

Artificial-intelligence-driven discovery of catalyst “genes” with application to CO₂ activation on semiconductor oxides

Aliaksei Mazheika (✉ alex.mazheika@gmail.com)

The NOMAD Laboratory at the Fritz-Haber-Institut der Max-Planck-Gesellschaft <https://orcid.org/0000-0002-4705-1804>

Yanggang Wang

The NOMAD Laboratory at the Fritz-Haber-Institut der Max-Planck-Gesellschaft

Rosendo Valero

Universitat de Barcelona <https://orcid.org/0000-0002-4617-0721>

Francesc Vines

University of Barcelona <https://orcid.org/0000-0001-9987-8654>

Francesc Illas

Universitat de Barcelona <https://orcid.org/0000-0003-2104-6123>

Luca Ghiringhelli

The NOMAD Laboratory at the Fritz-Haber-Institut der Max-Planck-Gesellschaft

Sergey Levchenko

Skolkovo Institute of Science and Technology, Skolkovo Innovation Center

Matthias Scheffler

Fritz-Haber-Institut der <https://orcid.org/0000-0002-1280-9873>

Article

Keywords: artificial intelligence, CO₂, subgroup discovery

Posted Date: August 26th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-845882/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Communications on January 20th, 2022. See the published version at <https://doi.org/10.1038/s41467-022-28042-z>.

Artificial-intelligence-driven discovery of catalyst “genes” with application to CO₂ activation on semiconductor oxides

A. Mazheika^{1,*}, Y. Wang¹, R. Valero², F. Viñes,² F. Illas², L. M. Ghiringhelli^{1,3}, S. V. Levchenko^{4,*}, M. Scheffler^{1,3}

¹The NOMAD Laboratory at the Fritz-Haber-Institut der Max-Planck-Gesellschaft, 14195 Berlin-Dahlem, Germany

²Departament de Ciència de Materials i Química Física and Institut de Química Teòrica i Computacional (IQTCUB),

Universitat de Barcelona, c/ Martí i Franquès 1, Barcelona 08028, Spain

³The NOMAD Laboratory at the Humboldt University of Berlin, 12489 Berlin, Germany

⁴Skolkovo Institute of Science and Technology, Skolkovo Innovation Center, 3 Nobel Street, 143026 Moscow, Russia

*corresponding authors: alex.mazheika@gmail.com, mazheika@fhi-berlin.mpg.de; levchenko@fhi-berlin.mpg.de

keywords: CO₂ activation, first-principles calculations, subgroup discovery, heterogeneous catalysis, data-driven materials design

Abstract

Using subgroup discovery, an artificial intelligence (AI) approach that identifies statistically exceptional subgroups in a dataset, we develop a strategy for a rational design of catalytic materials. We identify “materials genes” (features of catalyst materials) that correlate with mechanisms that trigger, facilitate, or hinder the activation of carbon dioxide (CO₂) towards a chemical conversion. The approach is used to address the conversion of CO₂ to fuels and other useful chemicals. The AI model is trained on high-throughput first-principles data for a broad family of oxides. We demonstrate that bending of the gas-phase linear molecule, previously proposed as the indicator of activation, is insufficient to account for the good catalytic performance of experimentally characterized oxide surfaces. Instead, our AI approach identifies the common feature of these surfaces in the binding of a molecular O atom to a surface cation, which results in a strong elongation and therefore weakening of one molecular C-O bond. The same conclusion is obtained by using the bending indicator only when in

combination with the Sabatier principle. Based on these findings, we propose a set of new promising oxide-based catalyst materials for CO₂ conversion, and a recipe to find more. Our analysis also reveals advantages of local pattern discovery methods such as subgroup discovery over standard global regression approaches in discovering combinations of materials properties that result in a catalytic activation.

Introduction

The need for converting stable molecules such as carbon dioxide (CO₂), methane, or water into useful chemicals and fuels is growing quickly along with the depletion of fossil-fuel reserves and the pollution of the environment¹⁻³. Such a conversion does not have a satisfactory solution, so far.

The general understanding in heterogeneous catalysis is that a stable molecule needs to be “prepared”, before its catalytic conversion occurs. This leads to the notion of molecular activation⁴. However, on one hand, this notion encompasses a very wide variety of processes (adsorption, photo-excitation, application of electric field, etc.) and materials (including compositional and structural variability), and it remains unclear which properties of the catalytic material and the adsorbed molecule determine the final chemistry, what is the relationship between the two sets of properties, and how general this relationship may be. On the other hand, finding the set of descriptive parameters of a catalytic material that characterize the catalytic performance in a particular process, or even in general for a given reactant, would be very valuable, because it would allow us to quickly search for promising candidate catalysts using rational design⁵⁻¹². We call these properties “materials genes”. One way to find such properties for a given reaction is to explore the free-energy surface for each catalyst candidate, which is a slow and resource-consuming process, and currently computationally unfeasible for many materials on a high-throughput basis. An alternative approach consists in searching for a correlation between experimentally determined material’s properties and its catalytic performance. Such a strategy requires consistent experimental measurements at well-defined conditions for a set of materials. To the best of our knowledge, such consistent data have not been reported so far for CO₂ conversion on oxides. Moreover, available publications usually do not report unsuccessful experimental results. These issues and a strategy to address them have been recently discussed in our publication¹³.

Yet another strategy is to find an *indicator* of activation, namely, a property of the system that directly indicates certain catalytic performance of the material⁵. We distinguish here “indicator” from “materials genes” based on a qualitatively different level of computational complexity. The indicator can still be unfeasible or hard for a high-throughput study of hundreds of thousands or millions of materials. However, when it can be calculated for a few tens or hundreds of materials in a reasonable time, these data can then be used to find materials’ genes that control the value of *the indicator*. Since a direct search for a relationship between the indicator and catalytic performance of a material would also require a consistent set of data of turnover frequency (TOF), selectivity, and yield values, one could instead consider several most promising indicators, find out which materials are good catalysts, and then check which indicators correlate with this observation. This approach also addresses the problem of defining activation in terms of the adsorbed-molecule properties as potential indicators of

catalytic activity.

In this work, we focus on the CO₂ conversion as one of the most important societal and technological challenges^{1,2,14-18}, and oxide materials as candidate catalysts. Oxides are structurally and compositionally stable under realistic temperatures and can be less expensive than the traditional precious metal-containing catalysts¹⁹⁻²¹. An additional challenge is to ensure that the useful products as well as the surface catalytic activity are preserved under the conditions of activation and subsequent conversion. While the strong C-O double bonds in CO₂ can be weakened or even broken by adsorption at a solid surface at an elevated temperature, this may also lead to a too strong adsorption or further dissociation of the molecule, so that the catalytic surface is poisoned by carbonate or carbon deposits. A weak adsorption, on the other hand, means no activation. We take this into account by ensuring that the adsorption energy is not too large or too small (Sabatier principle). Our overall approach combines first-principles calculations with an artificial-intelligence (AI) method – subgroup discovery (SGD). The latter is used for identification of pristine materials properties that optimize indicators. Moreover, SGD allows identifying one or more distinct combinations of materials features (genes) that promote catalytic reactant activation. In order to ensure reproducibility of our AI data analysis, we provide all necessary metadata (input parameters) and workflow in the easily accessible form of a Jupyter notebook²². We argue that, with the ever-growing importance and complexity of AI, such detailed and tutorial documentation is a necessity of good scientific practice.

Our approach is applicable to a wider class of materials and molecules, not limited to oxides or CO₂. Our study by no means encompasses all possible mechanisms of CO₂ conversion on oxide surfaces, but it offers a clear design path among many possible ones.

Results

CO₂ activation. We find that on semiconductor oxide surfaces CO₂ is chemisorbed exclusively when the carbon atom binds to surface O-atoms. All other minima of the potential-energy surface are found to be either metastable or correspond to physisorption. Therefore, there are as many different potential chemisorption sites as there are unique O atoms at the surface. The data set includes all non-equivalent surface O atoms on the 141 considered surfaces of 71 materials, which sum up to 255 unique adsorption sites. Among these sites on about 4% (10 out of 255) CO₂ prefers to physisorb, *i.e.* any chemisorbed state is metastable with respect to the physisorbed one. The physisorption can be easily identified by an almost linear geometry of the adsorbed molecule, and a C-O bond distance very close to the C-O bond length in a gas-phase CO₂ molecule, 1.17 Å.

We considered six different candidate indicators of CO₂ activation, including OCO-angle and C-

O bond distance. Bending of the OCO angle in the adsorbed CO_2 molecule relative to the gas-phase value of 180° (linear configuration) has been previously proposed²⁴ and widely accepted as a good indicator of activation. For gas-phase CO_2 , it is understood that the C-O double bond is weakened when an electron is added to the lowest unoccupied orbital, because it is of antibonding (π^*) character with a concomitant bending of the molecule. There is a one-to-one mapping between the C-O bond length $l(\text{C-O})$ and the OCO angle in gas-phase $\text{CO}_2^{\delta-}$ for a range of $\delta > 0$ (red curve in Fig. 1). However, this is not the case for the adsorbed CO_2 (dots in Fig. 1). There is a subset of adsorbed CO_2 that is close to the red line, but there are many cases where $l(\text{C-O})$ is substantially larger for a given OCO angle. This is in contrast to metal alloy nanoparticle catalysts, where there is a better correlation between OCO angle and $l(\text{C-O})$ ²⁵. Also, a longer C-O bond reflects a weakening and readiness for further chemical transformations. Thus, the bond elongation itself may be an alternative indicator of activation. A look at the adsorbed CO_2 structures reveals that, on sites following the gas-phase correlation, the molecule adsorbs in nearly symmetric adsorption structures with nearly equal length of the two C-O bonds. In the other cases one O atom of CO_2 is close to surface cation(s), leading to a pronounced asymmetry of the adsorbed molecule.

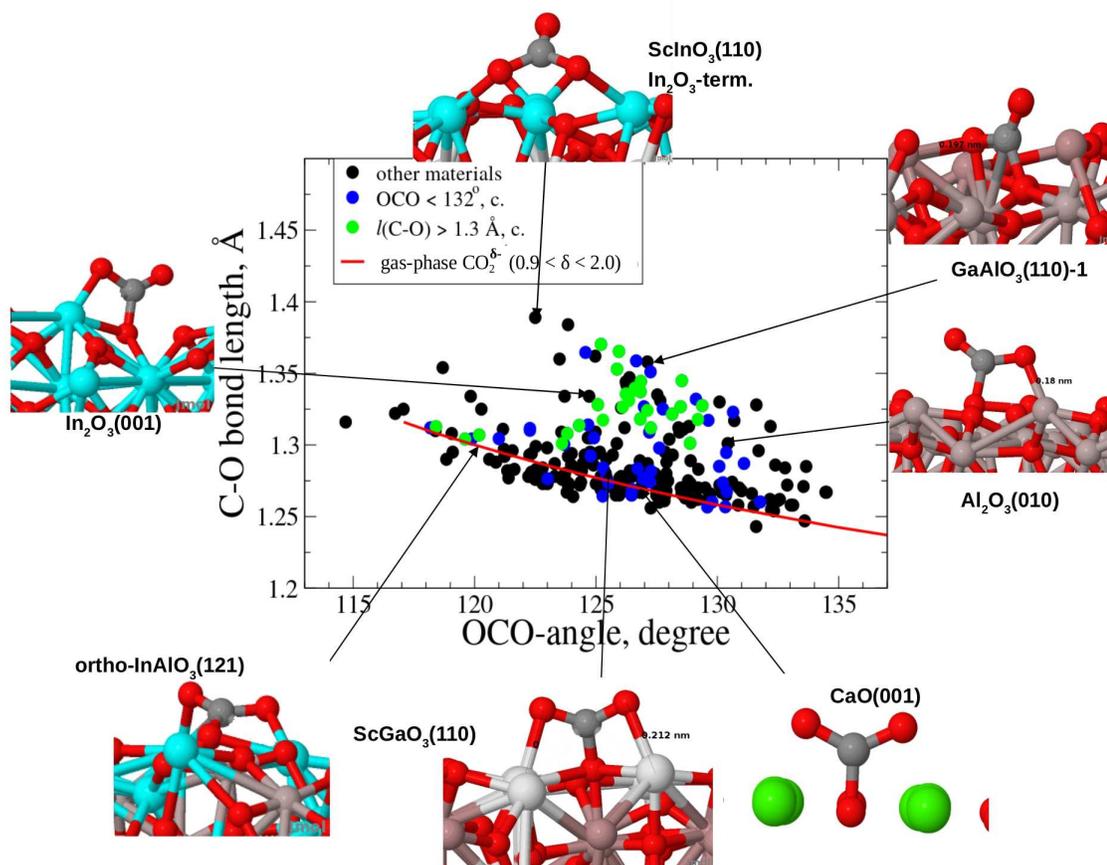


Figure 1. The correlation between the larger of the two C-O bond lengths (in case the two bond

lengths are different) and the OCO-angle in charged gas-phase (red line) and adsorbed CO₂ (dots). Colored dots: blue – adsorption sites from the unconstrained subgroup with OCO < 132°, green – subgroup of sites with $l(\text{C-O}) > 1.30 \text{ \AA}$, black – the remaining samples (see the text). The subgroups obtained with Sabatier principle constraint are marked with “c.”.

Other considered potential indicators of activation include Hirshfeld charge²⁸ of adsorbed CO₂ (a direct indicator of the charge transferred to CO₂), dipole moment of the surface along the surface normal per adsorbed CO₂ molecule (includes charge transfer to the molecule, as well as adsorption induced surface relaxation), the difference in Hirshfeld charges of C and O atoms in an adsorbed CO₂ molecule (indicates the ionicity of C-O bonds), and the difference in Hirshfeld charges of the O atoms in the adsorbed molecule (indicates asymmetry of the adsorbed molecule).^{28,4}

Subgroup discovery. To find out which properties (features) of the clean surfaces determine when a given activation indicator is maximized or minimized, we employ the subgroup-discovery (SGD) approach²⁹⁻³³. Given a dataset and a target property known for all data points, the SGD algorithm identifies subgroups with “outstanding characteristics” (see further for the criteria for being outstanding) and describes them by means of conjunction of basic propositions (*selectors*) of the kind “($f_1 < a$) AND ($f_2 \geq b$) AND ...”, where f_i is a feature and a, b are threshold values also found by SGD. In the framework of SGD, we call the selected primary features $\{f_1, f_2, \dots\}$ “materials genes”. Thus, SGD identifies both the outstanding subgroups and the relevant materials genes for a given target property.

Obviously, the selectors should only contain features that are much easier to evaluate than the target property. In presented work, the considered features include properties of gas-phase atoms that build the material, and properties of the pristine material (properties of the bulk phase and of the pristine relaxed surface). Overall 46 primary features have been considered. The full list is presented in SI. The features selected by the SGD are summarized in Table 1.

Table 1. Features that appear in the top SGD selectors (see text).

symbol	Meaning
$IP_{\min/\max}$	ionization potential, minimal and maximal in the pair of atoms A and B ; calculated as $E_{\text{atom}} - E_{\text{cation}}$
$EA_{\min/\max}$	electron affinity, minimal and maximal in the pair of atoms A and B ; calculated as $E_{\text{anion}} - E_{\text{atom}}$
$EN_{\min/\max}$	Mulliken electronegativity, minimal and maximal in the pair of gas-phase atoms A and B
$r_{-1}^{\min}, r_{-1}^{\max}$	radii of the maximum value of the Kohn-Sham radial wave functions of the spin-unpolarized

	spherically symmetric atom for HOMO-1, maximum (max) and minimum (min) in the pair of atoms <i>A</i> and <i>B</i>
$r_{+1}^{min}, r_{+1}^{max}$	radii of the maximum value of the Kohn-Sham radial wave functions of the spin-unpolarized spherically symmetric atom for LUMO, maximum (max) and minimum (min) in the pair of atoms <i>A</i> and <i>B</i>
<i>M</i>	energy at which the surface O 2 <i>p</i> -band projected density of states (PDOS) is maximal
d_1, d_2, d_3	distances from surface O-atom to the first-, second-, and third-nearest cations
<i>W</i>	work function <i>W</i> , as the negative of the valence band maximum ($W = -VBM$) with respect to vacuum level
q_{min}, q_{max}	minimal and maximal Hirshfeld charges of cations in the pair <i>A</i> and <i>B</i> , calculated as an average for all surface cations of a given type
Δ	band gap
<i>CBM</i>	conduction band minimum
Q_5, Q_6	local-order parameter with $l = 5$ or 6
<i>PC</i>	weighted surface O 2 <i>p</i> -band center
α_0, C_6^0	polarizability and C_6 -coefficient for surface O-atom obtained from many-body dispersion scheme
$\alpha_{min}, \alpha_{max}, C_6^{min}, C_6^{max}$	polarizability and C_6 -coefficient for cations, minimal and maximal in the pair <i>A</i> and <i>B</i> , calculated as an average for all surface cations of a given type
q_0	Hirshfeld charge of O-atom at the surface
<i>wid</i>	square-root of the second moment of surface O 2 <i>p</i> -band
wid_{min}, wid_{max}	square root of the second moment of PDOS of cations within valence band, minimal and maximal in the pair <i>A</i> and <i>B</i> , calculated as an average for all surface cations of a given type
c_{min}, c_{max}	first moment for PDOS of cation within valence-band, minimal and maximal in the pair <i>A</i> and <i>B</i> , calculated as an average for all surface cations of a given type
$\varphi_{1.4}, \varphi_{2.6}, \varphi_{1.4} - \varphi_{2.6}$	electrostatic potentials above surface O-atom at 1.4 and 2.6 Å and their difference. 1.4 Å corresponds to the average length of the bond between C and surface O, 2.6 Å is the minimal distance from surface O to C-atom of physisorbed carbon-dioxide molecule as observed from our calculations
L_{min}, L_{max}	energy of lowest unoccupied projected eigenstate of surface cations, minimal and maximal in the pair <i>A</i> and <i>B</i> , calculated as an average for all surface cations of a given type
<i>kurt</i>	kurtosis of surface O 2 <i>p</i> -band PDOS
<i>U</i>	eigenstate with least negative value in surface O 2 <i>p</i> -band
<i>BV</i>	bond-valence value of surface O-atom

The outstanding subgroup should satisfy several criteria. It should be statistically relevant; therefore the subgroups of too small size should be penalized. Target property values (OCO-angle, C-O bond length, etc.) for subgroup samples should be as different as possible from corresponding gas-

phase values since their change upon adsorption indicates CO₂ activation³². To achieve this, two requirements are imposed simultaneously: (i) The target-property values for subgroup members should be smaller or larger (depending on the target) than a certain value (a cutoff), and (ii) the target-property values are minimized or maximized within the cutoff. The latter condition gives preference to subgroups with smaller or larger target property values among similarly sized subgroups within the cutoff. The value of the cutoff is a parameter. As it approaches the optimal value of an activation indicator among all data points, additional or alternative materials genes and their combinations leading to stronger activation are identified. We explore the whole range of the parameter for each target property (for OCO-angle – 123°, 124°, 126°, 128°, 130°, and 132°; for l(C-O) – 1.26 Å, 1.28 Å, and 1.30 Å).

In addition to these criteria, we consider the requirement that adsorption energies are not too strong and not too weak for most of the samples in a subgroup. Strong activation (*i.e.*, strong weakening of the C-O bonds) can be achieved by strong binding to the surface. It is well known that good catalytic performance requires a balanced adsorption strength. This is known as Sabatier principle. In addition to the practical value of identifying subgroups that satisfy this principle, comparison of subgroup selectors obtained with and without this requirement helps to identify combinations of materials features that promote desired changes in target properties and at the same time yield intermediate adsorption energies.

Sabatier principle is reflected by a characteristic volcano-type behavior of catalytic activity as a function of adsorption energy of reactants and intermediates. Position of the top of volcano depends on particular reaction and conditions. It can be estimated from condition $|\Delta G| \sim 0$, where ΔG is the Gibbs free energy of adsorption. For CO₂ adsorption at room temperature and partial CO₂ pressure of 1 atm this condition corresponds to about -0.5 eV adsorption energy²⁶. At temperatures around 450 °C (typical conditions for CO₂ methanation²⁷) $\Delta G = 0$ corresponds to adsorption energy -1.7 eV²⁷. Therefore, for catalytic conversion at low or moderate temperatures this implies that CO₂ adsorption energies should be in the range from between -2.0 and -0.5 eV.

These requirements are implemented in the following quality functions that are maximized during the search for subgroups. In particular, for OCO-angle minimization we use:

$$F(Z) = \theta_{cut} \left[\frac{s(Z)}{s(Y)} \cdot \left(\frac{\max(Z) - \alpha_g}{\min(Y) - \alpha_g} \right) \cdot u(p) \right] \quad (1)$$

and for C-O bond maximization the following quality function was applied:

$$F(Z) = \theta_{\text{cut}} \left[\frac{s(Z)}{s(Y)} \cdot \left(\frac{\min(Z) - l_g}{\max(Y) - l_g} \right) \cdot u(p) \right] \quad (2)$$

where Y is the whole data set, Z – a subgroup, s – size (number of data points), \min and \max – minimal or maximal value of the target property, α_g and l_g are the gas-phase values of OCO-angle and C-O bond distance, 180° and 1.17 \AA respectively, and θ_{cut} is the Heaviside step function which is equal 1 if all data points in the subgroup satisfy the cutoff condition and 0 otherwise. Thus, larger values of the quality function $F(Z)$ are obtained for those subgroups in which minimal (maximal) value of a target property is close to the maximal (minimal) value of the whole sampling with respect to the gas-phase value of CO_2 molecule. The use of maximum/minimum instead of a median is done to ensure that a target property is optimal for as many members of a subgroup as possible. The gas-phase reference values are usually significantly different from the “chemisorption” subset. Therefore, the term in squared brackets in eq. (1) and (2) can noticeably contribute only when sizes of candidate subgroups are similar.

The term $u(p)$ in eq. (1) and (2) is added in order to account for Sabatier principle in SGD framework. We have implemented a multitask quality function, where a factor $u(p)$ increases quality of subgroups with adsorption energies falling within this range. This is formulated in terms of the information gain³³, *i.e.*, reduction of the normalized Shannon entropy. We perform the SGD for each target property both explicitly accounting for the Sabatier principle and without it. The latter case is equal to $u(p) = 1$ in eq. (1) and (2).³³

We note that SGD is qualitatively different from machine-learning classification/regression techniques such as neural networks, kernel regression methods, or decision-tree regression (DTR³⁴) (*e.g.*, random forest). SGD is typically referred to as a *supervised descriptive rule-induction* technique³⁵, *i.e.*, it uses the labels assigned to the data points (the values of the target property) in order to identify patterns in the data distribution (the statistically exceptional data groups) and the rules defining them (the selectors), by optimizing a quality function which is a functional of the distribution of values of the target property³⁵. While there are apparent similarities between SGD and DTR as both methods yield models in terms of physically interpretable selectors (usually, inequalities) on a selected subset of the input features, the analogy stops at this level, as SGD focuses at (and only at) subgroups from the very beginning and says nothing about the data that are not in the subgroup. In contrast, DTR determines a global partitioning of the input space by minimizing a global quality function, *i.e.*, the quality of a single subset is secondary with respect to the resulting quality of all subsets partitioning the whole data set. In other words, for finding distinct combinations of materials genes driving desirable

changes in a particular target property (possibly different combinations leading to the same result), the SGD approach has a significantly higher flexibility and reliability. This is demonstrated below for a DTR analysis for our target properties.

The metadata and workflow for the AI analysis is documented in the Jupyter notebook²².

Results of the subgroup discovery.

The SGD for OCO angles was done with eq. (1) for the quality function, and OCO as a target property, since smaller angles indicate larger charge transferred to the molecular π^* orbital. The subgroup selectors obtained with different OCO angle cutoffs (126°, 128°, 130°, and 132°) with or without the adsorption energy constraint are listed in Table 2 (for more details see the Table S4 in SI). Analysis of these subgroups reveals that the angle reduction is determined by an interplay of several factors: an electron transfer from the cations to surface O atoms, delocalization of electron density between cations and O atoms, and coordination of the surface O atoms. Without the Sabatier principle constraint, the OCO angle reduction below 132° is mainly due to the electron accumulation at the O atom of the clean surface. This is expressed by the conditions of more negative Hirshfeld charge on O-atoms ($q_o < \dots$), not very low IP of at least one cation ($IP_{\max} > \dots$), and increased polarizability of the surface O atom on which CO₂ is adsorbed ($C_6^O > \dots$). Upon adsorption of CO₂, this charge on the surface O-atom is readily available for transfer to CO₂. When the Sabatier principle constraint is introduced, the OCO < 132° subgroup also includes sites with a pronounced electron transfer to CO₂, but with a lower-energy O 2p band maximum ($M < \dots$) with respect to vacuum level, and a larger kurtosis ($kurt > \dots$). These conditions imply reduced inter-electronic repulsion around the surface O atom achieved by partial *delocalization* of the charge density.

At lower OCO cutoffs, the subgroup selectors include coordination descriptors Q_i , $i = 5, 6$. Without Sabatier principle, sites with larger Q_i are selected, and *vice versa*. Larger Q_i indicate lower coordination of the O atom. This reduces electron repulsion and therefore facilitates electron transfer to the O atom of the clean surface. However, this also increases the bonding strength of CO₂ to the surface. This explains why selectors of subgroups obtained with Sabatier principle include the opposite conditions ($Q_5 < \dots$).

Other surface features describing electron distribution are related to Madelung potential: electrostatic potential and field ($\varphi_{1.4}$, $\varphi_{2.6}$, and $\Delta\varphi = \varphi_{1.4} - \varphi_{2.6}$) and distances between the O atom and surface cations. More open surface structure with larger distances between cations at the O site facilitates charge transfer to adsorbed CO₂ molecule, since the Madelung potential from the nearby cations is reduced. This is reflected in appearance of propositions involving features d_1 , d_2 , and d_3 . For example, for the OCO $\leq 130^\circ$ subgroups, imposing energy constraint changes proposition ($d_1 > \dots$) to

($d_1 < \dots$), which implies an increased energy cost for transferring electrons to CO₂. Larger electric fields $\Delta\phi$ around the adsorption site imply stronger localization of electron density on O atoms, and thus also improve efficiency of charge transfer to the adsorbed molecule.

The smaller OCO subgroups with Sabatier principle also include propositions implying increased polarizability of both cations ($C_6^{\min} > \dots$). Another support-defining condition is that the radius of the lowest unoccupied orbital for the metal atoms should not be small ($r_{+1} \geq \dots$). This requirement is true for most cations with negative electron affinities (Figure S2). Analysis of adsorbed CO₂ structures and Hirshfeld charges reveals that this condition together with the higher polarizability of cations at the *pristine* surface encompass two scenarios: (i) additional electron transfer to CO₂ upon adsorption and (ii) stronger binding between O atoms in CO₂ and surface cations. When scenario (ii) dominates, CO₃^{δ-} anion lies nearly horizontally at the surface, and is bound with nearby cations by chemical bonds via its oxygen atoms. Such a structure leads to small OCO angles in CO₃^{δ-} (around 120°), even if charge transfer is limited. Thus, increased bending of adsorbed CO₂ occurs due to charge transfer over larger distances and/or distortion of the adsorbed molecule and the surface, both leading to a weaker adsorption. The cases where both scenarios are active include the same sites as in the subgroups with elongated $l(\text{C-O})$, as described below.

Table 2. Top subgroups and their selectors obtained by minimization of OCO-angle and maximization of $l(\text{C-O})$ with/out Sabatier principle (energies are in eV, distances are in Å, charges are in units of absolute electron charge, polarizabilities are in Bohr³). Proposition replacements that do not change the support are shown in parentheses.

cutoff	size	selector
OCO minimization without Sabatier principle constraint		
126	19	$L_{\max} > -2.70$ ($L_{\min} > -2.19$, $CBM > -3.40$, $r_{+1}^{\max} \leq 2.83$, $W < 5.80$, $U > -5.61$) $IP_{\max} \geq -6.05$ $\alpha_{\max} \leq 184.5$ $\Delta\phi > 1.33$ $q_{\max} \leq 0.59$ $wid \leq 1.59$ $wid \geq 0.58$
128	44	$EA_{\max} \geq -0.43$ $Q_6 \geq 0.51$ $\alpha_{\max} \geq 50.4$ ($C_6^{\max} \geq 389.5$, $\alpha_0 \leq 2.70$) $\Delta\phi \geq 1.00$ $q_{\min} \leq 0.49$

130	77	$L_{\max} \geq -5.23$ $EA_{\max} \leq 0.16$ ($C_6^{\max} \geq 389.5$, $IP_{\max} \geq -7.00$) $d_1 \geq 1.82$ $d_2 > 2.10$
132	139	$IP_{\max} \geq -6.99$ $q_0 \leq -0.32$ $C_6^0 \geq 10.36$
OCO minimization with Sabatier principle constraint		
126	15	$L_{\min} \geq -5.1085$ $\varphi_{2,6} \leq 0.3033$ $\Delta\varphi \leq 1.0622$ ($c_{\max} \leq -8.5915$) $d_1 \geq 1.82$ $d_2 \geq 2.005$ $r_{+1}^{\max} > 2.83$
128	30	$C_6^{\min} \geq 369.5$ $L_{\max} \geq -4.73$ ($r_{+1}^{\min} \leq 2.82$, $IP_{\min} \leq -5.83$, $r_{\text{HOMO}}^{\min} \leq 1.41$) $Q_5 \leq 0.83$ $\Delta\varphi \geq 0.60$ $r_{+1}^{\max} \geq 2.80$ $C_6^0 \leq 12.10$
130	40	$\varphi_{2,6} \geq -0.15$ $\Delta\varphi \geq 0.73$ $d_1 \leq 2.01$ $d_2 \geq 1.96$ $d_3 \geq 2.025$ ($c_{\min} \leq -9.07$, $W \geq 5.10$) $q_{\min} \leq 0.49$ $r_{+1}^{\min} \geq 1.94$
132	58	$q_0 \leq -0.3386$ $M \leq -6.292$ $kurt \geq 2.1035$ $IP_{\max} \geq -6.2085$ $r_{\text{HOMO}}^{\min} \leq 1.407$ ($IP_{\min} \leq -5.91$, $r_{+1}^{\min} \leq 2.82$)
l(C-O) maximization without Sabatier principle constraint		
1.26	121	$C_6^{\min} \geq 343.5$ $\varphi_{2,6} \leq 0.66$ $Q_5 \leq 0.83$ $M \geq -8.05$ ($PC \geq -9.32$)
1.28	38	$EA_{\max} \leq 0.005$ $d_2 > 2.22$ $M \leq -4.12$
1.30	27	$U \leq -5.34$

		$d_2 > 2.14$ $q_{\min} < 0.48$ $kurt \geq 2.10$ ($q_{\max} \geq 0.47$)
<i>l</i> (C-O) maximization with Sabatier principle constraint		
1.26	56	$CBM \geq -5.17$ ($L_{\min} \geq -5.11$) $\Delta\phi \leq 1.13$ $PC \geq -8.62$ $d_3 \leq 2.48$ $M \leq -6.06$
1.28	30	$W \geq 5.10$ ($M \leq -5.19$, $U \leq -4.92$, $PC \leq -7.21$) $d_2 > 2.14$ $q_{\min} < 0.48$
1.30	27	$EA_{\max} \leq 0.005$ ($W \geq 5.10$, $M \leq -5.19$, $U \leq -4.92$, $PC \leq -7.21$) $EN_{\min} \leq -3.19$ ($W \geq 5.10$, $q_O \geq -0.45$, $c_{\max} \leq -7.18$, $r_{\text{HOMO}}^{\min} \leq 1.41$, $\phi_{1.4} \leq 2.40$, $c_{\min} \leq -8.135$, $q_{\max} \geq 0.47$, $M \leq -5.19$, $IP_{\min} \leq -5.91$, $wid \geq 0.58$, $U \leq -4.92$, $r_{-1}^{\max} \geq 0.97$, $PC \leq -7.21$, $\Delta\phi \leq 1.81$) $d_2 > 2.14$ $q_{\min} < 0.48$ $kurt \geq 2.51$

In order to obtain the subgroups of adsorption sites with larger *l*(C-O), we performed the SGD with the quality function eq. (2) and *l*(C-O) as target property. The results for *l*(C-O) cutoffs 1.26, 1.28, and 1.30 Å are summarized in Table 2 and S5 in SI. In contrast to OCO, the analysis of the obtained top subgroups shows a much less pronounced or no effect of imposing Sabatier principle on the distribution of adsorption energies within the subgroups. This is because sites with too strong adsorption are excluded based on *l*(C-O) threshold alone, without the need to introduce the energy constraint. For example, the range of *l*(C-O) for the top *l*(C-O) > 1.26 Å subgroup without constraining adsorption energies is the same as for the top OCO < 130° subgroup, but it contains significantly more sites with intermediate adsorption energies.

Electron transfer to an adsorbed CO₂ molecule increases both the OCO bending and C-O bond elongation. The main difference between OCO and *l*(C-O) subgroups is that in the latter an additional mechanism of increasing *l*(C-O) is in effect, namely a covalent bonding between one O atom of the CO₂ molecule and the nearest surface cation. This can be concluded from the analysis of adsorption geometries, and correlates with the presence of proposition ($EA_{\max} \leq 0.005$ eV), selecting cation species that can accept electron density, e.g. from an O atom in adsorbed CO₂ molecule. Other proposition that appears in most selectors of top subgroups is ($d_2 > 2.14$ Å) or ($d_2 > 2.22$ Å) – larger distances to the second nearest cation from an O-atom. Larger elongation of the C-O bond is achieved by asymmetry of

the cation types at the surface, where one can bind an O atom of the adsorbed CO₂, while the other (located further away) cannot. An example asymmetric CO₂ adsorption structure is shown in Fig. S3 in SI.

Other propositions indicate a moderate charge transfer to adsorbed CO₂ molecule as in the case of OCO subgroups with adsorption energy constraint. Propositions ($M \geq -8.05$ eV), ($PC \geq -9.32$ eV) in $l(\text{C-O}) < 1.26$ Å subgroups imply enhanced charge density on the surface O-atoms, since electron-electron repulsion raises energies of O 2p band states. However, at larger $l(\text{C-O})$ cutoffs the electron transfer is balanced by such propositions as ($M \leq -5.19$ eV), ($U \leq -4.92$ eV), and ($W \geq 5.10$ eV) indicating limited electron transfer. These propositions point to a more covalent bonding between cations and surface O atom. Rather persistent proposition observed in many selectors of $l(\text{C-O})$ subgroups is the limit of minimal charge on surface cations ($q_{\min} < 0.48e$). It also shows the limitation of the charge transfer from one type of cations to surface oxygen atoms.

In general, we find that subgroups obtained with smaller cutoffs do not have a strong overlap with subgroups with larger cutoffs for OCO. In particular, for subgroups with close cutoffs the overlap can be smaller than 50% of the smaller subgroup (but is never below 30%). Interestingly, for $l(\text{C-O})$ the situation is opposite: subgroups with tighter cutoffs are mostly contained in the subgroups for more relaxed constraints. This means that, while larger values of $l(\text{C-O})$ are mainly controlled by the same or *additional* genes, smaller values of OCO are due to *alternative* genes. The overlap of OCO subgroups becomes even smaller when Sabatier principle is included, confirming the absence of a universal mechanism for OCO angle reduction that is compatible with moderate adsorption energy.

In summary, we find that, while an increased electron density at the O adsorption site is necessary for chemisorption and leads to both OCO bending and C-O bond elongation in an adsorbed CO₂ molecule, there are additional actuators for these effects that are different for different target properties. The OCO angle is in general minimized by increasing electron transfer to the O site. However, this also leads to a strong adsorption for many materials (Figure 2). To satisfy Sabatier principle, the electron transfer to CO₂ must be moderate. This is achieved by delocalization of charge density around O sites and/or by distortion of the adsorbed molecule due to formation of covalent bonds between O atoms in CO₂ and surface cations. The largest C-O bond elongations are achieved when both charge transfer to adsorbed CO₂ and the covalent interaction are present, and local geometry around surface O-atom provides the asymmetry in adsorption structure. This mechanism automatically fulfills the Sabatier principle.

The subgroups found by SGD for the dipole moment induced by CO₂ adsorption, its total Hirshfeld charge, and the difference of charges on C and O atoms significantly overlap with the

subgroup of smaller OCO-angles. The subgroup found by maximizing the difference of Hirshfeld charges on O-atoms of an adsorbed CO₂ largely overlaps with the subgroup of sites delivering larger $l(\text{C-O})$. In general, these indicators are not better than OCO or $l(\text{C-O})$. Therefore, below we focus on OCO angle and $l(\text{C-O})$ as indicators of CO₂ activation. More details about the other indicators can be found in SI.

Comparison with experimental results.

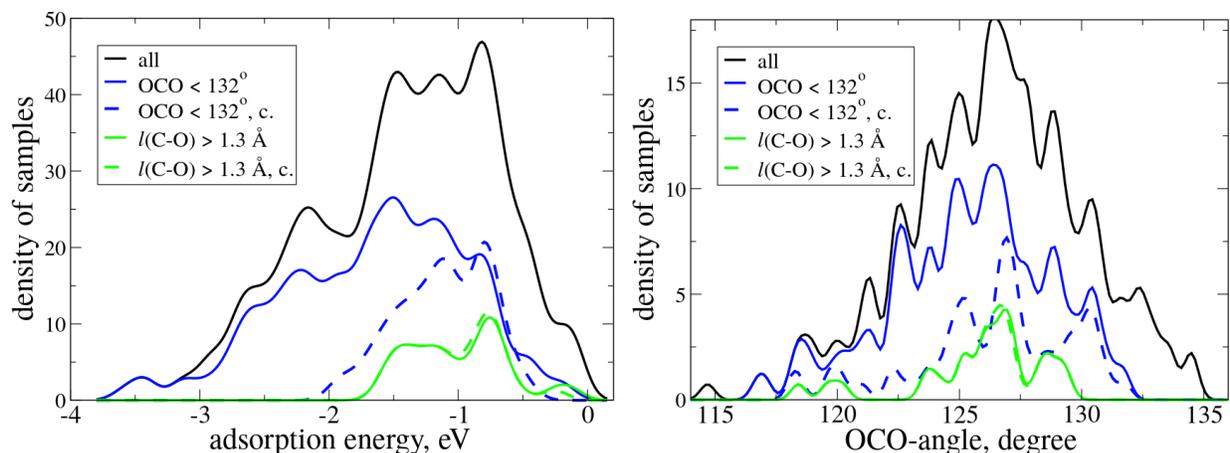


Figure 2. Distribution of adsorption energies (left) and OCO-angles (right) for the whole data set, the top subgroups of sites with OCO < 132° angles (blue) and $l(\text{C-O}) > 1.30$ Å (green). The subgroups obtained with adsorption energy constraint are marked with “c.” and shown with dashed lines. The adsorption energy E_{ads} is defined as the difference between the total energy of the slab with adsorbed CO₂ and the sum of total energies of the clean slab and an isolated CO₂ molecule.

To address the question which of the discussed properties can serve as indicator of the catalytic activity, we compare our predictions to reported experimental results (Table 3). It should be stressed that the available experimental data are scarce, and results are difficult to compare quantitatively. We consider thermally- and, for completeness, some photo-driven catalysis and thus also include supported metal catalysts with the considered oxides as support. Despite possibly different mechanisms for CO₂ conversion in the different types of catalysis, we believe that the properties of adsorbed CO₂ molecule can still serve as indicators of the catalytic activity. Thus, it is possible that under such daunting situation a reliable indicator of CO₂ activation can still be identified. As described below, our analysis confirms this hope.

Table 3. The catalytic performance of materials which contain the sites from larger $l(\text{C-O})$ or/and

smaller OCO subgroups.

material	catalytic reaction	CO ₂ adsorption energies, eV	belong to subgroups
NaNbO ₃	photocatalytic CO ₂ reduction with ~70% of CO selectivity ^{37,39}	-0.77 – -0.81	materials with sites from $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroup and $\text{OCO} < 132^\circ$ subgroup <i>with</i> Sabatier principle constraint
LaAlO ₃	dry reforming of methane with Ni-nanoparticles; performance is higher than for Ni-La ₂ O ₃ and Ni-Al ₂ O ₃ ⁴⁰	-1.17	
KNbO ₃	photocatalytic reduction of CO ₂ into CH ₄ as a composite with Pt/g-C ₃ N ₄ ; significant improvement of activity when compared to Pt/g-C ₃ N ₄ ; Pt-KNbO ₃ is ~2.5 times more photoactive than Pt-NaNbO ₃ ^{37,38}	-0.56 – -0.68	
CaTiO ₃	CO ₂ hydrogenation under UV-irradiation, although activity is not very high ^{43,46} ; twice higher activity with Ni nanoparticles ⁴⁶	up to -2.70	materials with sites from $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroups and from $\text{OCO} < 132^\circ$ subgroup <i>without</i> Sabatier principle constraint
CaZrO ₃ , SrZrO ₃ , BaZrO ₃ , SrTiO ₃	reverse water gas shift reaction (RWGS) under 700-1100°C ⁴²	up to -2.75	
SrTiO ₃	photocatalytic CO ₂ methanation with Pt, Au-nanoparticles, significant decrease of activity during reaction ⁴⁷	up to -2.40	
YInO ₃ *	no activity observed in photocatalytic CO ₂ conversion ³⁶	-1.16 – -1.47	materials with sites only from $\text{OCO} < 132^\circ$ subgroup <i>without</i> Sabatier principle constraint
CaO, SrO, BaO, Na ₂ O	strong carbonation, candidate materials for carbon capture and storage (CCS) ⁴¹	-1.60 – -3.57	
La ₂ O ₃	dry reforming of methane with supported Ni- nanoparticles; lower performance than on Ni-LaAlO ₃ ⁴⁰ and on some other supported catalysts ⁴⁴ at 700 and 250°C correspondingly	-2.14 – -3.11	
CaO	twice smaller reaction rate in CO ₂ reforming of methane reaction with supported Ni nanoparticles than on Ni-La ₂ O ₃ ⁴⁵ at 750°C	-1.60 – -3.42	
Ga ₂ O ₃	electrochemical reduction of CO ₂ to formic acid ⁴⁸ ; (photo)catalytic hydrogenation of CO ₂ ⁴⁹	-0.74 – -1.34	materials with sites from $\text{OCO} < 132^\circ$ subgroup <i>with</i> Sabatier principle constraint
Al ₂ O ₃	dry reforming of methane with supported Ni- nanoparticles ⁵⁰ ; lower performance than on Ni-LaAlO ₃ ⁴⁰	-0.87	

*material with sites also from $\text{OCO} < 132^\circ$ subgroup *with* Sabatier principle constraint

First we consider materials with the sites from subgroups obtained by minimization of OCO-angle without Sabatier principle constraint.²⁴ For quite many materials from these subgroups, independent on the cutoff value, there are no reports of successful CO_2 conversion, even when they are used as supports for metal nanoparticles (Table 3). This is explained by the fact that absolute adsorption energies for these materials are above 2 eV (Fig. 2 left, Table S4), indicating that their surfaces will be permanently poisoned by carbonate species at low or intermediate temperatures. This means that on materials with these sites hardly any reaction of CO_2 conversion can proceed at low, especially room temperature. Moreover, as shown in Table 3, even at increased temperatures, 700-750 °C, the activity of these materials is low. Some of them have been considered as candidates for carbon capture and storage (CaO, SrO, BaO, Na_2O)⁴¹, which implies formation of stable carbonates rather than CO_2 transformation. Thus, we conclude that OCO-angle alone is not a good indicator of enhanced catalytic activity in CO_2 conversion.

On the other hand, *several* of the materials with sites from $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroups (independent on either with or without Sabatier principle constraint) are known as good materials for CO_2 conversion (Table 3) in different reactions proceeding at room or higher temperatures. For these sites, the absolute adsorption energies already satisfy the Sabatier principle (Fig. 2, left), as discussed above. We note that, contrary to what one may expect, there is no correlation between the adsorption energy and the value of $l(\text{C-O})$ (see Fig. S3 in SI). Although there is a general trend, there are also significant variations in $l(\text{C-O})$ for a given adsorption energy.

Interestingly, some of the materials with sites in the $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroups were studied as supports for metallic nanoparticles. For instance, Ni/LaAlO₃ is a catalyst for dry reforming of methane⁴⁰ at 700°C. It was shown that its catalytic performance is higher in terms of CO_2 and CH_4 conversion rates compared to Ni/La₂O₃ and Ni/Al₂O₃⁴⁰. All sites on considered lanthanum (III) oxide surfaces belong to the subgroup of $\text{OCO} < 132^\circ$ without Sabatier constraint, whereas the sites on Al₂O₃ do not enter any of the two subgroups. KNbO₃ has been studied only with Pt nanoparticles and as a composite with g-C₃N₄ in photocatalytic reduction of CO_2 into CH_4 ^{37,38}. Pt-KNbO₃ is ~2.5 times more photoactive than Pt-NaNbO₃³⁷, whereas the NaNbO₃ is known to be photoactive even without nanoparticles³⁹. This seems to suggest that $l(\text{C-O})$ is a good indicator of CO_2 activation for both unsupported and supported catalysts even at increased temperatures. Hence, the other materials with the sites from this subgroup are promising new candidates for this task. The most promising materials identified in this work are CsNbO₃, CsVO₃, RbVO₃, LaScO₃, RbNbO₃, and NaSbO₃ as they have the

sites from the larger $l(\text{C-O})$ subgroups satisfying the above-mentioned criteria.

There is also a set of materials [ternaries $A^{2+}B^{4+}O_3$ ($A = \text{Ca, Sr, Ba}$, $B = \text{Zr, Ti, Ge, Sn, Si}$) with a perovskite structure] containing both the surfaces with sites from the smaller OCO subgroups without Sabatier constraint and the surfaces with sites from the larger $l(\text{C-O})$ subgroups (Table 3). These two types of sites are located on different surfaces. Thus, based on the above results, a material for which a surface with sites from the $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroups has lower formation energy and is more abundant than the surface with sites from smaller OCO subgroups without Sabatier constraint is expected to be a good catalyst. To explore this possibility, we analyze the surfaces of these materials in more detail. Their most stable surfaces are AO-terminated (001) facets containing sites from the smaller OCO subgroup. The formation energies of ABO_3 -terminated (110) surfaces with larger $l(\text{C-O})$ sites are higher: for BaZrO_3 , SrZrO_3 , CaZrO_3 , and SrTiO_3 the differences in formation energies are 0.049, 0.027, 0.013 and 0.037 eV/\AA^2 , respectively. The zirconates and SrTiO_3 were found to catalyze the water gas-shift reaction under increased temperatures, 700-1100 $^\circ\text{C}$ ⁴². At room temperature the photocatalytic activity of SrTiO_3 was found to be significantly decreased⁴⁷. We attribute the latter finding to the strong carbonation of its most stable surface, which is consistent with the calculated high absolute value of CO_2 adsorption energy (-2.4 eV) for this surface. Thus, the activity of SrTiO_3 at 700 $^\circ\text{C}$ and higher temperatures is consistent with the estimates of the CO_2 chemical potential given above. The difference in formation energies of the most stable CaO-terminated (001) surface and the stoichiometric (110) surface for CaTiO_3 is less pronounced compared to zirconates and other titanates (CaO-terminated (001) is more stable than the (110) surface by only 0.009 eV/\AA^2). Thus, the (110) facets, which contain sites from the long $l(\text{C-O})$ subgroup, may be present on catalyst particles at the reaction conditions. This can explain the observed activity of CaTiO_3 in CO_2 conversion not only at high but also at room temperature. We note that the activity of this material was also attributed to the presence of TiO_2 nanoparticles on the surface⁴³ at reaction conditions.

The OCO subgroup that includes most of the known good catalysts and a minimal number of inactive materials is $\text{OCO} < 132^\circ$ with Sabatier principle. It contains the sites on discussed above LaAlO_3 , KNbO_3 and NaNbO_3 catalysts, but also on non-active YInO_3 according to Ref. 36 (Table 3). This subgroup contains in addition the sites on a well-known CO_2 conversion catalyst Ga_2O_3 . We should mention that catalytic activity of Ga_2O_3 has been attributed to its reducibility. According to Pan and coworkers⁵¹ CO_2 molecules are activated via dissociation on surface O-vacancies. However, in ref. 52 only one Ga_2O_3 (100) surface was considered for which no energetically stable CO_2 chemisorption structures were obtained with the PBE functional. We show in Table S1 and Figure S1 that this functional underestimates CO_2 adsorption energies. Moreover, in our study we considered also other

surfaces and found stable CO₂ chemisorption structures on these surfaces. Thus, activation of CO₂ on Ga₂O₃ can indeed proceed on O-atoms as discussed in our study, even without surface O vacancies. The subgroups with small OCO cutoffs, 123° and 124°, do not contain any sites on known active or non-active catalysts.

OCO < 132° subgroup with Sabatier principle contains a large number of sites with elongated C-O bonds. The overlap of this subgroup with $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroups is 19 samples (70% of the latter).

To demonstrate advantages of SGD over DTR in finding materials genes and their optimal combinations, we have done comparison of found SGD subgroups with DTR performance for $l(\text{C-O})$. DTR terminal nodes (leaves) with largest average $l(\text{C-O})$ (see SI, Fig. S4 and S5) include surface sites on materials prone to extremely strong carbonation (Table 2), and also sites at which CO₂ prefers to physisorb, with $l(\text{C-O}) = 1.17 \text{ \AA}$. Also, one cannot check the effect of imposing the constraint as there is no standard way to mix regression and classification in DTR. Thus, DTR in contrast to SGD is not able to separate different activation modes and even fails sometimes in distinguishing activation from non-activation.

Best materials for CO₂ reduction among calculated ones.

Now that good indicators of activation are identified (OCO with Sabatier principle and $l(\text{C-O})$), all calculated materials can be ranked according to the value of these indicators (smaller OCO or larger $l(\text{C-O})$ indicate C-O bond weakening and therefore higher catalytic activity, provided adsorption energy is moderate). The resulting list of most promising catalysts for CO₂ conversion is presented in Table 4. Each surface is characterized by maximum $l(\text{C-O})$ and minimum OCO among all inequivalent sites on that surface. The materials with $l(\text{C-O}) > 1.30 \text{ \AA}$ are listed in the order of decreasing $l(\text{C-O})$. Materials with OCO < 132° but $l(\text{C-O}) < 1.30 \text{ \AA}$ are appended at the bottom of the list in the order of increasing OCO.

Table 4. Best materials and surface cuts for CO₂ activation according to the $l(\text{C-O})$ and OCO indicators.

material	surface cut	$l(\text{C-O}), \text{ \AA}$	OCO, degree	$E_{\text{ads}}, \text{ eV}$	in $l(\text{C-O}) > 1.30 \text{ \AA}$ subgroup	in OCO < 132° c. subgroup
NaSbO ₃	100	1.370	125.21	-1.32	yes	yes
Ga ₂ O ₃	212	1.365	124.57	-1.34		yes
NaSbO ₃	010	1.365	125.95	-1.09	yes	yes

LiSbO ₃	010	1.359	126.66	-1.04		yes
NaNbO ₃	100	1.353	125.87	-0.78	yes	yes
ScAlO ₃	010	1.351	127.25	-1.18		yes
KSbO ₃	110	1.345	128.54	-0.72	yes	yes
LiNbO ₃	100	1.344	126.23	-0.87		
NaNbO ₃	010	1.344	126.85	-0.77	yes	yes
InScO ₃	121	1.342	126.26	-1.23		
CsNbO ₃	100	1.34	126.6	-0.87	yes	
RbNbO ₃	111	1.338	126.61	-1.37	yes	yes
CsNbO ₃	010	1.336	126.23	-1.11	yes	
MgSnO ₃	100	1.334	119.84	-1.58		
GaAlO ₃	100	1.332	129.12	-1.02		yes
CaGeO ₃	001(GeO ₂ -term.)	1.331	127.65	-0.75		
InAlO ₃ -or.	121	1.33	130.09	-1.02		
ScAlO ₃	121	1.328	131.61	-0.86		
GaInO ₃	110	1.327	126.98	-1.34	yes	
LaAlO ₃	110	1.327	129.38	-1.17	yes	yes
CsVO ₃	110	1.327	126.1	-0.72	yes	
KNbO ₃	110	1.327	128.49	-0.68	yes	yes
RbVO ₃	110	1.326	126.04	-1.14		
Ga ₂ O ₃	110	1.325	127.76	-1.09		yes
NaVO ₃	110	1.324	127.12	-0.755	yes	
NaNbO ₃	110	1.322	128.14	-0.805	yes	yes
InAlO ₃ -rh.	110	1.318	126.83	-0.73	yes	yes
LaGaO ₃	100	1.317	125.29	-0.97	yes	
ScGaO ₃	010	1.314	124.68	-1.06		yes
GaInO ₃	120	1.313	118.41	-1.43	yes	yes
MgGeO ₃ - tetr.	001(GeO ₂ -term.)	1.312	126.18	-1.35		
ScAlO ₃	100	1.312	122.28	-1.89		yes
YAlO ₃	011	1.312	127.26	-1.18	yes	yes
InScO ₃	110	1.31	122.28	-1.54		yes
In ₂ O ₃	111	1.309	128.44	-0.65		
InAlO ₃ -or.	110	1.309	127.2	-0.66		yes
YAlO ₃	100	1.308	123.82	-1.305	yes	yes
InScO ₃	110(In ₂ O ₃ -term.)	1.305	124.92	-1.57		yes
YGaO ₃	100	1.305	124.76	-1.23		
In ₂ O ₃	110	1.301	125.86	-1.00		
Sc ₂ O ₃	111	1.301	130.43	-0.885		
LaGaO ₃	110	1.301	128.88	-0.83	yes	yes
LaScO ₃	100	1.301	123.6	-1.53	yes	
CaSiO ₃	001(CaO-term.)	1.290	118.84	-1.54		
SrSiO ₃	001(SrO-term.)	1.295	119.10	-1.66		
CaGeO ₃	001(CaO-term.)	1.288	120.88	-1.94		
Ga ₂ O ₃	212	1.297	121.21	-1.53		
InScO ₃	110	1.292	121.23	-1.88		

InScO ₃	100	1.277	121.40	-1.74		
RbVO ₃	100	1.283	121.64	-0.53		
In ₂ O ₃	110	1.280	122.52	-1.57		
InScO ₃	110(In ₂ O ₃ -term.)	1.284	122.80	-1.78		
SrGeO ₃	100(SrO-term.)	1.277	122.90	-1.70		
TiO ₂ -rutile	100	1.276	123.61	-1.05		
ZrO ₂	111	1.280	123.72	-0.92		
BaSnO ₃	001(BaO-term.)	1.267	123.80	-1.89		
ScGaO ₃	110	1.292	123.85	-1.22		
ZrO ₂	011	1.264	124.06	-0.72		
LiVO ₃	110	1.295	124.76	-0.70		
NaNbO ₃	010	1.273	125.00	-1.66		
MgTiO ₃	012	1.295	125.16	-1.47		
InAlO ₃ -or.	010	1.284	125.30	-0.82		yes
YInO ₃	100	1.293	125.69	-1.47		
KNbO ₃	010	1.277	125.97	-1.52		
InAlO ₃ -or.	110	1.278	126.04	-0.90		
ScAlO ₃	110	1.277	126.10	-1.33		
Al ₂ O ₃	012	1.265	126.46	-0.87		yes
Sc ₂ O ₃	110	1.265	126.47	-1.14		
CaSiO ₃	110(CaO-term.)	1.278	126.49	-1.44		
LaInO ₃	100	1.287	127.13	-1.27		
Sc ₂ O ₃	111	1.265	127.49	-0.95		
YInO ₃	110	1.298	127.61	-1.22		yes
ScAlO ₃	121	1.268	127.73	-0.755		
MgTiO ₃	001	1.265	127.85	-1.37		
BaGeO ₃	001(BaO-term.)	1.270	128.50	-1.80		
SrTiO ₃	001(TiO ₂ -term.)	1.266	128.53	-1.92		
ZnO	10-10	1.270	128.60	-1.005		
YGaO ₃	110	1.263	128.68	-1.60		
SrSnO ₃	001(SnO ₂ -term.)	1.273	128.90	-1.64		
Sc ₂ O ₃	001	1.289	128.90	-1.70		
MgGeO ₃	001	1.260	128.93	-1.09		
CaO	001	1.262	129.20	-1.60		
Al ₂ O ₃	001	1.283	129.22	-1.315		
BaSnO ₃	001(SnO ₂ -term.)	1.270	129.50	-1.87		
CaSnO ₃	001(SnO ₂ -term.)	1.272	130.09	-1.32		
KVO ₃	010	1.267	130.17	-0.55		
CaZrO ₃	101(ZrO ₂ -term.)	1.265	130.36	-1.86		
CaSnO ₃	110(SnO ₂ -term.)	1.272	130.50	-1.44		
SrGeO ₃	100(GeO ₂ -term.)	1.270	130.90	-1.515		
CaTiO ₃	101(TiO ₂ -term.)	1.266	131.42	-1.505		
SnO ₂	100	1.257	131.50	-0.85		
BaSiO ₃	100	1.243	131.60	-0.75		
MgO	111	1.296	131.70	-1.24		

As can be seen from this list, not all promising materials belong to one of the found subgroups.

This means that there are other optimal materials gene combinations that are not identified by SGD as statistically significant based on the current data set. Such combinations may be unique for a given material, or they may be found when more data for different materials are considered.

Discussion

We have developed the subgroup discovery strategy for finding improved oxide-based catalysts for the conversion of chemically inert molecules such as CO₂ into useful chemicals or fuels. For this purpose we identified a new indicator of CO₂ activation, namely the large C-O bond distance of the adsorbed molecule. This artificial-intelligence approach identifies the materials genes that correlate most strongly with the activation of the adsorbed molecule. Specifically these are the following clean surface properties: Hirshfeld charges of O atom at which CO₂ adsorbs (q_o) and of surface cations (q_{\min} , q_{\max}), surface geometric features [coordination descriptors Q_i , $i = 5, 6$, distances between the surface O atom and the nearest surface cations (d_i , $i = 1-3$)], electrostatic potential and electric field above the adsorption site ($\Delta\phi$, $\phi_{2,6}$), polarizability and C_6 coefficients for surface atoms (C_6^{\min} , C_6^O , α_{\max}), radii of HOMO and LUMO of the cation species (r_{+1}^{\max} , r_{+1}^{\min} , r_{HOMO}^{\min}), ionization potential, electron affinity, and electronegativity of surface cation species (IP_{\max} , EA_{\max} , EN_{\min}), features of O 2p DOS ($kurt$, M , PC , U), conduction band minimum (CBM), energies of the lowest unoccupied projected eigenstates of surface cation species (L_{\max} , L_{\min}), and surface work function (W). The found subgroup selectors predict whether a given candidate material belongs to the class of promising catalysts. The peculiarity of the large C-O bond indicator is that it automatically satisfies Sabatier principle for low and middle temperature CO₂ conversion.

The present study shows also that the previously proposed indicator for CO₂ activation, the decrease of the OCO angle²⁴, is not appropriate and even correlates with a strong adsorption, so that poisoning by carbonation is likely which may be useful for carbon capture and storage (CCS) but not for carbon capture and utilization (CCU). When Sabatier principle is purposely included in the SGD search for small OCO, found subgroups substantially overlap with large $l(\text{C-O})$ subgroups (70%), although still contain a few sites on inactive materials for CO₂ conversion.

The subgroup analysis revealed an alternative mechanism of CO₂ activation by adsorption, namely bonding of an O atom in CO₂ with a surface cation(s), combined with only moderate electron transfer from the surface to the molecule, which results not only in reduction of OCO-angles, but also in pronounced elongation and weakening of the C-O bond. Although the latter can be achieved also by a larger charge transfer, it results in stronger binding of CO₂ molecule to the surface and poisoning of the catalyst, contrary to the new mechanism. The same new mechanism is revealed when Sabatier

principle is included when searching for small OCO subgroups.

We also demonstrated that a standard regression technique (DTR), which gives prediction models in a physically interpretable form similar to subgroup discovery (selectors based on identified descriptor), fails to identify the optimal combinations of materials genes and the activation in general. This failure is traced back to the fact that DTR is a global approach, which minimizes error in prediction of the value of a target property for the whole data set. As a result, different combinations of genes leading to optimal value of the same target property are intermixed, and the combination that leads to the most optimal value is not identified. On the contrary, subgroup discovery finds unique local subsets in the data independent on the rest of the data. This makes it more suitable for identifying different combinations of materials genes that result in activation.

The other four considered potential indicators (charge at the adsorbed CO₂, adsorption induced dipole moment, difference of charges on O-atoms and on C and O atoms of adsorbed CO₂) were found to reproduce the results of SGD obtained for OCO-angles or C-O bond distances with significant overlap with corresponding subgroups.

Based on our results, we propose several new promising oxide-based catalysts for CO₂ conversion (Table 4). Although the present work has focused on oxides only, the overall strategy is general and can be applied to any other family of materials. This work also emphasizes the importance of documenting metadata and workflows for AI data analysis in materials science in order to ensure reproducibility of AI models and data analysis results.

Methods

Ab initio calculations. The calculations are performed using density-functional theory (DFT) with the PBEsol exchange-correlation functional²³ as implemented in FHI-aims code⁵² using ‘tight’ basis sets. The functional is chosen based on a comparison of calculated bulk lattice constants²³ and CO₂ adsorption energy to the available experimental results and high-level calculations (CCSD(T) and validated hybrid); see Supporting Information (SI) for more details on the computational setup. Nevertheless, it is expected that, because of the large set of systems inspected and the small variations introduced by the functional choice, the main trends will hold even when using another functional.

Studied materials. The data set includes 71 semiconductor oxide materials, with 141 surfaces. The materials are ternary (ABO₃) and binary oxides with metal cations *A* and *B* from groups 1 to 5 (including La) and groups 12 to 15 of the periodic table. The full list of materials and surface cuts is given in SI. In this study we considered only stoichiometric surface reconstructions obtained by atomic relaxation of stoichiometric bulk-like initial surface geometries. While this seems to be a limitation, our results show that indicators of activation calculated with this assumption correlate with experimental activity for known good oxide catalysts. This does not imply that surfaces of these materials do not reconstruct, but that the properties of unreconstructed surfaces can be used as descriptors for catalysis at reconstructed and defected surfaces under realistic conditions. Inclusion of surface reconstructions in the training data will further improve the predictions and will be a subject of future work.

The details of SGD. The SGD was done with the RealKD code (<https://bitbucket.org/realKD/>), modified to include quality functions described by eq. (1) and (2) in which the information gain was defined as:

$$u(p) = 1 - \left(\frac{-1}{\ln 2} \right) (p \cdot \ln(p) + (1-p) \cdot \ln(1-p)) \quad (3)$$

here p is the number of samples in a subgroup within required adsorption energy range divided by the total number of samples in the subgroup. Since Shannon entropy is a symmetric parabola-like function around 0.5, we set here $F(Z) = 0$ for $p \leq 0.5$. Also, $x \cdot \ln(x) = 0$ for $x = 0$. The search of subgroups is performed using a Monte Carlo scheme adapted for these tasks³³.

The cut-off values x, y, \dots used for setting propositions (feature-1 < x , feature-2 $\geq y$, etc.) are obtained by k -means clustering, as implemented within RealKD. That is, for a desired number $n = k - 1$ of cut-off values a set of k representative values of a given feature and k groups (clusters) of the data points are determined that minimize the deviation of all the feature values from the representative

values. Thus, each value of the feature in the data set is assigned to a particular cluster, and the cut-offs are determined as the arithmetic mean between the closest feature values in neighboring clusters. The number k is a parameter, and different k -values can in principle result in different cut-off values. It is worth noting that, due to the stochastic Monte-Carlo sampling, the exact definitions of the subgroups may vary for consecutive runs of the SGD algorithm. We have tested $k = 12, 14,$ and 16 and rerun the algorithm several times for each k . While the results indeed depend on the run and on the k value, the subgroups maximizing the quality function have largely or entirely overlapping populations, and selectors with the same or similar propositions. Here we report selectors that appear most often and have a high population and quality function values.

Decision-tree regression. The DTR analysis was performed using Python scikit-learn libraries. DTR is a supervised learning method in which the training set is repeatedly split into patterns (so called leaves) by means of propositions built from primary features. The fitting of a model is done with respect to the cost function, which encloses the deviation of fitted values of a target property from the actual values. In this study we considered two cost functions – mean squared error (MSE) and mean absolute error (MAE). The search of the most optimal partitioning (the so-called tree) is done with the greedy algorithm. To obtain the most optimal TR model, we used a standard approach for supervised machine learning – leave-one-out cross-validation with respect to the hyperparameters – minimal size of a leaf, maximal depth. Minimal size of a leaf is a bottom threshold of the population of a pattern, since too small size might result in overfitting. Maximal depth is a limit for the maximal number of splits in a tree.

Supporting Information

Ab initio methodology; studied materials and surface terminations; options applied in the program RealKD (Creedo); full list of used primary features; alternative subgroups found by SGD; density of samples with respect to the larger of the two C-O bond lengths; typical CO₂ adsorption structure from the subgroup with larger $l(\text{C-O})$; dependence of CO₂ adsorption energy on C-O bond length; SGD for alternative indicators; DTR results for all indicators

Acknowledgements

We thank Mario Boley for fruitful discussions on SGD and for providing the RealKD (for SGD) code. We also thank Yoshi Tateyama and Xinyi Lin for helping to generate the bulk oxide models and Helena Muñoz Galan and Oriol Lamiel Garcia for preliminary calculations. This project has received funding

from the European Union's Horizon 2020 research and innovation program (#951786: The NOMAD European Center of Excellence and the ERC grant #740233: TEC1p), and the Spanish MICIUN/FEDER RTI2018-095460-B-I00 and *María de Maeztu* MDM-2017-0767 grants and, in part, by *Generalitat de Catalunya* 2017SGR13 grant, plus a generous allocation of computational time provided by the *Red Española de Supercomputación* - RES (QCM-2017-3-0006, QCM-2017-2-0005, QCM-2016-3-0005, QCM-2016-2-0007).

References

1. Arakawa, H. et al. Catalysis research of relevance to carbon management: progress, challenges, and opportunities. *Chem. Rev.* **101**, 953–996 (2001).
2. Olah, G. A. Beyond oil and gas: the methanol economy. *Angew. Chem. Int. Ed.* **44**, 2636–2639 (2005).
3. Olah, G. A., Goeppert, A. & Surya Prakash, G. K. Chemical recycling of carbon dioxide to methanol and dimethyl ether: from greenhouse gas to renewable, environmentally carbon neutral fuels and synthetic hydrocarbons. *J. Org. Chem.* **74**, 487–498 (2009).
4. Somorjai, G. A. & Li, Y. Introduction to Surface Chemistry and Catalysis, 2nd Edition. John Wiley & Sons. 1-800 (2010).
5. Nørskov, J. K., Studt, F., Abild-Pedersen, F. & Bligaard, T. Fundamental Concepts in Heterogeneous Catalysis (2014).
6. Thornton, A. W., Winkler, D. A., Liu, M. S., Haranczyk, M. & Kennedy, D. F. Towards computational design of zeolite catalysts for CO₂ reduction. *RSC Adv.* **5**, 44361 (2015).
7. Duyar, M. S. et al. Discovery of a highly active molybdenum phosphide catalyst for methanol synthesis from CO and CO₂. *Ang. Chem. Int. Ed.* **57**, 15045–15050 (2018).
8. Peterson, A. A. & Nørskov, J. K. Activity Descriptors for CO₂ Electroreduction to Methane on Transition-Metal Catalysts. *J. Phys. Chem. Lett.* **3**, 251–258 (2012).
9. Liu, X. et al. Understanding trends in electrochemical carbon dioxide reduction rates. *Nat. Commun.* **8**, 15438 (2017).
10. Schlexer Lamoureux, P. et al. Machine learning for computational heterogeneous catalysis. *ChemCatChem.* **11**, 3581-3601 (2019).
11. Kitchin, J.P. Machine learning in catalysis. *Nat. Catal.* **4**, 230-232 (2018).
12. Medford, A.J., Kunz, M.R., Ewing, S.M., Borders, T., Fushimi, R. Extracting knowledge from data through catalysis informatics. *ACS Catal.* **8**, 7403-7429 (2018).
13. Foppa, L. et al. Materials genes of heterogeneous catalysis from clean experiments and artificial

intelligence. arXiv:2102.08269v.

14. Martens, J. A. et al. The Chemical Route to a Carbon Dioxide Neutral World. *ChemSusChem*. **10**, 1039-1055 (2017).
15. Klankermayer, J., Wesselbaum, S., Beydoun, K. & Leitner, W. Selective catalytic synthesis using the combination of carbon dioxide and hydrogen: catalytic chess at the interface of energy and chemistry. *Angew. Chem. Int. Ed.* **55**, 7296-7343 (2016).
16. Artz, J. et al. Sustainable conversion of carbon dioxide: an integrated review of catalysis and life cycle assessment. *Chem. Rev.* **118**, 434-504 (2018).
17. Li, W. et al. A short review of recent advances in CO₂ hydrogenation to hydrocarbons over heterogeneous catalysts. *RSC Adv.* **8**, 7651-7669 (2018).
18. Singh, A. K., Montoya, J. H., Gregoire, J. M. & Persson, K. A. Robust and synthesizable photocatalysts for CO₂ reduction: a data-driven materials discovery. *Nat. Commun.* **10**, 443 (2019).
19. Richter, N. A., Sicolo, S., Levchenko, S. V., Sauer, J. & Scheffler, M. Concentration of vacancies at metal-oxide surfaces: case study of MgO(100). *Phys. Rev. Lett.* **111**, 045502 (2013).
20. Arndt, S. et al. A critical assessment of Li/MgO-based catalysts for the oxidative coupling of methane. *Cat. Rev. Sci. Eng.* **53**, 424-514 (2011).
21. Yan, Z., Chinta, S., Mohamed, A. A., Fackler, J. P. & Goodman, D. W. The role of f-centers in catalysis by Au supported on MgO. *J. Am. Chem. Soc.* **127**, 1604-1605 (2005).
22. <https://nomad-lab.eu/services/AIToolkit> → Subgroup discovery for carbon-dioxide activation
23. Perdew, J. P. et al. Restoring the density-gradient expansion for exchange in solids and surfaces. *Phys. Rev. Lett.* **100**, 136406 (2008).
24. Freund, H.-J. & Roberts M. W. Surface chemistry of carbon dioxide. *Surf. Sci. Rep.* **25**, 225-273 (1996).
25. N. Austin, B. Butina, G. Mpourmpakis. *Progress in Natural Science: Materials International*. **2016**, 26, 487-492.
26. Stull, D. R. & Prophet, H. JANAF thermochemical tables. *J. Phys. Chem.* **78**, 2496-2506 (1974).
27. Wang, W. & Gong, J. Methanation of carbon dioxide: an overview. *Front. Chem. Sci. Eng.* **5**, 2–10 (2011).
28. Hirshfeld, F. L. Bonded-atom fragments for describing molecular charge densities. *Theor. Chim. Acta*, **44**, 129–138 (1977).
29. Wrobel, S. An algorithm for multi-relational discovery of subgroups. *European Symposium on Principles of Data Mining and Knowledge Discovery (Springer)*. 78–87 (1997).
30. Friedman, J. H. & Fisher, N. I. Bump hunting in high-dimensional data. *Statistics and Computing*.

9, 123-143 (1999).

31. Atzmueller, M.. Subgroup discovery, *Data Min. Knowl. Disc.* **5**, 35-49 (2015).

32. Boley, M., Goldsmith, B., Ghiringhelli, L. M. & Vreeken, J. Identifying consistent statements about numerical data with dispersion-corrected subgroup discovery. *Data Min. Knowl. Disc.* **31**, 1391–1418 (2017).

33. Goldsmith, B., Boley, M., Vreeken, J., Scheffler, M. & Ghiringhelli, L. M. Uncovering structure-property relationships of materials by subgroup discovery. *New J. Phys.* **19**, 013031 (2017).

34. Breiman, L., Friedman, J., Olshen, R. & Stone, C. Classification and regression trees. Wadsworth, New York (1984).

35. Novak, P. K., Lavrač, N. & Webb, G. I. Supervised descriptive rule discovery: A unifying survey of contrast set, emerging pattern and subgroup mining. *J. Mach. Learn. Res.* **10**, 377-403 (2009).

36. Khraisheh, M., Khazndar, A. & Al-Ghouti, M. A. Visible light-driven metal-oxide photocatalytic CO₂ conversion. *Int. J. Energy Res.* **39**, 1142-1152 (2015).

37. Shi, H. & Zou, Z. Photophysical and photocatalytic properties of ANbO₃ (A=Na, K) photocatalysts. *J. Phys. and Chem. Sol.* **73**, 788-792 (2012).

38. Shi, H., Zhang, C., Zhou, C. & Chen, G. Conversion of CO₂ into renewable fuel over Pt-g-C₃N₄/KNbO₃ composite photocatalyst. *RSC Adv.*, **5**, 93615-93622 (2015).

39. Fresno, F. et al. CO₂ reduction over NaNbO₃ and NaTaO₃ perovskite photocatalysts. *Photochem. Photobiol. Sci.*, **16**, 17-23 (2017).

40. Kathiraser, Y., Thitsartarn, W., Sutthiumporn, K. & Kawi, S. Inverse NiAl₂O₄ on LaAlO₃-Al₂O₃: unique catalytic structure for stable CO₂ reforming of methane. *J. Phys. Chem. C* **117**, 8120–8130 (2013).

41. Dunstan, M. T. et al. Large scale computational screening and experimental discovery of novel materials for high temperature CO₂ capture. *Energy Environ. Sci.* **9**, 1346-1360 (2016).

42. Saito. Y. *Patent No.:* US 8,540,898 B2; Sep. 24 (2013).

43. Sub Kwak, B. & Kang, M. Photocatalytic reduction of CO₂ with H₂O using perovskite Ca_xTi_yO₃. *Appl. Surf. Sci.* **337**, 138-144 (2015).

44. Muroyama, H. et al. Carbon dioxide methanation over Ni catalysts supported on various metal oxides. *J. Catal.* **343**, 178–184 (2016).

45. Zhang, Z., Verykios, X. E., MacDonald, S. M. & Affrossman, S. Comparative study of carbon dioxide reforming of methane to synthesis gas over Ni/La₂O₃ and conventional nickel-based catalysts. *J. Phys. Chem.* **100**, 744-754 (1996).

46. Lee, J. H. Cost-effective and dynamic carbon dioxide conversion into methane using a CaTiO₃@Ni-

- Pt catalyst in a photo-thermal hybrid system. *J. Photochem. Photobiol. A: Chem.* **364**, 219-232 (2018).
47. Zeng, S., Kar, P., Thakur, U. K. & Shankar, K. A review on photocatalytic CO₂ reduction using perovskite oxide nanomaterials. *Nanotechnology* **29**, 052001 (2018).
48. Sekimoto, T. Electrochemical application of Ga₂O₃ and related materials: CO₂-to-HCOOH conversion. *Jpn. J. Appl. Phys.* **55**, 1202 (2016).
49. Teramura, K., Tsuneoka, H., Shishido, T. & Tanaka, T. Effect of H₂ gas as a reductant on photoreduction of CO₂ over a Ga₂O₃ photocatalyst. *Chem. Phys. Lett.* **467**, 191-194 (2008).
50. Tang, S. et al. CO₂ Reforming of Methane to Synthesis Gas over Sol-Gel-made Ni/γ-Al₂O₃ Catalysts from Organometallic Precursors. *J. Catal.* **194**, 424-430 (2000).
51. Pan, Y.-X., Liu, C.-J., Mei, D. & Ge, Q. Effects of hydration and oxygen vacancy on CO₂ adsorption and activation on β-Ga₂O₃(100). *Langmuir* **26**, 5551 (2010).
52. Blum, V. et al. Ab initio molecular simulations with numeric atom-centered orbitals. *Comput. Phys. Comm.* **180**, 2175-2196 (2009).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SI.pdf](#)