggin, A. Tao, A. G. Maloney, P. A. Wood, S. C. Ward and D. Fairen-Jimenez, *Chem. Mater.*, 2017, **29**, 2618–2625.
- 45 Y. G. Chung, E. Haldoupis, B. J. Bucior, M. Haranczyk, S. Lee, H. Zhang, K. D. Vogiatzis, M. Milisavljevic, S. Ling, J. S. Camp, B. Slater, J. I. Siepmann, D. S. Sholl and R. Q. Snurr, *J. Chem. Eng. Data*, 2019, **64**, 5985–5998.
- 46 M. Anselmi, G. Slabaugh, R. Crespo-Otero and D. D. Tommaso, *Digit. Discov.*, 2024, **3**, 1048–1057.
- 47 P. E. Blöchl, *Phys. Rev. B: Condens. Matter Mater. Phys.*, 1994, **50**, 17953.
- 48 V. Wang, N. Xu, J. C. Liu, G. Tang and W. T. Geng, *Comput. Phys. Commun.*, 2021, **267**, 108033.
- 49 A. S. Rosen, V. Fung, P. Huck, C. T. O'Donnell, M. K. Horton, D. G. Truhlar, K. A. Persson, J. M. Notestein and R. Q. Snurr, *npj Comput. Mater.*, 2022, **8**(1), 112.
- 50 Y. Cui, A. Labidi, X. Liang, X. Huang, J. Wang, X. Li, Q. Dong, X. Zhang, S. I. Othman, A. A. Allam, D. W. Bahnemann and C. Wang, *ChemSusChem*, 2024, **17**, e202400551.
- 51 J. Wu, Y. Huang, W. Ye, Y. Li, J. Wu, Y. Huang, W. Ye and Y. Li, *Advanced Science*, 2017, **4**, 1700194.
- 52 S. Abedi, M. T. Ahmadpour, S. Baninajarian, H. Kahnouji, S. J. Hashemifar, Z. K. Han and S. V. Levchenko, *J. Chem. Phys.*, 2023, **158**, 184109.
- 53 K. T. Butler, C. H. Hendon and A. Walsh, *J. Am. Chem. Soc.*, 2014, **136**, 2703–2706.
- 54 K. Trepte and S. Schwalbe, *J. Comput. Chem.*, 2021, **42**, 630–643.
- 55 H. Lin and P. A. Maggard, *Inorg. Chem.*, 2007, **46**, 1283–1290.
- 56 J. Y. Zhang, Y. Y. Xing, Q. W. Wang, N. Zhang, W. Deng and E. Q. Gao, *J. Solid State Chem.*, 2015, **232**, 19–25.
- 57 G. G. Sezer, M. Arıcı, I. Erucar, O. Z. Yeşilel, H. U. Özel, B. T. Gemici and H. Erer, *J. Solid State Chem.*, 2017, **255**, 89–96.
- 58 Q. Lin, C. Mao, A. Kong, X. Bu, X. Zhao and P. Feng, *J. Mater. Chem. A*, 2017, **5**, 21189–21195.
- 59 P. G. Boyd, A. Chidambaram, E. García-Díez, C. P. Ireland, T. D. Daff, R. Bounds, A. Gładysiak, P. Schouwink, S. M. Moosavi, M. M. Maroto-Valer, J. A. Reimer, J. A. Navarro, T. K. Woo, S. Garcia, K. C. Stylianou and B. Smit, *Nature*, 2019, **576**(7786), 253–256.
- 60 Y. Lan, X. Han, M. Tong, H. Huang, Q. Yang, D. Liu, X. Zhao and C. Zhong, *Nat. Commun.*, 2018, **9**(1), 5274.
- 61 R. Grau-Crespo, A. Aziz, A. W. Collins, R. Crespo-Otero, N. C. Hernández, L. M. Rodríguez-Albelo, A. R. Ruiz-Salvador, S. Calero and S. Hamad, *Angew. Chem., Int. Ed.*, 2016, **55**, 16012–16016.

