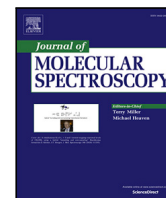




Contents lists available at ScienceDirect

## Journal of Molecular Spectroscopy

journal homepage: [www.elsevier.com/locate/jmmsp](http://www.elsevier.com/locate/jmmsp)

## Machine learning isotope shifts in molecular energy levels

Marco G. Barnfield<sup>a</sup>, Oleg L. Polyansky<sup>a,b</sup>, Sergei N. Yurchenko<sup>a</sup>, Jonathan Tennyson<sup>a</sup> <sup>\*</sup><sup>a</sup> Department of Physics and Astronomy, University College London, London, WC1E 6BT, UK<sup>b</sup> Institute of Applied Physics, Russian Academy of Sciences, Ulyanov Street 46, Nizhny Novgorod, 603950, Russia

## ARTICLE INFO

Dataset link: [https://github.com/mbarnfield63/ML\\_Isotopologue\\_Extrapolation.git](https://github.com/mbarnfield63/ML_Isotopologue_Extrapolation.git), [www.exomol.com](http://www.exomol.com)

## Keywords:

Carbon dioxide  
Carbon monoxide  
Isotopologues  
Machine learning

## ABSTRACT

Recent advances in the use of High-Resolution Cross-Correlation Spectroscopy (HRCCS) to detect molecular species in exoplanet atmospheres, presents a new challenge for the accuracy of reference spectroscopic line lists. While parent isotopologues of key atmospheric tracers are often well-characterized, minor isotopologues, crucial for diagnosing planetary formation histories and evolution, suffer from a scarcity of experimental data, often leading to reliance on less accurate theoretical predictions. In this work, a comprehensive machine learning framework is designed to mitigate these inaccuracies by modelling the residual errors of the isotopologue extrapolation (IE) method used within the ExoMol project. A fully connected neural network architecture for carbon dioxide (CO<sub>2</sub>) is shown to predict energy corrections with high fidelity, reducing the mean absolute error (MAE) relative to the original IE approach for more than 87% of the levels when benchmarked against empirical (MARVEL) energies. Furthermore, development of a novel hybrid, molecule-aware transfer learning architecture is presented that successfully propagates correction patterns from the data-rich CO<sub>2</sub> system to the data-poor carbon monoxide (CO) system. This transfer learning approach yields MAE improvements in over 93% of CO samples, demonstrating that physical correction factors related to isotopic substitution can be generalized across chemically related molecular systems. Updated and improved line lists are presented for 11 CO<sub>2</sub> isotopologues and energy levels for excited states of CO isotopologues are predicted. The methodology establishes a scalable, data-driven paradigm for refining molecular line lists, helping to bridge the gap between theoretical calculations and experimental precision.

## 1. Introduction

The field of exoplanetary science has undergone a rapid and transformative evolution over the past three decades. Since the detection of 51 Pegasi b in 1995 [1], which marked the first confirmation of a planet orbiting a main-sequence star other than our Sun, the catalogue of confirmed exoplanets has expanded rapidly. The number of confirmed exoplanets has now passed 6000, a milestone achieved through the tireless operations of space-based missions such as Kepler [2] and the Transiting Exoplanet Survey Satellite (TESS) [3], alongside ground-based radial velocity surveys [4]. This initial era of discovery, focused primarily on demographics and orbital dynamics, has now given way to a new phase: atmospheric characterization.

Beyond addressing long-standing questions about habitability and the potential for life, the chemical composition of an exoplanetary atmosphere also provides a fossil record of the planet's formation. It encodes information about where the planet formed within the protoplanetary disk, how it migrated to its current orbit, and the nature of its interaction with the host star [5]. For instance, the carbon-to-oxygen (C/O) ratio is widely used as a diagnostic of whether a planet

formed beyond the snow lines of key volatile species such as water, carbon dioxide, or carbon monoxide [6]. For this two complementary spectroscopic techniques are actively being pursued [7].

Transit Spectroscopy has been the primary tool for this chemical forensic work. By analysing the wavelength-dependent absorption or emission of light from the planet, astronomers can identify the unique spectral fingerprints of atmospheric molecules and retrieve temperature–pressure profiles, chemical abundances, and dynamical information such as winds and rotation rates [8]. The recent launch and successful deployment of the James Webb Space Telescope (JWST) has provided the community with high signal-to-noise, moderate- to high-resolution infrared spectra of exoplanet atmospheres from space, enabling multi-molecule detections and joint retrievals of composition and thermal structure [9].

In parallel, the next generation of  $\geq 30$  m ground-based facilities, such as the European Southern Observatory's Extremely Large Telescope (ELT) [10], is under construction and will host high-resolution spectrographs operating at the resolving powers  $R = \lambda/\Delta\lambda \approx 10^5$ . These instruments will exploit high-resolution cross-correlation spectroscopy

\* Corresponding author.

E-mail address: [j.tennyson@ucl.ac.uk](mailto:j.tennyson@ucl.ac.uk) (J. Tennyson).<https://doi.org/10.1016/j.jms.2026.112084>

Received 22 February 2026; Accepted 31 March 2026

Available online 7 April 2026

0022-2852/© 2026 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

(HRCCS), wherein by leveraging the orbital Doppler shifts of the planet relative to the quasi-stationary stellar and telluric spectra, HRCCS effectively separates the planetary component, thus delivering the highest effective spectral resolving power among current techniques for exoplanet atmosphere characterization [11]. This enables constraints on molecular abundances, atmospheric dynamics, and even 3D structure that are inaccessible at low resolution.

In the regime of high-resolution spectroscopy, the requirements for reference data precision become exceptionally stringent. While standard, computed line lists are often sufficient for lower-resolution transit studies, where broad absorption bands dominate the signal, HRCCS relies on matching thousands of individual spectral lines. At  $R \geq 10^5$ , the small shifts in energy levels caused by isotopic substitution become resolvable and, consequently, scientifically critical. Isotopologues are molecules that differ only in their isotopic composition, such as the substitution of a carbon-12 atom for carbon-13 in carbon monoxide ( $^{12}\text{C}^{16}\text{O}$  vs.  $^{13}\text{C}^{16}\text{O}$ ). Despite typically low abundances, these species provide invaluable diagnostics. The ratio of  $^{12}\text{C}$  to  $^{13}\text{C}$ , for example, is not affected by chemical processing in the same way as elemental ratios and can serve as a robust tracer of the primary isotope reservoir from which the planet accreted. Recent observational successes, such as the detection of  $^{13}\text{C}^{16}\text{O}$  in the atmosphere of the super-Jupiter TYC 8998-760-1 b [12] and the hot-Jupiter WASP-39 b [13] using cross-correlation methods, confirm that isotopologue detections are not merely theoretical in exoplanet atmospheres, but practically achievable with current technology.

However, the detectability of isotopologues is severely limited by the quality of the underlying spectroscopic data. Simulation studies have demonstrated that HRCCS is acutely sensitive to inaccuracies in line positions [14]. If the theoretical template used for cross-correlation is shifted even slightly from the true physical spectrum, the cross-correlation signal is degraded or lost entirely. Biases in the line lists can lead to erroneous retrievals of atmospheric abundances, temperatures, and velocities. Consequently, the promise of isotopic diagnostics in exoplanets depends directly on our ability to model the spectroscopy of minor isotopologues with the same fidelity as the parent species.

To address the spectroscopic needs of astrophysics, the ExoMol project was established with the mandate to produce comprehensive line lists for all molecules likely to be observable in hot atmospheres [15]. The ExoMol methodology combines sophisticated *ab initio* quantum mechanical calculations with laboratory data to generate line lists that are complete (often covering billions of transitions) and accurate.

A critical component in refining these theoretical models is the use of empirical energy levels derived from a large set of experimental data. The Measured Active Rotational–Vibrational Energy Levels (MARVEL) algorithm [16] is a powerful and widely used method for this purpose. MARVEL systematically processes all available assigned transitions for a given molecule, often collated from dozens of different experimental studies. These transitions are used to construct a spectroscopic network, where the quantum states are represented as nodes (or vertices) and the observed transitions between them are the links (or edges). This network-based approach is used to validate the input transitions, check for inconsistencies, and identify outliers or misassigned lines from the experimental data. Following this validation, a weighted least-squares fit is performed on the network to determine a self-consistent set of empirical energy levels with robustly determined uncertainties. These MARVEL-derived energy levels provide the high-quality benchmark data used to both refine the *ab initio* potential energy surface in the ExoMol procedure and to explicitly update individual energy levels, helping to ensure that the final line lists meet the stringent accuracy demands of modern astrophysics [17].

Despite the success of ExoMol, which currently provides comprehensive line lists for almost hundred parent molecules, a significant data gap exists for minor isotopologues. Experimental data for these

species is often sparse, with laboratory spectroscopy historically focused on the most abundant isotopologues due to signal-to-noise limitations. Without a dense network of experimental energy levels to refine the theoretical potential energy surfaces (PES), the calculations for minor isotopologues rely on extrapolations that inevitably introduce errors.

The “isotopologue extrapolation” (IE) method was introduced by Polyansky et al. [18] for water isotopologues, then formalized by McKemmish et al. [19] and applied to TiO, MgO and VO, leading to significant improvements in the accuracy as shown by its use to detect of TiO isotopologues in stellar spectra [20]. The method leverages experimental residuals ( $E_{\text{exp}} - E_{\text{calc}}$ ) of a parent isotopologue to correct calculated energy levels of minor isotopologues:

$$E_{\text{IE}}^{\text{iso}} = E_{\text{Ca}}^{\text{iso}} + \left( E_{\text{Ma}}^{\text{parent}} - E_{\text{Ca}}^{\text{parent}} \right), \quad (1)$$

where  $E_{\text{Ca}}^{\text{parent}}$  is the variationally calculated energy of the parent isotopologue,  $E_{\text{Ma}}^{\text{parent}}$  is the empirically derived MARVEL [16] energy,  $E_{\text{Ca}}^{\text{iso}}$  is the calculated variational energy of the given isotopologue, and  $E_{\text{IE}}^{\text{iso}}$  is the extrapolated isotopologue energy. In essence, the residuals in energy between the empirical and calculated line lists are assumed to be constant for all isotopologues [18,19]. While this is a reasonable first-order approximation, it fails to account for the subtle, mass-dependent effects of the breakdown of the Born–Oppenheimer approximation. These residual errors, while small, are large enough to hamper high-resolution studies.

The limitations of the constant-shift assumption in standard IE become apparent when considering the physical origins of the discrepancy between variational calculations and experimental data. These discrepancies are not random numerical artefacts; rather, they arise primarily from the breakdown of the Born–Oppenheimer (BO) approximation [21]. The BO approximation relies on the separation of electronic and nuclear motion based on their disparate masses. However, neglecting the coupling between these motions introduces specific error terms, adiabatic and non-adiabatic corrections, that scale inversely with nuclear mass and evolve with internal energy. Consequently, the residual encapsulates these neglected physical interactions. Because these effects follow deterministic physical laws dependent on quantum numbers and mass, the residuals are structured and predictable. This physical determinism is the crucial prerequisite that renders these residuals “learnable”, thus providing the theoretical justification for a machine learning approach.

Machine Learning (ML), specifically deep learning, has emerged as a powerful paradigm for modelling complex, non-linear physical phenomena [22,23]. The confluence of large-scale spectroscopic datasets and advanced computational techniques presents a unique opportunity. The ExoMol project provides an extensive database of transitions and energy levels that can be utilized by ML, with uses already in line broadening by Guest et al. [24] and automated quantum number assignment in the CO<sub>2</sub> “Dozen” line list [25].

In this paper a shift from purely physical extrapolation to a data-driven correction framework is proposed from the original IE method. By training neural networks to map quantum numbers and physical descriptors to these residuals, high-accuracy energy corrections can be predicted for unmeasured states in the minor isotopologues of carbon dioxide (CO<sub>2</sub>) and carbon monoxide (CO), two molecules of immense astrophysical importance. Furthermore, transfer learning is explored by using the knowledge learned from the data-rich CO<sub>2</sub> system to improve predictions for the data-poor CO system. This work aims to establish a robust, scalable, and transferable pipeline for generating high-fidelity spectroscopic data, thereby supporting the next generation of exoplanet observations.

**Table 1**

Summary of the empirical MARVEL data available for the minor isotopologues of CO<sub>2</sub>, detailing the maximum vibrational ( $\nu_i^{\max}$  in Herzberg notation) and rotational ( $J^{\max}$ ) quantum numbers, maximum energy ( $E^{\max}$  in cm<sup>-1</sup>), and total count of levels used.

Isotopologue	$\nu_1^{\max}$	$\nu_2^{\max}$	$\nu_3^{\max}$	$J^{\max}$	$E^{\max}$ (cm <sup>-1</sup> )	Count
<sup>16</sup> O <sup>12</sup> C <sup>17</sup> O (627)	4	5	5	109	13 469.81	4639
<sup>16</sup> O <sup>12</sup> C <sup>18</sup> O (628)	6	5	5	111	13 465.76	6008
<sup>13</sup> C <sup>16</sup> O <sub>2</sub> (636)	5	18	9	113	19 985.94	6918
<sup>16</sup> O <sup>13</sup> C <sup>17</sup> O (637)	4	4	4	99	9 178.13	2714
<sup>16</sup> O <sup>13</sup> C <sup>18</sup> O (638)	4	5	5	105	9 758.55	3881
<sup>12</sup> C <sup>17</sup> O <sub>2</sub> (727)	4	4	5	89	13 169.70	1851
<sup>17</sup> O <sup>12</sup> C <sup>18</sup> O (728)	4	3	4	90	9 105.74	2799
<sup>13</sup> C <sup>17</sup> O <sub>2</sub> (737)	2	5	3	56	7 175.70	541
<sup>17</sup> O <sup>13</sup> C <sup>18</sup> O (738)	3	2	3	83	8 594.01	963
<sup>12</sup> C <sup>18</sup> O <sub>2</sub> (828)	6	0	3	100	13 514.70	4663
<sup>13</sup> C <sup>18</sup> O <sub>2</sub> (838)	4	6	3	91	9 467.74	1703

## 2. Method

Although in principle machine learning (ML) can be employed to directly predict the energy levels of minor isotopologues from those of the parent isotopologue, in practice this approach yields larger residuals than the original IE method. Instead, this work presents a proof of concept that residuals of the original IE method against empirical (MARVEL) levels can be used in a ML framework. We show that the residual at each energy level ( $\Delta E_i$ ) is learned to produce corrections to energy levels across all isotopologues as follows

$$E_{\text{ML}}^{\text{iso}} = E_{\text{IE}}^{\text{iso}} + \Delta E_i^{\text{iso}}. \quad (2)$$

The residual itself is defined by

$$\Delta E_i^{\text{iso}} = E_{\text{Ma}}^{\text{iso}} - E_{\text{IE}}^{\text{iso}}. \quad (3)$$

### 2.1. Data curation & preparation

The CO<sub>2</sub> data was derived from the ‘‘Dozen’’ line list [25], a comprehensive update to the ExoMol database for carbon dioxide. The isotopologues’ energy levels were extracted in a series of MARVEL studies which considered the parent species <sup>12</sup>C<sup>16</sup>O<sub>2</sub> as well all the 11 stable minor isotopologues [26–33], such as <sup>13</sup>C<sup>16</sup>O<sub>2</sub>, <sup>16</sup>O<sup>12</sup>C<sup>18</sup>O, and <sup>13</sup>C<sup>18</sup>O<sub>2</sub> which we denote 636, 628, 638, respectively, below. For the parent <sup>12</sup>C<sup>16</sup>O and the five stable CO isotopologues (<sup>12</sup>C<sup>17</sup>O, <sup>12</sup>C<sup>18</sup>O, <sup>13</sup>C<sup>16</sup>O, <sup>13</sup>C<sup>17</sup>O, <sup>13</sup>C<sup>18</sup>O, henceforth denoted 27, 28, 36, 27 and 38), updated energy levels derived from recent MARVEL analyses were adopted [34,35]. For both CO<sub>2</sub> and CO, the MARVEL datasets contain a significant number of validated, empirical energy levels: 8268 levels spanning energies up to 20654 cm<sup>-1</sup> with  $J \leq 118$  for <sup>12</sup>C<sup>16</sup>O<sub>2</sub> and 2293 levels spanning energies up to 67 148 cm<sup>-1</sup> with  $\nu \leq 41$  and  $J \leq 123$  for <sup>12</sup>C<sup>16</sup>O; specific quantum number coverage and energy limits for the individual minor isotopologues are detailed in Tables 1 and 2. The number of available MARVEL-derived energy levels for the minor CO isotopologues is significantly lower than for CO<sub>2</sub> (Tables 1 and 2); this data scarcity necessitates the transfer learning approach described in Section 2.4. Only energy levels possessing values for all four necessary components ( $E_{\text{Ma}}^{\text{parent}}$ ,  $E_{\text{Ca}}^{\text{parent}}$ ,  $E_{\text{Ca}}^{\text{iso}}$ ,  $E_{\text{Ma}}^{\text{iso}}$ ) were retained in our ML procedure. This intersection is necessary to calculate both the input features (the original IE prediction) and the target label (the true residual to the MARVEL-derived energy level) for supervised learning.

### 2.2. Feature engineering & importance

To enable the neural network to learn the physics of the residuals, a rich ML feature set was constructed representing the quantum state and isotopic properties of each energy level. The features (Table 3) were divided into continuous variables (standardized) and categorical

**Table 2**

Summary of the empirical MARVEL data available for the minor isotopologues of CO, detailing the maximum vibrational ( $\nu^{\max}$ ) and rotational ( $J^{\max}$ ) quantum numbers, maximum energy ( $E^{\max}$  in cm<sup>-1</sup>), and total count of levels used.

Isotopologue	$\nu^{\max}$	$J^{\max}$	$E^{\max}$ (cm <sup>-1</sup> )	Count
<sup>12</sup> C <sup>17</sup> O (27)	4	21	8 689.70	33
<sup>12</sup> C <sup>18</sup> O (28)	22	38	40 467.91	498
<sup>13</sup> C <sup>16</sup> O (36)	27	45	48 166.08	737
<sup>13</sup> C <sup>17</sup> O (37)	1	22	2 964.58	45
<sup>13</sup> C <sup>18</sup> O (38)	13	30	25 080.73	345

**Table 3**

Feature set used for the CO<sub>2</sub> Neural Network.

Feature name	Description
<i>Energies</i>	
$E_{\text{Ca}}^{\text{iso}}$	Calculated energy (minor isotopologue)
$E_{\text{Ca}}^{\text{parent}}$	Calculated energy (parent)
$E_{\text{Ma}}^{\text{parent}}$	MARVEL energy (parent)
$E_{\text{IE}}^{\text{iso}}$	Original IE calculated energy (minor isotopologue)
<i>Quantum numbers</i>	
$J$	Total rotational quantum number
$g_{\text{tot}}$	Total state degeneracy
$\nu_1, \nu_2, \nu_3$	Normal-mode vibrational quantum numbers (Herzberg notation)
$m_1, m_2, l_2, m_3, r$	AFGL vibrational quantum numbers
$l_1, l_2, l_3$	TROVE local-mode quantum numbers
<i>Reduced masses</i>	
$\mu_1, \mu_2, \mu_3, \mu_{\text{all}}$	Reduced masses (sym., asym., and bend) and respective ratios to the parent isotopologue. See Appendix A for definitions.
<i>Boolean flags</i>	
Atomic masses	e.g.: Presence of <sup>13</sup> C atom; Presence of <sup>16</sup> O in position 1
e/f	Parity labels
$A'/A''$	TROVE symmetry labels

variables (encoded). Notably, states with  $J + \nu_3$  odd, denoted B' and B'' symmetries by TROVE, appear in the 727 and 737 isotopologues; these levels are excluded from the training set as there are no corresponding experimental levels within the parent isotopologue 626 due to nuclear spin statistics.

Feature importance was subsequently analysed via test-set permutation scoring to aid interpretability regarding which spectroscopic inputs most strongly influenced the learned CO<sub>2</sub> corrections. In this approach, each input feature is randomly shuffled in the held-out data and the resulting degradation in predictive accuracy is measured. Consequently, features whose perturbation most harms performance are identified as the most influential for the model’s corrected CO<sub>2</sub> energy predictions.

### 2.3. Network architecture I: CO<sub>2</sub>

For the data-rich CO<sub>2</sub> dataset, a feed forward neural network was implemented using the PyTorch [36] library. A full representation of the structure is shown in Fig. 1, containing six hidden layers, each with a decreasing number of units from 1024 to 32, with a single final output layer to predict the correction.

The activation function used was the Gaussian Error Linear Unit (GELU) [37]. Unlike the standard Rectified Linear Unit (ReLU) [38] which has a sharp discontinuity at zero, GELU, is a smooth, probabilistic activation function. This smoothness is advantageous for regression tasks in physics, where the target function (in this case the energy correction) is continuous and differentiable [37].

Dropout layers were interleaved between the dense layers which randomly zero out a fraction of neurons during training, thus preventing the network from relying too heavily on any single feature or memorizing noise in the training data. This promotes generalization to the unseen energy levels.

The optimizer used was Adam (adaptive moment estimation) [39] with a learning rate of  $1 \times 10^{-3}$  and a weight decay of  $1 \times 10^{-6}$ . The

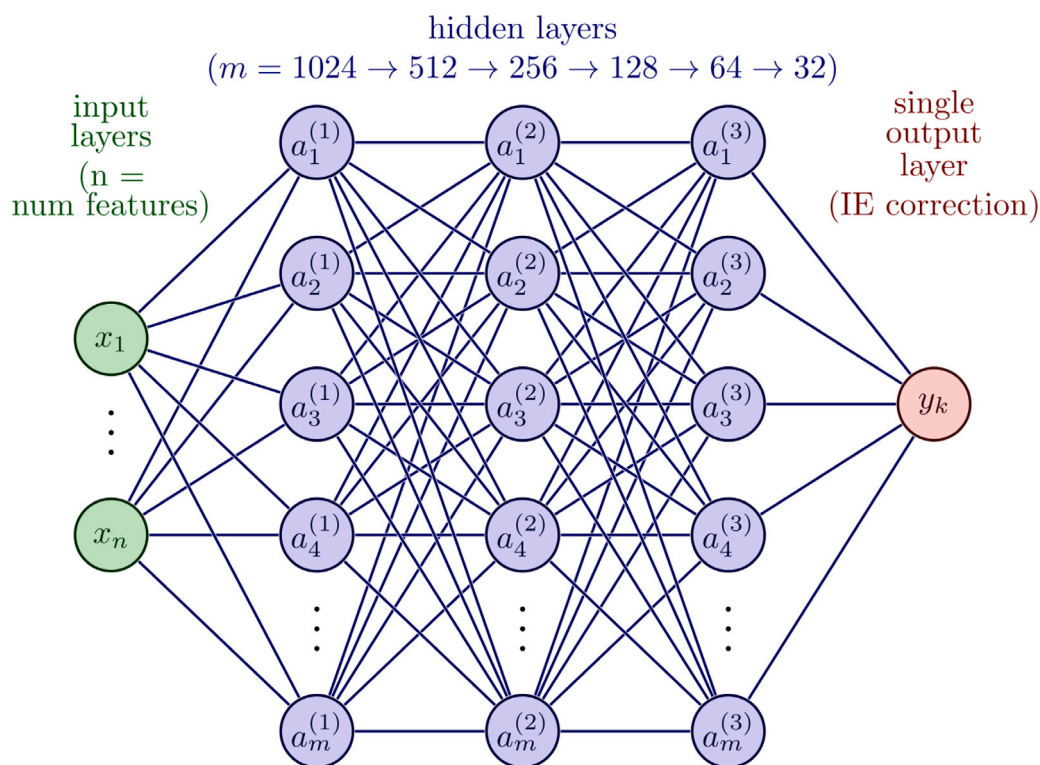


Fig. 1. Neural network structure for CO<sub>2</sub> IE corrections.

loss function was a Huber Loss [40], which behaves similarly to mean squared error (MSE) for small errors and mean absolute error (MAE) for large errors. Additionally it is robust to outliers, essential given that experimental spectroscopic data can contain misassignments which would skew a pure MSE loss.

The training protocol involved a standard 70/10/20 split between training, validation and test sets, ensuring all isotopologues were represented in each split. Final evaluation used only the untouched test split, reporting loss, root mean square error (RMSE) and MAE. After inference, the model's predicted correction  $\Delta E^{\text{iso}}$  was combined with the baseline calculated energy to compute corrected energies (Eq. (2)) and their improvement relative to the uncorrected values ( $E_{\text{IE}}^{\text{iso}}$ ).

#### 2.4. Network architecture II: Transfer learning for CO<sub>2</sub>

A hybrid, molecule-aware architecture is used for CO that allows the network to learn generalized correction patterns for the data-rich CO<sub>2</sub> dataset while utilizing individual “adapter heads” for each isotopologue to capture the specific spectroscopic nuances of CO. This transfer-learning approach ensures that the model benefits from the large scale statistics of CO<sub>2</sub> without losing specificity for CO. A full representation of the structure is shown in Fig. 2.

The shared trunk now consists of just two dense layers (256 → 128) with LayerNorm [41] to normalize the inputs across the features, GELU activations and dropout; a learned gating mechanism blends the shared and isotope-specific outputs to give the final scalar correction  $\Delta E^{\text{ML}}$  at each isotopologue specific head. This design preserves the ability to transfer broad chemical patterns across data, while still allowing per-isotopologue and per-molecule adjustments when required. Training used Adam optimization with a lower learning rate of  $5 \times 10^{-4}$  with  $1 \times 10^{-6}$  weight decay and Huber loss as before. For stability, the training gradients were clipped to a maximum norm of 1.0 and the model output bias was initialized to the mean target of the training subset to speed convergence.

A random weighted sampler was employed to ensure that CO samples were drawn more frequently during training, counteracting CO<sub>2</sub>

dominating learning and predictions, thus preventing the network from ignoring the minority class. The training strategy also contained adaptive weighting, i.e., if a specific isotopologue's performance lagged, its weight in the loss function was dynamically increased, forcing the network to focus on the “harder” cases. Given the small sample size, a stratified 5-fold cross-validation (CV) was also implemented to ensure the performance in all available CO data points was adequate. To further improve reliability and reduce dependence on random initialization, the complete cross-validation process was repeated across five independent random seeds. For each seed, improvement in prediction accuracy relative to the original IE residuals was quantified using per-isotopologue MAE. The final reported improvement is given as a mean  $\pm$  standard deviation across seeds, providing a clear measure of both the overall correction achieved and the variability of the results across runs.

### 3. Results & discussion

#### 3.1. Global performance summary

To provide a unified quantitative comparison across both molecular systems, the mean absolute errors (MAEs) for the raw variational calculations from the “Dozen” line list (Ca) [25], the original isotopologue extrapolation method (IE), and the machine-learning corrected IE energies (ML) are summarized in Table 4. The metrics show that, while the original IE method yields only marginal improvement over the raw calculated values, the machine-learning models deliver a substantial reduction in MAE for both CO<sub>2</sub> and CO. This demonstrates that the residual-learning approach captures structured, isotopologue-dependent deviations that neither variational calculations nor constant-shift extrapolation adequately model.

The improvement is especially notable for CO, where the ML correction reduces the MAE by approximately an order of magnitude compared with the original IE method. For CO<sub>2</sub>, the numerical improvements are smaller because the extensive MARVEL dataset already

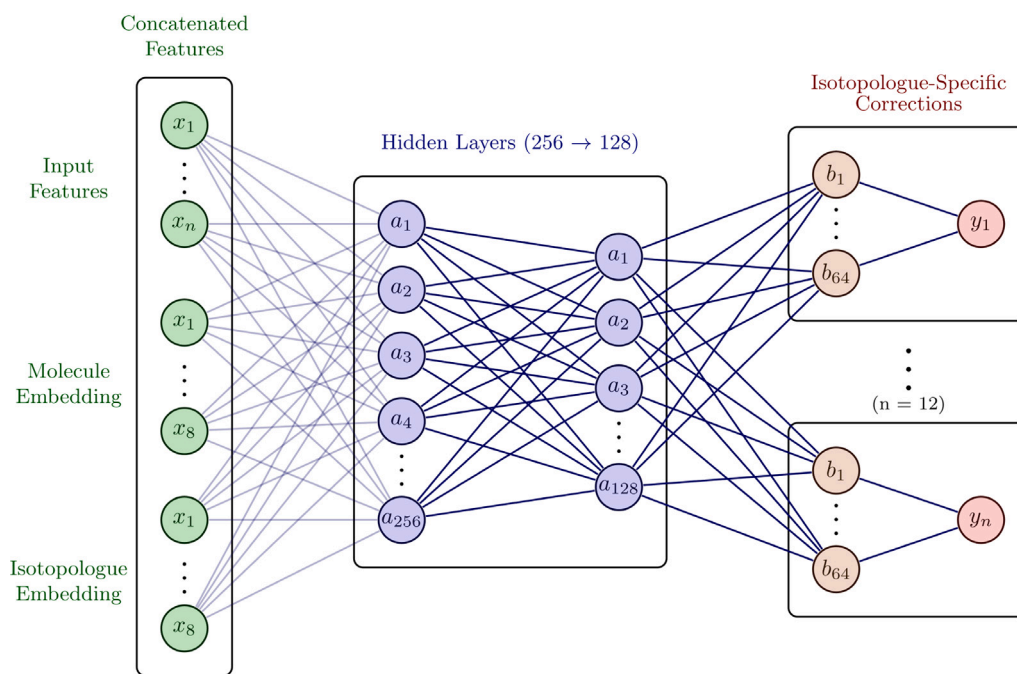


Fig. 2. Neural network structure for CO IE corrections.

Table 4

Mean absolute errors (MAE) in  $\text{cm}^{-1}$  for variational calculations (Ca), original isotopologue extrapolation (IE), and ML-corrected-IE (ML) predictions for  $\text{CO}_2$  and CO energy levels when compared with their respective empirical ( $\text{MARVEL}$ ) levels (previously presented in Tables 1 and 2).

Molecule	Ca	IE	ML
$\text{CO}_2$	0.01395	0.01394	0.00232
CO	0.03007	0.02896	0.00524

makes the IE predictions highly accurate; however, the ML corrections still deliver a clear and consistent reduction in error beyond what IE alone achieves. This cross-molecule behaviour strengthens the case for employing machine-learning residual correction as a general framework for isotopologue refinement, irrespective of data abundance.

### 3.2. $\text{CO}_2$

While initial training included all eleven minor  $\text{CO}_2$  isotopologues, the  $^{13}\text{C}^{16}\text{O}_2$  (636) isotopologue was subsequently excluded from the training set. This decision is empirically justified by analyzing the baseline performance of the original IE method. As shown in Fig. 3, the original IE method yields an MAE for 636 of  $<0.002 \text{ cm}^{-1}$ , nearly an order of magnitude lower than the problematic species in 628 or 838.

It is worth noting that when included, the neural network did successfully learn the residual physics for this species, achieving a further 33.12% reduction in error as illustrated in the residual distribution in Fig. 4. However, this comparatively marginal gain for an already accurate species came at a significant cost. Because 636 represents a large proportion of the available data ( $\approx 18\%$ ), including these low-error targets dominated the loss function. This hampered the network's ability to predict corrections for the remaining isotopologues, resulting in a lower global accuracy. Consequently, 636 was removed to allow the model to focus on the species where the standard approximation fails significantly.

When 636 was included, the overall MAE improvement across isotopologues was 85.93%, with 83.14% of individual samples showing improvement. After removing 636, the overall MAE improvement

increased to 89.27%, and the proportion of samples showing improvement rose to 91.62% over the original IE method.

The overall model performance following this adjustment is shown in Fig. 5, which compares the mean absolute error (MAE) and root mean squared error (RMSE) across the remaining  $\text{CO}_2$  isotopologues before and after neural-network correction. Both metrics decrease substantially following correction, confirming that the model successfully compensates for systematic discrepancies present in the original IE calculations. The reductions are consistent across isotopologues, suggesting that the network captures generalizable correction patterns rather than overfitting to any specific isotopic composition. The greater reduction observed in RMSE compared with MAE indicates that the model is particularly effective at correcting large outliers, which tend to dominate the squared-error measure.

The distribution of residuals across all  $\text{CO}_2$  isotopologues before and after correction in the unseen test set is shown in Fig. 6. Prior to correction, residuals are widely spread and biased to the right, indicating that the original IE method systematically overestimates corrections to the calculated values. This behaviour can also be observed in the individual isotopologue residual plots shown in Fig. 7, with the minor exception of  $^{13}\text{C}^{17}\text{O}_2$  (737), which is routinely underestimated. After correction, the residuals form a much narrower and more symmetric distribution centred around zero, and distinct patterns are no longer discernible within the individual isotopologue subplots. This demonstrates the network's ability to remove systematic bias while substantially reducing errors across the full energy-level range.

#### 3.2.1. Feature analysis

Feature importance was examined using the ablation approach described in Section 2.2, where individual features were removed and the corresponding change in MAE was recorded. Table 5 shows that all importance values were positive, indicating that removing any single feature reduced the overall predictive accuracy. The most substantial contributions to minimizing the mean absolute error (MAE) were derived from the inclusion of the isotopic masses and key spectroscopic predictors,  $J$  and vibrational quantum numbers from TROVE, Herzberg, and AFGL. Standard normal-mode vibrational quantum numbers are  $\nu_1, \nu_2, l_2, \nu_3$  in Herzberg notation; AFGL ( $m_1, m_2, l_2, m_3, r$ ) refers to the

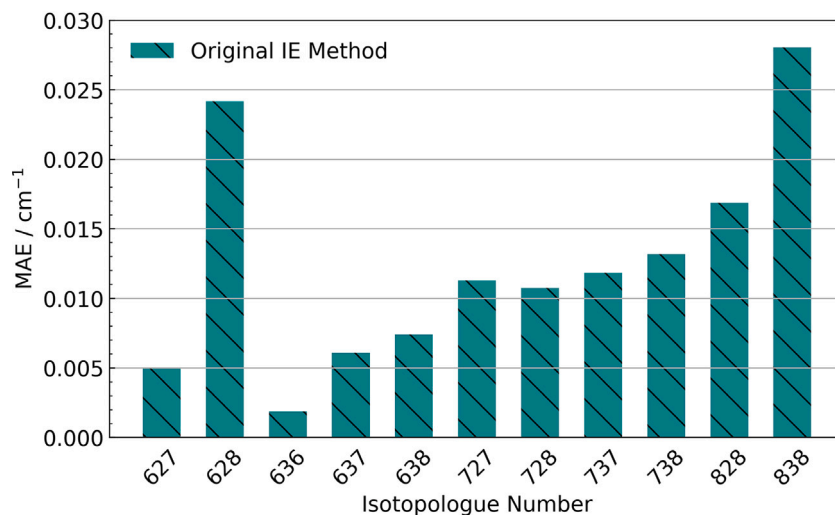


Fig. 3. The mean absolute error (MAE) of the original Isotopologue Extrapolation (IE) method for each CO<sub>2</sub> minor isotopologue.

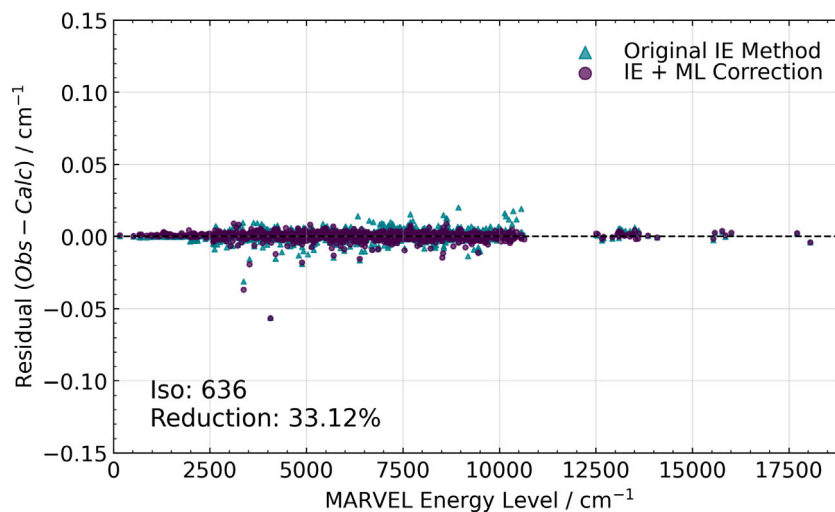


Fig. 4. <sup>13</sup>C<sup>16</sup>O residuals before and after ML correction.

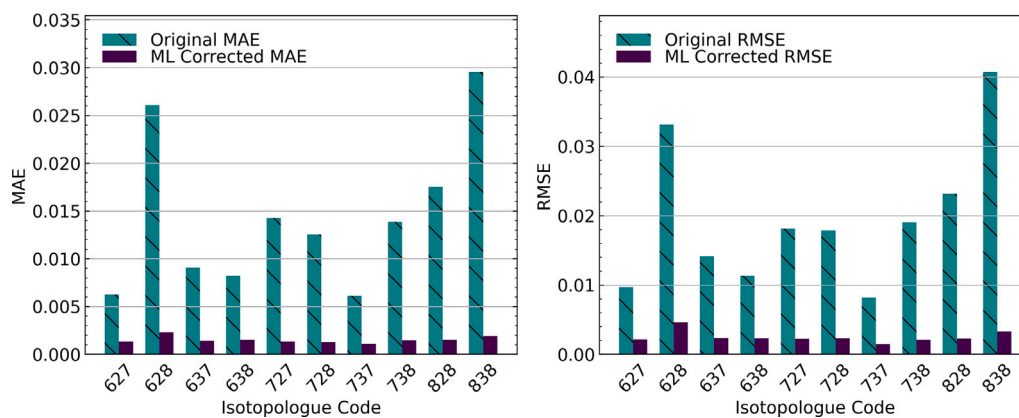
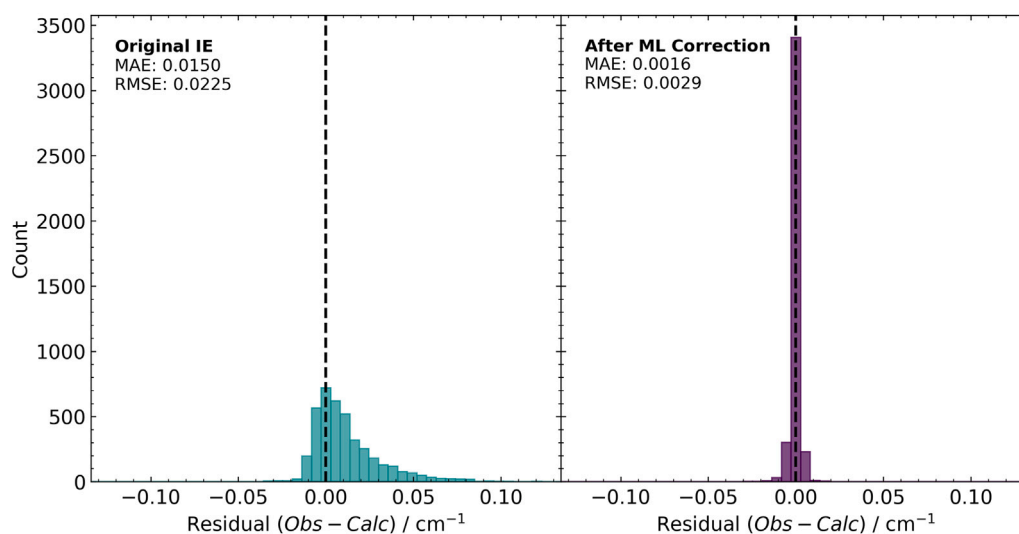
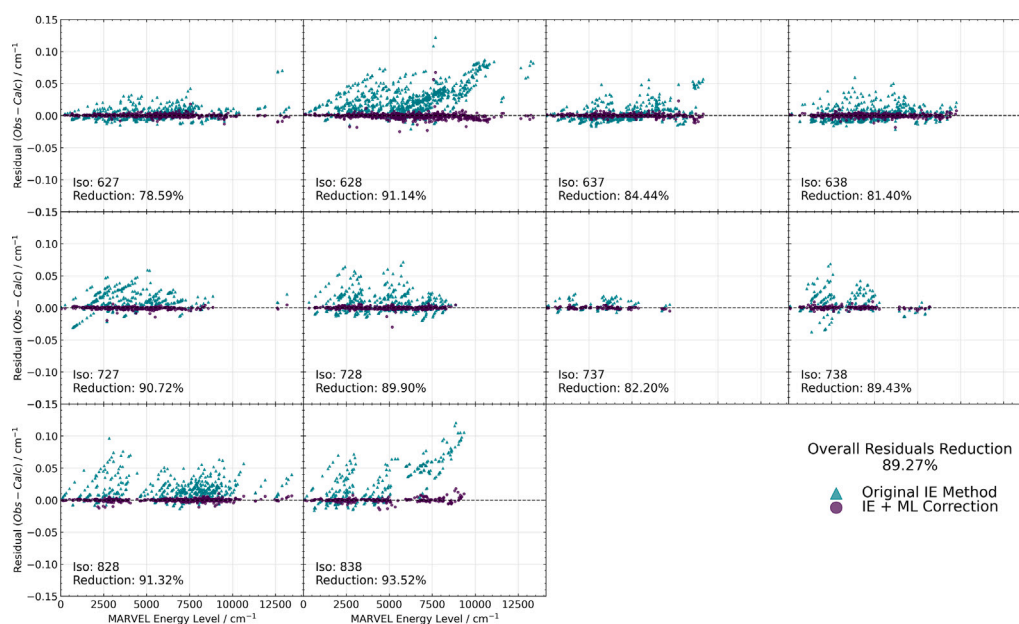


Fig. 5. Mean absolute error (MAE) and root mean square error (RMSE) across minor CO<sub>2</sub> isotopologues before and after the ML correction.



**Fig. 6.** Distribution of residuals for all  $\text{CO}_2$  isotopologues, representing the discrepancy between empirical MARVEL energy levels and the IE-calculated energies before and after ML correction.



**Fig. 7.** Residuals for individual  $\text{CO}_2$  isotopologues plotted against MARVEL empirical energy levels before and after the ML correction.

Air Force Geophysics Laboratory notation standard [42]; TROVE quantum numbers  $(t_1, t_2, t_3)$  correspond to local-mode assignments from the variational calculations.

The strong importance of explicit isotopic masses and rovibrational quantum numbers, beyond what is captured by simple reduced-mass scaling, is consistent with the known isotopic dependence of non-adiabatic (Born–Oppenheimer breakdown) corrections in high-precision molecular spectra. Taken together, these results indicate that the network relies on physically meaningful quantities and that the learned corrections encode genuine isotopic dependencies rather than fortuitous correlations.

### 3.3. Carbon monoxide (CO)

Much like  $\text{CO}_2$ , the CO network demonstrates a clear improvement in predictive accuracy following neural-network correction, as shown in Fig. 8, yielding an average MAE improvement of 87.82% over the original IE method. In total, 91.37% of samples showed improved agreement

with experiment after the ML corrections were applied. An example of these improvements is presented in Table 6. These results confirm that the hybrid molecule-aware architecture effectively generalizes across isotopologues while retaining molecule-specific correction capability.

The residual distributions in Fig. 9 mirror those observed for  $\text{CO}_2$ , narrowing significantly and centring around zero after correction. The network successfully removes the negative bias of the original IE calculations in CO, indicating that the inclusion of  $\text{CO}_2$  data helped stabilize learning and extend the correction trends to CO. The inset metrics summarize the overall error reduction and confirm the improvement seen across isotopologues.

Fig. 10 presents the residuals for each CO isotopologue plotted against the empirical MARVEL energy levels. As with  $\text{CO}_2$ , pre-correction residuals display systematic structure, particularly at higher energies, which largely disappears after correction. The uniform reduction in error across isotopologues demonstrates the success of transferring learned correction behaviour from the  $\text{CO}_2$  dataset to CO.

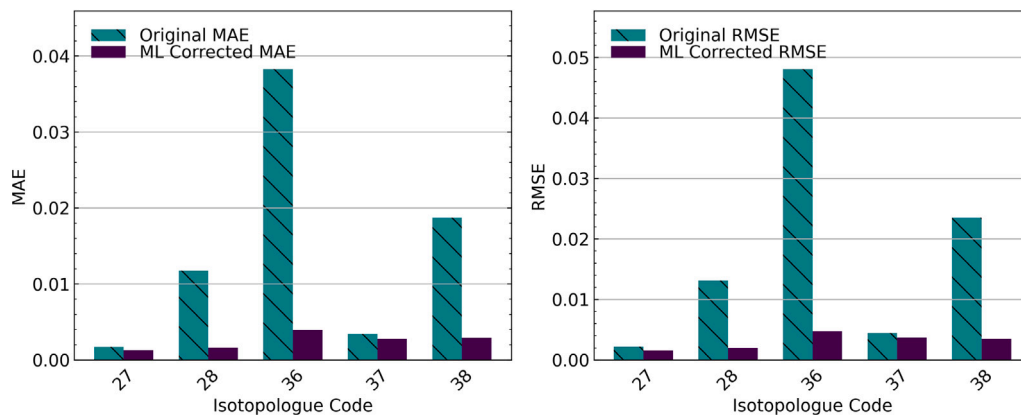


Fig. 8. Mean absolute error (MAE) and root mean square error (RMSE) across minor CO isotopologues before and after the ML correction.

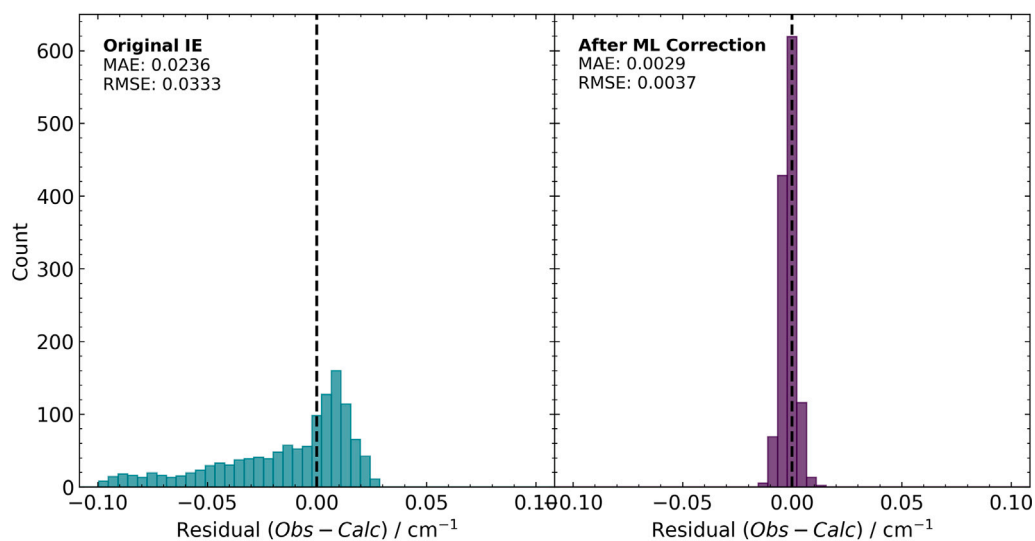


Fig. 9. Distribution of residuals for CO isotopologues before and after the ML correction.

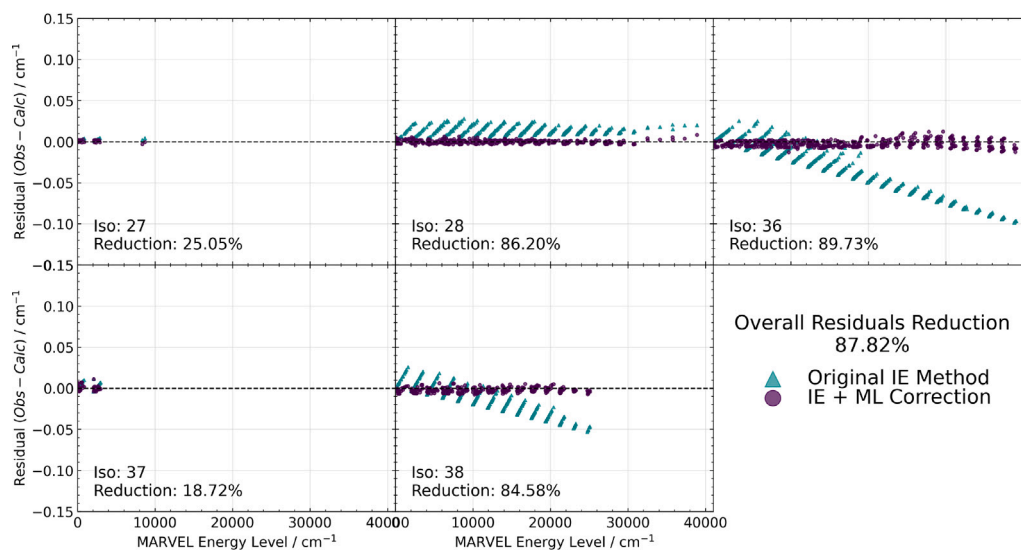


Fig. 10. Residuals for individual CO isotopologues plotted against empirical MARVEL energy levels before and after the ML correction.

**Table 5**

Feature importance for the CO<sub>2</sub> correction network, based on MAE increase upon ablation, highlighting the dominant role of isotopic masses and rovibrational spectroscopic predictors. (Oxygen mass position is denoted in brackets).

Feature	Importance	Feature	Importance
$J$	6.68e-03	<sup>12</sup> C	8.02e-04
<sup>17</sup> O (2)	2.84e-03	AFGL I2	7.30e-04
TROVE $\nu_3$	2.26e-03	Symmetry $A_2$	6.69e-04
Herzberg $\nu_3$	2.23e-03	Herzberg I2	6.53e-04
Herzberg $\nu_1$	2.06e-03	<sup>18</sup> O (1)	6.06e-04
<sup>13</sup> C	1.94e-03	Symmetry $A''$	5.87e-04
AFGL $m_3$	1.75e-03	AFGL $m_2$	5.84e-04
$E_{Ca}^{parent}$	1.70e-03	Symmetry $A'$	3.68e-04
AFGL $m_1$	1.66e-03	Symmetry $A_1$	3.62e-04
Herzberg $\nu_2$	1.53e-03	f	1.18e-04
$E_{Ca}^{iso}$	1.43e-03	$\mu_3$	7.56e-05
$E_{Ma}^{parent}$	1.33e-03	e	6.39e-05
<sup>17</sup> O (1)	1.28e-03	$\mu_2$	2.60e-05
Trove $\nu_2$	1.20e-03	Trove coefficient	1.09e-05
<sup>16</sup> O (1)	1.17e-03	$\mu_1$	8.98e-06
AFGL $r$	1.16e-03	$\mu_3$ ratio	4.13e-06
$\mu_{all}$	1.13e-03	$\mu_{all}$ ratio	2.46e-06
$g_{tot}$	1.12e-03	$\mu_2$ ratio	1.69e-06
$E_{IE}^{iso}$	9.80e-04	$\mu_1$ ratio	1.01e-06
Trove $\nu_1$	9.75e-04	<sup>16</sup> O (2)	1.00e-06
<sup>18</sup> O (2)	8.08e-04		

### 3.4. Inference

Inference for all CO<sub>2</sub> isotopologues was performed on calculated energy levels up to 12 500 cm<sup>-1</sup>, consistent with the energy range of the training set. Applying this cutoff resulted in the update of 36 795 energy levels; a breakdown of these totals for each isotopologue is provided in Table 7. All levels have been used to update the “Dozen” line lists [25].

Inference for all CO isotopologues was performed up to 40 000 cm<sup>-1</sup>, consistent with the energy range of the training set. Applying this cutoff resulted in the update of 3348 energy levels; a breakdown of these totals for each isotopologue is provided in Table 8. These CO data will be used for future isotopologue studies within the ExoMol project.

### 4. Conclusion

The CO<sub>2</sub> and CO isotopologue extrapolation correction networks demonstrate that neural networks can effectively learn physically meaningful relationships between isotopic composition and energy-level deviations. Moving beyond traditional methods that rely on simple, global numerical corrections, this work implements a per-energy-level correction scheme. This granular approach resulted in the correction of 36 795 energy levels for CO<sub>2</sub> and 3348 energy levels for CO, yielding significant improvements in accuracy over previous uniform scaling methods.

The CO<sub>2</sub> model established a robust baseline for isotopologue-specific correction, while the CO network extended this approach by integrating cross-molecular information through a hybrid, molecule-aware architecture. The CO results closely follow the high-accuracy trends of CO<sub>2</sub>, with the key distinction that the hybrid design enables the effective transfer of learned isotopic correction trends between related molecular systems. This transfer learning approach improves model robustness, accelerates convergence, and enhances correction accuracy even for isotopologues with limited training data.

Together, these results confirm that carefully structured neural networks can both replicate and generalize spectroscopic corrections across families of chemically related molecules. Future work will focus on extending these capabilities to a broader range of molecules, particularly those containing hydrogen atoms where energy level shifts

are significantly larger and more heavily influenced by non-Born-Oppenheimer effects. This work ultimately supports the production of more accurate and complete spectroscopic data for use across high-resolution astrophysical applications.

### CRedit authorship contribution statement

**Marco G. Barnfield:** Data curation, Formal analysis, Investigation, Methodology, Validation, Visualization, Writing – original draft. **Oleg L. Polyansky:** Conceptualization, Methodology. **Sergei N. Yurchenko:** Data curation, Project administration. **Jonathan Tennyson:** Conceptualization, Funding acquisition, Methodology, Project administration, Writing – original draft.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgements

This work was supported by the UCL Center for Doctoral Training in Data Intensive Science, funded by the STFC training grant reference ST/P006736/1, as well as by the STFC Projects ST/Y001508/1 and UKRI/ST/B001183/1, and ERC Advanced Investigator Project 883830 (ExoMolHD).

### Appendix A. Reduced mass equations for CO<sub>2</sub>

Effective reduced mass for symmetric stretch mode.

$$\mu_1 = \left( \frac{\text{mass}_O^1 \times \text{mass}_C}{\text{mass}_O^1 + \text{mass}_C} + \frac{\text{mass}_O^2 \times \text{mass}_C}{\text{mass}_O^2 + \text{mass}_C} \right) \div 2 \quad (\text{A.1})$$

Reduced mass for bending mode (Oxygen–Oxygen pair).

$$\mu_2 = \frac{\text{mass}_O^1 \times \text{mass}_O^2}{\text{mass}_O^1 + \text{mass}_O^2} \quad (\text{A.2})$$

Effective reduced mass for asymmetric stretch mode.

$$\mu_3 = \frac{(\text{mass}_O^1 + \text{mass}_O^2) \times \text{mass}_C}{\text{mass}_O^1 + \text{mass}_O^2 + \text{mass}_C} \quad (\text{A.3})$$

Overall reduced mass of the triatomic system.

$$\mu_{all} = \frac{\text{mass}_O^1 \times \text{mass}_C \times \text{mass}_O^2}{\text{mass}_O^1 + \text{mass}_C + \text{mass}_O^2} \quad (\text{A.4})$$

### Appendix B. Additional features for combined CO & CO<sub>2</sub> dataset

Feature name	Description	Data type
$\nu$	Vibrational energy level	Float
$\mu$	Reduced Mass	Float
$\mu$ ratio	Ratio of isotopologue’s $\mu$ to the parent isotopologue’s $\mu$	Float

All features not applicable to CO, e.g., AFGL/Herzberg/TROVE quantum numbers,  $\mu_{1-3}$ , and secondary oxygen masses, were passed into the model as zeros.

### Appendix C. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jms.2026.112084>.

**Table 6**

Comparison of empirical (MARVEL), variational (Calc.), and machine learning corrected (ML) energy levels for the  $\nu = 10$  vibrational state of the  $^{12}\text{C}^{18}\text{O}$  isotopologue (28). All energy and residual values are given in  $\text{cm}^{-1}$ .

J	MARVEL (Ma)	Calc. (Ca)	Resid. (Ma – Ca)	ML	Resid. (Ma – ML)
0	19794.0437	19794.0350	$8.65 \times 10^{-3}$	19794.0423	$1.38 \times 10^{-3}$
1	19797.3802	19797.3718	$8.38 \times 10^{-3}$	19797.3804	$-2.01 \times 10^{-4}$
2	19804.0539	19804.0453	$8.63 \times 10^{-3}$	19804.0562	$-2.32 \times 10^{-3}$
3	19814.0635	19814.0551	$8.42 \times 10^{-3}$	19814.0649	$-1.44 \times 10^{-3}$
4	19827.4095	19827.4009	$8.64 \times 10^{-3}$	19827.4108	$-1.29 \times 10^{-3}$
5	19844.0905	19844.0821	$8.40 \times 10^{-3}$	19844.0913	$-7.52 \times 10^{-4}$
6	19864.1068	19864.0981	$8.64 \times 10^{-3}$	19864.1081	$-1.29 \times 10^{-3}$
7	19887.4564	19887.4482	$8.24 \times 10^{-3}$	19887.4553	$1.11 \times 10^{-3}$
8	19914.1399	19914.1312	$8.71 \times 10^{-3}$	19914.1397	$2.48 \times 10^{-4}$
9	19944.1543	19944.1463	$8.07 \times 10^{-3}$	19944.1549	$-5.64 \times 10^{-4}$
10	19977.5008	19977.4921	$8.69 \times 10^{-3}$	19977.5013	$-5.20 \times 10^{-4}$

**Table 7**

Number of previously variationally calculated energy levels corrected using the IE ML method for each minor  $\text{CO}_2$  isotopologue.

Isotopologue	No. energy levels
$^{16}\text{O}^{12}\text{C}^{17}\text{O}$ (627)	1791
$^{16}\text{O}^{12}\text{C}^{18}\text{O}$ (628)	1053
$^{16}\text{O}^{13}\text{C}^{17}\text{O}$ (637)	2917
$^{16}\text{O}^{13}\text{C}^{18}\text{O}$ (638)	2352
$^{12}\text{C}^{17}\text{O}_2$ (727)	7040
$^{17}\text{O}^{12}\text{C}^{18}\text{O}$ (728)	2805
$^{13}\text{C}^{17}\text{O}_2$ (737)	7041
$^{17}\text{O}^{13}\text{C}^{18}\text{O}$ (738)	3764
$^{12}\text{C}^{18}\text{O}_2$ (828)	3154
$^{13}\text{C}^{18}\text{O}_2$ (838)	4878

**Table 8**

Number of previously variationally calculated energy levels corrected using the IE ML method for each minor  $\text{CO}$  isotopologue.

Isotopologue	No. energy levels
$^{12}\text{C}^{17}\text{O}$	1152
$^{12}\text{C}^{18}\text{O}$	693
$^{13}\text{C}^{16}\text{O}$	545
$^{13}\text{C}^{17}\text{O}$	1144
$^{13}\text{C}^{18}\text{O}$	844

## Data availability

All data used for training are available in the cited papers, see Section 2.1. The code used in this work is publicly accessible in the GitHub project repository, [https://github.com/mbarnfield63/ML\\_Isotopologue\\_Extrapolation.git](https://github.com/mbarnfield63/ML_Isotopologue_Extrapolation.git). All updated energy levels for  $\text{CO}_2$  have been used to update the “Dozen” line list states file available on the ExoMol website, [www.exomol.com](http://www.exomol.com). The ML-corrected  $\text{CO}$  isotopologue levels, in  $\text{cm}^{-1}$ , are available as supplementary material.

## References

- [1] M. Mayor, D. Queloz, A jupiter-mass companion to a solar-type star, *Nature* 378 (1995) 355–359, <http://dx.doi.org/10.1038/378355a0>.
- [2] W.J. Borucki, D. Koch, G. Basri, N. Batalha, T. Brown, D. Caldwell, J. Caldwell, J. Christensen-Dalsgaard, W.D. Cochran, E. DeVore, E.W. Dunham, A.K. Dupree, T.N. Gautier, J.C. Geary, R. Gilliland, A. Gould, S.B. Howell, J.M. Jenkins, Y. Kondo, D.W. Latham, G.W. Marcy, S. Meibom, H. Kjeldsen, J.J. Lissauer, D.G. Monet, D. Morrison, D. Sasselov, J. Tarter, A. Boss, D. Brownlee, T. Owen, D. Buzasi, D. Charbonneau, L. Doyle, J. Fortney, E.B. Ford, M.J. Holman, S. Seager, J.H. Steffen, W.F. Welsh, J. Rowe, H. Anderson, L. Buchhave, D. Ciardi, L. Walkowicz, W. Sherry, E. Horch, H. Isaacson, M.E. Everett, D. Fischer, G. Torres, J.A. Johnson, M. Endl, P. MacQueen, S.T. Bryson, J. Dotson, M. Haas, J. Kolodziejczak, J. Van Cleve, H. Chandrasekaran, J.D. Twicken, E.V. Quintana, B.D. Clarke, C. Allen, J. Li, H. Wu, P. Tenenbaum, E. Verner, F. Bruhweiler, J. Barnes, A. Prsa, Kepler planet-detection mission: Introduction and first results, *Science* 327 (2010) 977–980, <http://dx.doi.org/10.1126/science.1185402>.
- [3] G.R. Ricker, J.N. Winn, R. Vanderspek, D.W. Latham, G.A. Bakos, J.L. Bean, Z.K. Berta-Thompson, T.M. Brown, L. Buchhave, N.R. Butler, R.P. Butler, W.J. Chaplin, D. Charbonneau, J. Christensen-Dalsgaard, M. Clampin, D. Deming, J. Doty, N. De Lee, C. Dressing, E.W. Dunham, M. Endl, F. Fressin, J. Ge, T. Henning, M.J. Holman, A.W. Howard, S. Ida, J.M. Jenkins, G. Jernigan, J.A. Johnson, L. Kaltenegger, N. Kawai, H. Kjeldsen, G. Laughlin, A.M. Levine, D. Lin, J.J. Lissauer, P. MacQueen, G. Marcy, P.R. McCullough, T.D. Morton, N. Narita, M. Paegert, E. Palle, F. Pepe, J. Pepper, A. Quirrenbach, S.A. Rinehart, D. Sasselov, B. Sato, S. Seager, A. Sozzetti, K.G. Stassun, P. Sullivan, A. Szentgyorgyi, G. Torres, S. Udry, J. Villaseñor, Transiting exoplanet survey satellite, *J. Astron. Telesc. Instruments Syst.* 1 (2014) 014003, <http://dx.doi.org/10.1117/1.JATIS.1.1.014003>.
- [4] E.L. Rice, T. Barman, I.S. Mclean, L. Prato, J.D. Kirkpatrick, Physical properties of young brown dwarfs and very low mass stars inferred from high-resolution model spectra, *Astrophys. J. Suppl. Ser.* 186 (2010) 63, <http://dx.doi.org/10.1088/0067-0049/186/1/63>.
- [5] N. Madhusudhan, Exoplanetary atmospheres: Key insights, challenges, and prospects, *Annu. Rev. Astron. Astrophys.* 57 (2019) 617–663, <http://dx.doi.org/10.1146/annurev-astro-081817-051846>.
- [6] P. Mollière, T. Molyarova, B. Bitsch, T. Henning, A. Schneider, L. Kreidberg, C. Eistrup, R. Burn, E. Nasedkin, D. Semenov, C. Mordasini, M. Schlecker, K.R. Schwarz, S. Lacour, M. Nowak, M. Schulik, Interpreting the atmospheric composition of exoplanets: Sensitivity to planet formation assumptions, *Astrophys. J.* 934 (2022) 74, <http://dx.doi.org/10.3847/1538-4357/ac6a56>.
- [7] S.N. Yurchenko, J. Tennyson, M. Brogi, High-resolution spectroscopy of exoplanets: Data challenges and prospects, *Nat. Rev. Phys.* (2025) <http://dx.doi.org/10.1038/s42254-025-00839-z>.
- [8] J.K. Barstow, S. Aigrain, P.G.J. Irwin, S. Kendrew, L.N. Fletcher, Transit spectroscopy with James Webb Space Telescope: Systematics, starspots and stitching, *Mon. Not. R. Astron. Soc.* 448 (2015) 2546–2561, <http://dx.doi.org/10.1093/mnras/stv186>.
- [9] L.S. Wiser, T.J. Bell, M.R. Line, E. Schlawin, T.G. Beatty, L. Welbanks, T.P. Greene, V. Parmentier, M.M. Murphy, J.J. Fortney, K. Arnold, N. Mehta, K. Ohno, S. Mukherjee, A precise metallicity and carbon-to-oxygen ratio for a warm giant exoplanet from its panchromatic JWST emission spectrum, *Proc. Nat. Acad. Sci. U.S.A.* 122 (2024) <http://dx.doi.org/10.1073/pnas.2416193122>.
- [10] I. Snellen, R. de Kok, J.L. Birkby, B. Brando, M. Brogi, C. Keller, M. Kenworthy, H. Schwarz, R. Stuik, Combining high-dispersion spectroscopy with high contrast imaging: Probing rocky planets around our nearest neighbors, *Astron. Astrophys.* 576 (2015) A59, <http://dx.doi.org/10.1051/0004-6361/201425018>.
- [11] I.A. Snellen, Exoplanet atmospheres at high spectral resolution, *Annu. Rev. Astron. Astrophys.* 63 (2025) 83–125, <http://dx.doi.org/10.1146/annurev-astro-052622-031342>.
- [12] Y. Zhang, I.A.G. Snellen, A.J. Bohn, P. Molliere, C. Ginski, H.J. Hoeijmakers, M.A. Kenworthy, E.E. Mamajek, T. Meshkat, M. Reggiani, F. Snik, The  $^{13}\text{C}$ -rich atmosphere of a young accreting super-Jupiter, *Nature* 595 (7867) (2021) 370–372, <http://dx.doi.org/10.1038/s41586-021-03616-x>.
- [13] E. Esparza-Borges, M. López-Morales, J.I. Adams Redai, J. Kirk E. Pallé, N. Casasayas-Barris, N.E. Batalha, B.V. Rackham, J.L. Bean, S.L. Casewell, L. Decin, L.A. Dos Santos, J. Harrington A.G. Muñoz, G. Heng, R. Hu, L. Mancini, K. Molaverdikhani, G. Morello, N.K. Nikolov, M.C. Nixon, S. Redfield, K.B. Stevenson, H.R. Wakeford, M.K. Alam, B. Benneke, J. Blečić, N. Crouzet, T. Daylan, J. Inglis, L. Kreidberg, D.J.M. Petit dit de la Roche, J.D. Turner, Detection of carbon monoxide in the atmosphere of WASP-39b applying standard cross-correlation techniques to JWST NIRSpec G395H data, *Astrophys. J. Lett.* 955 (2023) L19, <http://dx.doi.org/10.3847/2041-8213/acf27b>.
- [14] M. Brogi, M.R. Line, Retrieving temperatures and abundances of exoplanet atmospheres with high-resolution cross-correlation spectroscopy, *Astrophys. J.* 157 (2019) 114, <http://dx.doi.org/10.3847/1538-3881/aafdd3>.
- [15] J. Tennyson, S.N. Yurchenko, ExoMol: Molecular line lists for exoplanet and other atmospheres, *Mon. Not. R. Astron. Soc.* 425 (2012) 21–33, <http://dx.doi.org/10.1111/j.1365-2966.2012.21440.x>.

- [16] T. Furtenbacher, A.G. Császár, J. Tennyson, MARVEL: Measured active rotational-vibrational energy levels, *J. Mol. Spectrosc.* 245 (2007) 115–125, <http://dx.doi.org/10.1016/j.jms.2007.07.005>.
- [17] J. Tennyson, S.N. Yurchenko, J. Zhang, C.A. Bowesman, R.P. Brady, J. Buldyreva, K.L. Chubb, R.R. Gamache, M.N. Gorman, E.R. Guest, C. Hill, K. Kefala, A.E. Lynas-Gray, T.M. Mellor, L.K. McKemmish, G.B. Mitev, I.I. Mizus, A. Owens, Z. Peng, A.N. Perri, M. Pezzella, O.L. Polyansky, Q. Qu, M. Semenov, O. Smola, A. ov, W. Somogyi, A. Upadhyay, S.O.M. Wright, N.F. Zobov, The 2024 release of the ExoMol database: Molecular line lists for exoplanet and other hot atmospheres, *J. Quant. Spectrosc. Radiat. Transfer* 326 (2024) 109083, <http://dx.doi.org/10.1016/j.jqsrt.2024.109083>.
- [18] O.L. Polyansky, A.A. Kyuberis, L. Lodi, J. Tennyson, R.I. Ovsyannikov, N. Zobov, ExoMol molecular line lists XIX: High accuracy computed line lists for H<sub>2</sub><sup>17</sup>O and H<sub>2</sub><sup>18</sup>O, *Mon. Not. R. Astron. Soc.* 466 (2017) 1363–1371, <http://dx.doi.org/10.1093/mnras/stw3125>.
- [19] L.K. McKemmish, C.A. Bowesman, K. Kefala, A.N. Perri, A.M. Syme, S.N. Yurchenko, J. Tennyson, A hybrid approach to generating diatomic line lists for high resolution studies of exoplanets and other hot astronomical objects: Updates to ExoMol MgO, VO and TiO line lists, *RAS Tech. Instr.* 3 (2024) 565–583, <http://dx.doi.org/10.1093/rasti/rzae037>.
- [20] Y.V. Pavlenko, S.N. Yurchenko, L.K. McKemmish, J. Tennyson, Analysis of the TiO isotopologues in stellar optical spectra, *Astron. Astrophys.* 42 (2020) A77, <http://dx.doi.org/10.1051/0004-6361/202037863>.
- [21] D.W. Schwenke, Beyond the potential energy surface: Ab initio corrections to the born-oppenheimer approximation for H<sub>2</sub>O, *J. Phys. Chem. A* 105 (2001) 2352–2360, <http://dx.doi.org/10.1021/jp0032513>.
- [22] K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O. von Lilienfeld, K.-R. Müller, A. Tkatchenko, Machine learning predictions of molecular properties: Accurate many-body potentials and nonlocality in chemical space, *J. Phys. Chem. Lett.* 6 (2015) 2326–2331, <http://dx.doi.org/10.1021/acs.jpclett.5b00831>.
- [23] J. Westermayr, P. Marquetand, Machine learning spectroscopy to advance computation and analysis, *Chem. Sci.* 46 (16) (2025) 21660–21676, <http://dx.doi.org/10.1039/d5sc05628d>.
- [24] E.R. Guest, J. Tennyson, S.N. Yurchenko, Modelling the rotational dependence of line broadening using machine learning, *J. Mol. Spectrosc.* 401 (2024) 111901, <http://dx.doi.org/10.1016/j.jms.2024.111901>.
- [25] S.N. Yurchenko, M.G. Barnfield, C.A. Bowesman, R.P. Brady, E.R. Guest, K. Kefala, Q.-H. Ni, A.N. Perri, O.A. Smola, A. Solokov, C. Tao, J. Tennyson, ExoMol line lists – LXIII: ExoMol line lists for 12 isotopologues of CO<sub>2</sub>, *Mon. Not. R. Astron. Soc.* 545 (2026) staf2135, <http://dx.doi.org/10.1093/mnras/staf2135>.
- [26] M.T.I. Ibrahim, D. Alatoom, T. Furtenbacher, A.G. Császár, S.N. Yurchenko, A.A.A. Azzam, J. Tennyson, MARVEL analysis of high-resolution rovibrational spectra of <sup>13</sup>C<sup>16</sup>O<sub>2</sub>, *J. Comput. Chem.* 45 (2024) 969–984, <http://dx.doi.org/10.1002/jcc.27266>.
- [27] D. Alatoom, M.T.I. Ibrahim, T. Furtenbacher, A.G. Császár, M. Alghizzawi, S.N. Yurchenko, A.A.A. Azzam, J. Tennyson, MARVEL analysis of high-resolution rovibrational spectra of <sup>16</sup>O<sup>12</sup>C<sup>18</sup>O, *J. Comput. Chem.* 45 (2024) 2558, <http://dx.doi.org/10.1002/jcc.27453>.
- [28] A.A.A. Azzam, S.A.A. Azzam, K.A.A. Aburumman, J. Tennyson, S.N. Yurchenko, A.G. Császár, T. Furtenbacher, MARVEL analysis of high-resolution rovibrational spectra of <sup>18</sup>O<sup>12</sup>C<sup>18</sup>O, <sup>17</sup>O<sup>12</sup>C<sup>18</sup>O and <sup>18</sup>O<sup>13</sup>C<sup>18</sup>O isotopologues of carbon dioxide, *J. Mol. Spectrosc.* 405 (2024) 111947, <http://dx.doi.org/10.1016/j.jms.2024.111947>.
- [29] A.A.A. Azzam, B.M.J. Abou Doud, M.Q.A. Shersheer, B.K.M. Almasri, C.N.M. Bader, A.M.H.A. Baraa O.A. KH. Musleh, A.W.M. Al Shatarat, B.I.M. Qattan, L.H.M. Hamamsy, A.O.G. Saafneh, M.N.A. ALso'ub, M.M.A. Alkhashashneh, H.O.M. Al-Zawahra, D. Alatoom, M.T.I. Ibrahim, J. Tennyson, S.N. Yurchenko, T. Furtenbacher, A.G. Császár, The 626M24 dataset of validated transitions and empirical rovibrational energy levels of <sup>16</sup>O<sup>12</sup>C<sup>16</sup>O, *Sci. Data* 12 (2025) 532, <http://dx.doi.org/10.1038/s41597-025-04755-w>.
- [30] A.A.A. Azzam, J. Tennyson, S.N. Yurchenko, T. Furtenbacher, A.G. Császár, MARVEL analysis of high-resolution rovibrational spectra of <sup>16</sup>O<sup>13</sup>C<sup>18</sup>O, *J. Comput. Chem.* 46 (2025) e27541, <http://dx.doi.org/10.1002/jcc.27541>.
- [31] S.A.M. Obaidata, A.A.A. Azzam, J. Tennyson, S.N. Yurchenko, T. Furtenbacher, A.G. Császár, MARVEL analysis of high-resolution rovibrational spectra of <sup>16</sup>O<sup>12</sup>C<sup>17</sup>O, *J. Mol. Spectrosc.* 340 (2025) 109444, <http://dx.doi.org/10.1016/j.jqsrt.2025.109444>.
- [32] M.H.I. Mansour, A.A.A. Azzam, J. Tennyson, S.N. Yurchenko, T. Furtenbacher, A.G. Császár, MARVEL analysis of high-resolution rovibrational spectra of <sup>16</sup>O<sup>13</sup>C<sup>17</sup>O and <sup>17</sup>O<sup>13</sup>C<sup>17</sup>O, *Mol. Phys.* (2025) e2550568, <http://dx.doi.org/10.1080/00268976.2025.2550568>.
- [33] A.A.A. Azzam, J.M.A. AlAlawin, J. Tennyson, S.N. Yurchenko, T. Furtenbacher, A.G. Császár, MARVEL analysis of high-resolution rovibrational spectra of <sup>17</sup>O<sup>13</sup>C<sup>18</sup>O and <sup>17</sup>O<sup>13</sup>C<sup>17</sup>O, *J. Quant. Spectrosc. Radiat. Transfer* 343 (2025) 109485, <http://dx.doi.org/10.1016/j.jqsrt.2025.109485>.
- [34] S. Mahmoud, N. El-Kork, N. Abu Elkher, M. Almehairbi, M.S. Khalil, T. Furtenbacher, O.P. Yurchenko, S.N. Yurchenko, J. Tennyson, MARVEL analysis of the measured high-resolution spectra of <sup>12</sup>C<sup>16</sup>O, *Astrophys. J. Suppl. Ser.* 276 (2025) 66, <http://dx.doi.org/10.3847/1538-4365/ada3c9>.
- [35] T. Grigorev, Y. Dai, M. Potter, X. Xiang, K. Zhang, J. Tennyson, MARVEL analysis of the measured high-resolution spectra of CO isotopologues, *Astrophys. J. Suppl. Ser.* 283 (2026) 39, <http://dx.doi.org/10.3847/1538-4365/ae40f0>.
- [36] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala, Pytorch: An imperative style, high-performance deep learning library, in: *Advances in Neural Information Processing Systems* 32, Curran Associates, Inc., 2019, pp. 8024–8035.
- [37] D. Hendrycks, K. Gimpel, Gaussian error linear units (gelus), 2023, <http://dx.doi.org/10.48550/arXiv.1606.08415>, arXiv.
- [38] A.F. Agarap, Deep learning using rectified linear units (relu), 2019, <http://dx.doi.org/10.48550/arXiv.1803.08375>, arXiv.
- [39] D.P. Kingma, J. Ba, Adam: A method for stochastic optimization, 2017, <http://dx.doi.org/10.48550/arXiv.1412.6980>, arXiv.
- [40] P. Huber, Robust estimation of a location parameter, *Ann. Math. Stat.* 35 (1964) 73–101, <http://dx.doi.org/10.1214/aoms/1177703732>.
- [41] J.L. Ba, J.R. Kiros, G.E. Hinton, Layer normalization, 2016, <http://dx.doi.org/10.48550/arXiv.1607.06450>, arXiv.
- [42] L.S. Rothman, L.D.G. Young, Infrared energy levels and intensities of carbon dioxide-II, *J. Quant. Spectrosc. Radiat. Transfer* 25 (1981) 505–524, [http://dx.doi.org/10.1016/0022-4073\(81\)90026-1](http://dx.doi.org/10.1016/0022-4073(81)90026-1).