



VaPOrS v1.0.1: An automated model for estimating vapor pressure of organic compounds using SMILES notation

Mojtaba Bezaatpour^{1,a}, Miikka Dal Maso^a, Matti Rissanen^{2, a, b}

^a *Aerosol Physics Laboratory, Tampere University, 33720 Tampere, Finland*

^b *Department of Chemistry, University of Helsinki, 00014 Helsinki, Finland*

Abstract

Volatile organic compounds play a significant role in atmospheric chemistry, influencing air quality and climate change. Accurate prediction of their physical properties is essential for understanding their behavior. This paper introduces the VaPOrS (**V**apor **P**ressure in **O**rganics via **S**MILES) as a comprehensive tool designed to process SMILES notation of organic compounds, identify key functional groups, and calculate their saturation vapor pressure and enthalpy of vaporization at any specified temperature. While this first study focuses on applying the SIMPOL method for parameterization, VaPOrS is inherently adaptable to other structure-based parameterization approaches, such as group additivity and volatility basis set (VBS) methods by extracting substructure information from each string that is meaningful to property predictive techniques. It can also be extended to any thermodynamic property that relies on structural group-based parameterizations. In its current version, the tool automates the detection of 30 critical structural groups and has been validated against manually counted functional groups and experimental saturation vapor pressure data for a diverse set of compounds. The results demonstrate high accuracy, with the tool correctly identifying the same functional groups, followed by providing prompt saturation vapor pressure predictions according to the SIMPOL parameterization. The developed method can be integrated into large-scale simulation models targeting secondary aerosol formation and involving thousands of organic species at once. Thus, the developed tool offers a robust computational approach for research in atmospheric chemistry and environmental science, allowing to streamline the

¹ Corresponding author email address: mojtaba.bezaatpour@tuni.fi (M. Bezaatpour)

² Corresponding author email address: matti.rissanen@tuni.fi (M. Rissanen)



1 analysis of a large collection of organic compounds, aiding in the assessment of their climatic
2 impacts.

3 **Keywords:** Saturation vapor pressure, Volatile organic compounds, VaPOrS, SMILES
4 notation, Functional group, Atmospheric chemistry, Secondary organic aerosol.

5 **1. Introduction**

6 Volatile organic compounds (VOCs) are a diverse group of organic chemicals that significantly
7 impact atmospheric chemistry. Their volatility allows them to easily enter the atmosphere,
8 where they participate in complex chemical reactions that influence air quality, climate, and
9 human health (Mellouki, Wallington, and Chen 2015). VOCs are key precursors to both
10 ground-level ozone and secondary organic aerosols (SOA). Ground-level ozone, a harmful air
11 pollutant, forms through the photochemical reactions of VOCs with nitrogen oxides (NO_x) in
12 the presence of sunlight (Atkinson 2000). Ozone is a major component of urban smog and
13 poses serious health risks, including respiratory problems and cardiovascular disease (WHO
14 2005). Secondary organic aerosols, on the other hand, result from the oxidation of VOCs,
15 leading to the formation of particulate matter that can scatter sunlight and affect the Earth's
16 radiative balance (Jimenez et al. 2009). These particles contribute to atmospheric haze, impact
17 visibility, and have been linked to adverse health effects, including lung and heart diseases.
18 Understanding the formation, transformation, and fate of VOCs is thus crucial for predicting
19 their impact on both air quality and climate.

20 VOCs originate from a variety of sources, both natural and anthropogenic. Natural sources
21 include biogenic emissions from vegetation, such as isoprene and monoterpenes, which are
22 released in large quantities, especially in forested areas (Alex Guenther et al. 1995).
23 Anthropogenic sources are primarily related to industrial activities, transportation, fuel
24 combustion, and the use of solvents and consumer products (Goldstein and Galbally 2007;
25 McDonald et al. 2018). The structural diversity of VOCs poses a challenge for their
26 classification and analysis. They can range from simple hydrocarbons like methane to complex
27 multifunctional molecules containing oxygen, nitrogen, sulfur, and, for example, halogens.
28 This diversity affects their physical properties, reactivity, and environmental behavior, making
29 it necessary to study their properties comprehensively.



1 The physical properties of VOCs, particularly saturation vapor pressure and enthalpy of
2 vaporization, play a critical role in determining their volatility and phase behavior in the
3 atmosphere. Saturation vapor pressure is a measure of a substance's tendency to evaporate, and
4 it directly influences the distribution of VOCs between the gas and particulate phases (Mattila,
5 Kulmala, and Vesala 1997). Compounds with high saturation vapor pressure are more likely to
6 remain in the gas phase, whereas those with lower saturation vapor pressure can condense onto
7 particulate matter, contributing to SOA formation (Donahue et al. 2006). Enthalpy of
8 vaporization, the energy required to convert a substance from liquid to vapor, is an important
9 property that influences the temperature dependence of saturation vapor pressure. Accurate
10 knowledge of these properties is essential for modeling VOC emissions, their transport, and
11 their transformation in the atmosphere.

12 The oxidation of VOCs significantly changes their chemical properties. An extreme example
13 is provided by autoxidation, a chain-like oxidation process that propagates by sequential
14 additions of molecular oxygen and resulting peroxy radical rearrangement reactions, which
15 infuses multiple oxygen molecules to the hydrocarbon backbone. This process plummets the
16 saturation vapor pressures of the participating VOC and ultimately forms so-called Highly
17 Oxygenated organic Molecules (HOMs as expressed by Bianchi et al. (Bianchi et al. 2019;
18 Crounse et al. 2013)), crucial intermediates in the formation of SOA. HOM formation is a
19 common property of VOCs and is thus initiated by all common oxidants capable of forming
20 alkyl radicals, including hydroxyl radicals (OH), ozone (O₃), and nitrate radicals (NO₃) (Zhao
21 et al. 2021; Luo et al. 2023; Rissanen et al. 2014; Berndt et al. 2016) as well as by direct UV
22 photolysis. The resulting HOMs can rapidly condense onto existing particles or sometime even
23 form new particles through nucleation processes, thus contributing to the formation and growth
24 of SOA. The formation of HOMs is heavily influenced by the structure of the precursor VOCs
25 (Ehn et al. 2014). Recent advancements in mass spectrometry and atmospheric simulation
26 techniques have provided valuable insights into the chemical pathways leading to HOMs
27 formation and their role in atmospheric chemistry (Iyer 2023; Bianchi et al. 2019; Vereecken
28 et al. 2018; Michael E. Jenkin et al. 2019; Iyer et al. 2021a; 2021b). Understanding the volatility
29 of VOCs, and especially HOMs, is therefore essential for improving models of SOA formation.

30 Traditionally, the determination of saturation vapor pressure and enthalpy of vaporization has
31 relied on experimental methods, such as gas chromatography-mass spectrometry (GC-MS)
32 (Epping and Koch 2023). While these methods provide accurate measurements, they are time-
33 consuming, expensive, and often impractical for large-scale studies involving thousands of



1 compounds. Moreover, many of the condensable chemicals relevant to atmospheric SOA
2 formation are challenging to measure because their parent molecules are unstable, have never
3 been synthesized, or, importantly, cannot be synthesized due to their inherent instability. These
4 limitations highlight the critical need for approximative methods to estimate the physical
5 properties of these compounds, as direct measurements are often unfeasible.

6 Given the limitations of experimental approaches, there is an obvious need for computational
7 tools that can predict these properties efficiently. Such tools can complement experimental
8 methods by providing rapid estimates for a wide range of compounds, facilitating large-scale
9 atmospheric modeling studies. Various computational methods have been developed to
10 estimate the saturation vapor pressure of organic compounds, complementing or replacing
11 traditional experimental approaches. Among these, predictive models like COSMO-RS
12 (COnductor-like Screening MOdel for Real Solvents) use quantum chemistry calculations to
13 estimate thermodynamic properties, including saturation vapor pressure. COSMO-RS
14 simulates the solvent environment and molecular interactions, providing a theoretical basis for
15 property estimation. However, its computational intensity can be a drawback for large-scale
16 applications (Klamt 1995; Eckert and Klamt 2002; Klamt et al. 1998). On the other hand, group
17 contribution methods estimate saturation vapor pressure based on the contribution of structural
18 groups within a molecule. These methods leverage empirical correlations derived from
19 extensive datasets but often lack precision for complex or highly functionalized molecules
20 (Joback and Reid 1987; Myrdal and Yalkowsky 1997). The Nannoolal method is a
21 computational approach used to estimate the boiling points and saturation vapor pressures of
22 organic compounds based on their molecular structure. This method utilizes group contribution
23 techniques, where the overall properties of a molecule are determined by summing the
24 contributions of its individual structural groups, such as functional groups and bonding patterns
25 (Nannoolal, Rarey, and Ramjugernath 2008). The EVAPORATION method is a tool designed
26 for estimating the saturation vapor pressure of organic molecules, accounting for contributions
27 from the carbon skeleton, functional groups, and their interactions, and it adjusts for
28 functionalized diacids with empirical modifications. It predicts the saturation vapor pressure of
29 various compounds using only molecular structure as input, making it applicable to a wide
30 range of molecules (Compernelle, Ceulemans, and Müller 2011). The SIMPOL method,
31 developed by Pankow and Asher (Pankow and Asher 2008), is a widely used group contribution
32 method that correlates the presence of specific functional groups in a molecule to its vapor
33 pressure. By summing the contributions of individual functional groups, the method provides



1 an estimate of the compound's vapor pressure. The SIMPOL method has been validated against
2 experimental data and is recognized for its accuracy and applicability to a wide range of organic
3 compounds. Importantly, since all these methods employ group contribution techniques, they
4 could potentially be integrated into automated tools like VaPOrS, enabling the efficient
5 handling of large datasets and complex compounds with minimal manual effort.

6 The present study introduces a Python-based computational tool named VaPOrS (Vapor
7 Pressure in **O**rganics via **S**MILES) to process SMILES (Simplified Molecular Input Line Entry
8 System) notation of VOCs, identify key functional groups, and calculate their saturation vapor
9 pressure and enthalpy of vaporization at any specified temperature. A key distinction between
10 VaPOrS and existing tools such as UManSysProp (Topping et al. 2016) lies in their approach
11 to parsing SMILES notation for functional group identification. While both tools employ group
12 contribution methods for property estimation, VaPOrS explicitly searches for all possible
13 patterns of a specific functional group (e.g., aldehyde) directly from the SMILES string,
14 without relying on external libraries. More details on the specific patterns and their
15 implementation are provided in the Methodology section. In contrast, UManSysProp utilizes
16 SMARTS strings within the OpenBabel framework to identify molecular substructures relevant
17 to predictive methods. The advantage of VaPOrS' approach is that it ensures full control over
18 pattern-matching logic, potentially making it more adaptable to new group definitions without
19 dependency on predefined SMARTS rules. Currently, the tool automates the SIMPOL method
20 by detecting 30 functional groups, which are critical for determining the physical properties of
21 VOCs. This initial implementation focuses on the SIMPOL method as a starting point;
22 however, future iterations of VaPOrS aim to include additional parameterization approaches,
23 such as group additivity and volatility basis set (VBS) methods, further expanding its utility
24 and adaptability.

25 The development of this tool addresses several key challenges:

- 26 1. Automation of functional group detection: The tool eliminates the need for manual
27 identification of functional groups, reducing the potential for human error and
28 increasing efficiency.
- 29 2. Rapid and accurate property prediction: The tool leverages the SIMPOL method to
30 provide rapid and accurate predictions of saturation vapor pressure and enthalpy of
31 vaporization. This capability is particularly valuable for large-scale atmospheric
32 simulations targeting SOA formation involving thousands of compounds.



1 3. Scalability and flexibility: The tool is designed to handle large datasets, making it
2 suitable for high-throughput studies. It can process thousands of SMILES strings within
3 seconds, providing quick insights into the properties of VOCs.

4 The computational analysis of VOCs is integral to several established atmospheric databases
5 and models, which simulate the chemical and physical processes in the atmosphere.
6 Incorporating the VaPOrS developed in this study into them can significantly enhance their
7 efficiency and accuracy. As an instance, recently, Pichelstorfer et al. (Pichelstorfer et al. 2024)
8 developed a close-to-mechanistic approach called auto-APRAM-fw for predicting the
9 formation and general structure of HOMs during the autoxidation of initial radicals, outputting
10 the molecular structures in SMILES notation. Integrating this vast array of HOMs structures
11 into atmospheric simulation models presents significant challenges, particularly in identifying
12 the functional groups for each HOM and calculating their saturation vapor pressures. This
13 process is essential, as saturation vapor pressure is a critical parameter in understanding the
14 volatility and partitioning behavior of HOMs in atmospheric conditions, influencing SOA
15 formation and growth. The task, however, is resource-intensive when conducted manually or
16 through non-specialized methods. By automatically detecting functional groups and calculating
17 saturation vapor pressures for the generated SMILES strings, the VaPOrS tool can facilitate
18 the efficient integration of HOMs data from auto-APRAM-fw into atmospheric models. This
19 automation not only accelerates the data processing required for simulations but also improves
20 the accuracy of saturation vapor pressure predictions, thus enhancing the reliability of SOA
21 formation simulations. Consequently, VaPOrS holds significant potential for advancing
22 atmospheric chemistry research, particularly in extending auto-APRAM-fw mechanistic model
23 for broader applications under diverse atmospheric conditions. Some of the key atmospheric
24 databases and models that could benefit from this tool are discussed in the Results and
25 discussion section.

26 **2. Methodology**

27 The developed VaPOrS begins by processing the SMILES notation of the target chemical
28 compounds. SMILES strings are used as inputs to identify the structural components of the
29 molecules, such as number of carbon atoms, functional groups, and cyclic structures.
30 Altogether 30 functional groups parameterized in the SIMPOL method are identified. Several
31 functions are required to efficiently parse and interpret the structure of each compound. The
32 model consists of two main parts. The first part is dedicated to identifying and counting the



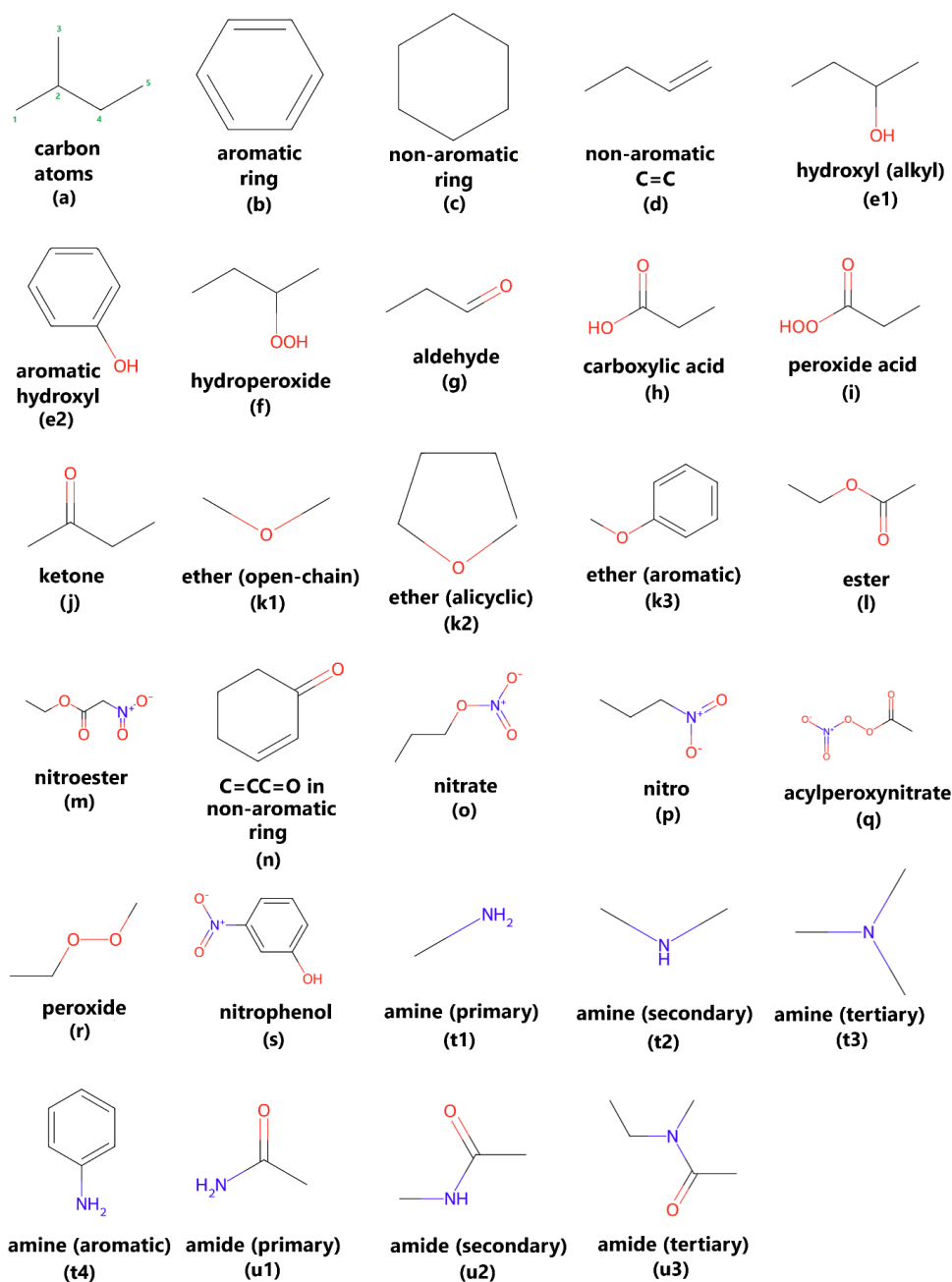
1 number of functional groups within the target while the second part focuses on calculating the
2 saturation vapor pressure as a function of temperature and providing the corresponding Antoine
3 equation parameters.

4 **2.1.Functional group identification**

5 The first part of the model identifies and counts 30 distinct structural groups in organic
6 compounds that are recognized by the SIMPOL method as influential on saturation vapor
7 pressure. These groups include functional components like carbonyl groups (-C=O), hydroxyl
8 groups (-OH), amines (-NR_3), ethers (-R-O-R-), and several others (see Figure 1). The model
9 processes the SMILES string of each compound and applies functions to detect the presence
10 and quantity of these groups. To achieve this, the SMILES string is parsed by several dedicated
11 functions, each corresponding to a specific functional group (e.g., O-atom in a carbonyl or
12 hydroxyl). The results are stored as variables that represent the number of occurrences of each
13 group within the molecule. These counts are then used as input parameters in the saturation
14 vapor pressure calculation. This identification process is critical, as the functional groups
15 detected are the most impactful to the molecule's saturation vapor pressure. By systematically
16 identifying the relevant groups for each compound, the model prepares the necessary data for
17 the next phase of thermodynamic calculations. When using SMILES notations for molecular
18 representation, it is important to exercise caution regarding the starting point of the notation.
19 In particular, the SMILES string must begin from an atom that initiates a branching structure
20 (e.g., position 1, 3, or 5 in Figure 1(a)) rather than from a middle atom (e.g., position 2 or 4 in
21 Figure 1(a)). This ensures that the VaPOrS code correctly interprets the molecular topology
22 and identifies functional groups accurately. This method of writing SMILES is the logical and
23 most straightforward approach that is seen in most databases, sources, and online tools. The
24 following subsections 2.1.1 to 2.1.19 provide detailed descriptions of the first steps of the
25 model computations.

26

27



1

2 **Figure 1.** Overview of chemical functional groups identified and analyzed by the VaPOrS, including number
3 of a) carbon atoms, b) aromatic and c) non-aromatic rings, d) non-aromatic double bonds (C=C), e) hydroxyl
4 (alkyl (e1) and aromatic (e2)), f) hydroperoxide, g) aldehyde, h) carboxylic acid, i) peroxy acid, j) ketone, k)
5 ethers (open-chain (k1), alicyclic (k2), and aromatic (k3)), l) ester, m) nitroester, n) C=CC=O in non-aromatic



1 rings, o) nitrate, p) nitro, q) acylperoxynitrate, r) peroxide, s) nitrophenol, t) amines (primary (t1), secondary
2 (t2), tertiary (t3), and aromatic (t4)), and u) amides (primary (u1), secondary (u2), and tertiary (u3)).

3

4 2.1.1. Carbon atoms

5 Carbon atoms are an essential component in all organic molecules, and their number forms the
6 basis for the calculations of molecular properties utilizing group additivity. The
7 `carbon_number(s)` function counts the occurrences of both uppercase ('C', representing
8 carbon atoms in alkyl chains) and lowercase ('c', representing carbon atoms in aromatic rings)
9 characters in the SMILES string. The total number of carbon atoms is then calculated by
10 summing these occurrences.

11 2.1.2. Identification of branches and rings in SMILES

12 Since SMILES notation uses parentheses to denote branching in molecular structures, and
13 numeric indicators to specify ring closure, the model incorporates specialized functions to
14 handle these aspects.

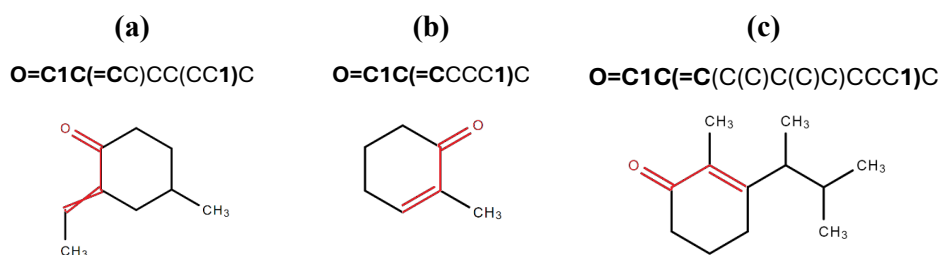
- 15 • The function `find_closing_parenthesis(s, open_index)` is designed to locate the
16 corresponding closing parenthesis for a given open parenthesis in the SMILES string.
17 This is crucial for identifying the boundaries of branched substructures. Similarly,
18 `find_opening_parenthesis(s, close_index)` identifies the opening parenthesis
19 corresponding to a closing parenthesis, enabling proper identification of substructures.
20 These functions ensure that the branching parts of the molecules are accurately
21 interpreted, which is essential for correctly assigning functional groups.
- 22 • The function `find_cycle_number(s)` is developed to parse through the SMILES string
23 and identify the largest numeric index used for ring closure. This index represents the
24 total number of rings in the compound. For example, if a compound contains three
25 rings, its SMILES notation will include ring closure indicators numbered 1, 2, and 3.
26 Each ring is assigned a specific index sequentially as it is encountered while writing the
27 SMILES. Therefore, the largest numeric index in the SMILES indicates the total
28 number of rings present in the structure.

29 To illustrate the importance of correctly matching parentheses in the algorithm, consider the
30 following example. The pattern `O=C1C(=C...1...)` at the beginning of a SMILES notation



1 may suggest the presence of a C=C–C=O group within a non-aromatic ring (such as in
2 cyclohex-2-enone, see example (n) in Figure 1). The parentheses must correspond to each other
3 in the above-mentioned pattern, enclosing the second occurrence of '1'. For instance, in the
4 SMILES string O=C1C(=CC)CC(CC1)C, although it initially appears to match the pattern, the
5 second '1' does not reside between the matching parentheses. As depicted in Figure 2(a), this
6 structure places the C=C bond outside the ring, diverging from the characteristics of the
7 intended functional group. By contrast, in the strings O=C1C(=CCCC1)C and
8 O=C1C(=C(C(C)C(C)C)CCC1)C, the second '1' is properly enclosed within the matching
9 parentheses (highlighted in bold), adhering to the pattern and satisfying the criteria for
10 detecting the C=C–C=O group within a non-aromatic ring. Notably, the latter example contains
11 more carbon atoms, making it easier to count and visually distinguish the structure's
12 complexity. This highlights the versatility of the algorithm in handling SMILES strings of
13 varying complexity, even when the carbon atom count increases.

14



15 **Figure 2.** Structural representations of SMILES strings demonstrating the importance of
16 correct parenthesis matching for the identification of C=C–C=O groups in non-aromatic
17 rings. Structure (a) (O=C1C(=CC)CC(CC1)C) incorrectly places the C=C bond outside the ring,
18 while structures (b) and (c) (O=C1C(=CCCC1)C and O=C1C(=C(C(C)C(C)C)CCC1)C) adhere to
19 the correct pattern.

20 2.1.3. Aromatic and non-aromatic rings

21 Organic molecules often contain cyclic structures, which can either be aromatic (such as
22 benzene rings) or non-aromatic (such as cyclohexane). The model uses two specific functions,
23 `aromatic_ring(s)` and `non_aromatic_ring(s)`, to identify and quantify these rings based
24 on the SMILES notation. The function `aromatic_ring(s)` utilizes the previously described
25 `find_cycle_number(s)` function to identify the largest ring number within the SMILES



1 string. This number represents the highest numeric indicator for cyclic structures (e.g.,
2 c1ccccc1 in SMILES denotes a benzene ring with the digit '1' marking the ring). The function
3 then constructs a list of potential aromatic ring representations, such as 'c1', 'c2', ..., 'cn',
4 where 'n' is the index of the last detected ring showing how many rings are in the structure. It
5 then iterates over the SMILES string, counting the occurrences of these aromatic rings. To
6 avoid double-counting, the final count is halved, since each ring closure is represented by two
7 digits (e.g., c1...c1). The function returns the total number of aromatic rings present in the
8 molecule.

9 Similarly, the `non_aromatic_ring(s)` function detects non-aromatic cyclic structures. Like
10 the aromatic ring detection, the function constructs list of possible ring representations,
11 including combinations such as 'C1', 'N1', 'O1' and so on, based on the largest ring index.
12 The function returns the total number of non-aromatic rings detected in the molecule.

13 2.1.4. *Non-aromatic double-bonded carbon atoms*

14 In addition to identifying cyclic structures, the model also detects non-aromatic double bonds
15 between two carbon atoms, a common feature in organic molecules that strongly influences
16 their chemical and physical properties. The function
17 `double_bound_nonaromatic_carbons(s)` is designed to identify occurrences of double-
18 bonded carbon atoms in non-aromatic structures. This function operates by scanning the
19 SMILES string for occurrences of =C, which indicates a double bond involving a carbon atom.
20 The following steps are used to ensure accurate detection of non-aromatic double bonds:

- 21 • Ring number detection: As with previous functions, `find_cycle_number(s)` is
22 employed to identify the largest numeric ring closure indicator, which is used to
23 distinguish between ring-bound and non-ring-bound carbons.
- 24 • Pattern matching: The function scans the SMILES string for =C, which denotes a double
25 bond with a carbon atom. Upon finding such occurrences, additional checks are
26 performed to ensure the carbon atoms involved in the double bond are not part of an
27 aromatic system:
 - 28 i. Straight chains: If the string preceding the double bond (=C) contains another
29 capital c, it is counted as a non-aromatic double bond (e.g., C=C in CC=CCCC).



- 1 ii. Ring systems: If the string preceding the double bond contains a numeric ring
2 closure indicator (e.g., C1=C in C1=CCCC1), it is also recognized as part of a
3 non-aromatic ring structure.
- 4 iii. Parentheses handling: The function is equipped to handle more complex
5 structures, such as those involving nested parentheses (e.g., C(...(...
6 (...))...)=C in CC(C(C(CC)C)C)=CC). In SMILES, nested parentheses
7 represent branching in a molecule occurring when an atom in the main chain of
8 the molecule is connected to one or more side chains. The parentheses indicate
9 the start and end of each branch, and nesting occurs when a branch contains
10 another branch. It uses the `find_opening_parenthesis` function to locate the
11 corresponding opening parenthesis and verify that the bond belongs to a non-
12 aromatic system.

13 This approach ensures that double bonds within non-aromatic systems are accurately
14 counted, even in cases where the SMILES notation involves branching or ring
15 structures.

16 2.1.5. Alkyl hydroxyl and hydroperoxide groups

17 The function `hydroxyl_group(s)` is designed to identify hydroxyl groups (-OH) present in
18 the SMILES representation of a compound. To ensure accurate detection, the function
19 incorporates specific conditions to avoid miscounting hydroxyl groups that are in specific
20 arrangements that do not denote targeted simple hydroxyl groups. For example, it ensures that
21 carboxylic acids with hydroxyl groups connected to carbonyls (C=O) are not mistakenly
22 counted as plain hydroxyls. One important condition is that the C(O) pattern (a hydroxyl group
23 in SMILES) must not be followed by =O or (=O), which would indicate a carboxylic acid
24 rather than a hydroxyl group. The following steps outline how this function operates to ensure
25 accurate detection of hydroxyl groups:

- 26 • Ring number detection: The `find_cycle_number(s)` is utilized to determine the
27 presence of cyclic structures to establish whether the hydroxyl group is part of a ring or
28 a straight-chain structure.
- 29 • Pattern matching: The function examines the SMILES string for various patterns that
30 denote hydroxyl groups, primarily focusing on terminal hydroxyls, where the function
31 checks if the hydroxyl group appears at the end of the SMILES string, represented as



1 ‘O’ or ‘C–O’ configurations (e.g., cccccO) or branching with parenthesis (e.g.,
2 ccc(C)O). It evaluates whether the hydroxyl group appears within a branching structure
3 or a cyclic component, recognizing patterns such as C(...)O or C1(...)O. Moreover,
4 the function is equipped to deal with complex SMILES representations involving
5 nested parentheses, ensuring that all possible hydroxyl configurations are evaluated
6 (e.g., (O) in cc(O)cc and cc(C)(O)cc).

7 • Conditions for counting: Specific conditions are implemented to avoid miscounting
8 hydroxyl groups that are either in specific arrangements that do not denote targeted
9 hydroxyl group or are redundant in cyclic structures. In this regard, the function ensures
10 that carboxylic acids with hydroxyl groups connected to carbonyls (C=O) are not
11 mistakenly counted as plain hydroxyls. For example, in the SMILES representation
12 C(O)=O, the hydroxyl group is not counted because it is directly bonded to a carbonyl
13 group. As an example, one condition to take into consideration is that the C(O) pattern
14 (a hydroxyl group in SMILES) does not proceed with =O and (=O). Additional criteria
15 influencing this classification are embedded in the code.

16 It's important to note that the detection of hydroperoxide groups (-OOH) in VaPOrS follows a
17 similar pattern as that used for hydroxyl groups (-OH). The key difference in recognizing
18 hydroperoxide groups is that instead of detecting single oxygen atoms (O), the algorithm would
19 look for two consecutive oxygen atoms (OO). This adjustment is made throughout the function
20 `hydroxyl_group(s)` and is provided in a new function as `hydroperoxide_group(s)` to
21 recognize and count hydroperoxide groups in addition to hydroxyl groups.

22 In cyclic structures, hydroxyl groups are excluded if they are connected directly to aromatic
23 rings (e.g., c1ccc(O)cc1). This is due to them being considered as another essential functional
24 group called aromatic hydroxyl in the SIMPOL method. The sequential subsection explains
25 how this functionality is described and detected by the developed algorithm.

26 2.1.6. Aromatic hydroxyl group

27 The function `aromatic_hydroxyl_group(s)` scans the SMILES string for the presence of
28 hydroxyl groups attached to aromatic systems. The function looks for specific patterns that
29 indicate an aromatic hydroxyl group, focusing on the position of the hydroxyl group in the
30 SMILES string:



- 1 i. Hydroxyl at the end: The function checks if the SMILES string ends with `O`,
2 ensuring that the preceding character(s) form part of an aromatic ring. For the
3 SMILES string such as `c1ccccc1O`, where the hydroxyl group is at the end of the
4 SMILES, it detects the hydroxyl group attached to an aromatic ring (`c1`).
5 ii. Hydroxyl at the start: The function checks if the SMILES starts with the `Oc1`
6 pattern, indicating a hydroxyl group attached to the first aromatic ring in the
7 compound. `Oc1cc` would be recognized as having an aromatic hydroxyl group
8 at the start.
9 iii. Hydroxyl in the middle: The function searches for occurrences of the pattern
10 `c(O)` or `c(...O)` where a hydroxyl group is attached to an aromatic ring
11 somewhere in the middle of the SMILES. For the SMILES `c1cc(O)cc1`, the
12 hydroxyl group is identified within the ring. In strings such as `c1c(c...c1O)C`,
13 the hydroxyl group connected to the aromatic ring `c1` is correctly identified.
14 iv. Hydroxyl at the end of a branch: The function checks for cases where a hydroxyl
15 group is part of a branch but still attached to an aromatic ring. For example, in
16 `c1c(c...)O`, where the hydroxyl group is at the end of a branch attached to an
17 aromatic ring, it detects the correct structure.
18 v. Handling nested parentheses: If multiple parentheses are involved, the function
19 ensures that the correct ring and attachment points are identified. For example,
20 the pattern `...c1...(...c1...)O)...` appearing in SMILES like
21 `CC(c1cc(c(c1)CC(C=O)C)O)C` identifies the aromatic hydroxyl correctly.
22 vi. Hydroxyl connected to the index of aromatic carbon: The function ensures that
23 hydroxyl groups are also counted if they are directly attached to index of carbon
24 atom in an aromatic ring (represented by `c` in SMILES notation). In
25 `c1(O)cccc1`, the hydroxyl group attached to `c1` is correctly counted.

26 As the final filtration step, after identifying a hydroxyl group in the aromatic ring, the algorithm
27 also checks for the presence of a nitro group within the same aromatic ring. If a nitro group is
28 found in addition to the hydroxyl group, the compound is no longer classified as an aromatic
29 hydroxyl group; instead, it is categorized as a nitrophenol compound. A detailed explanation
30 of this functional group is provided in a subsequent subsection.

31 2.1.7. Aldehyde group



1 The function `aldehyde_group(s)` is tasked with locating aldehyde groups (i.e., terminal
2 carbonyl groups) within a given SMILES string. The following steps illustrate how this
3 function operates:

4 • Pattern matching: The function checks the SMILES string for various patterns that
5 indicate the presence of aldehyde groups, focusing on:

6 i. Aldehyde at the beginning: The function checks if the SMILES starts with
7 common aldehyde patterns. For instance, In O=CC, C(=O)C, O=Cc, and C(=O)c
8 patterns appearing in the SMILES, the function counts these as aldehydes.

9 ii. Aldehyde at the end: The function checks if the SMILES ends with the
10 characteristic C=O pattern. For example, the last three characters are checked to
11 be C=O. If preceded by a c, such as in CC=O, it counts as an aldehyde. As another
12 example with cyclic structures, for C1C=O as the last characters of a string where
13 the aldehyde is linked to a cyclic structure, this is also counted.

14 iii. Aldehyde in the middle of the SMILES notation (not structure): The function
15 scans for C=O patterns within the SMILES string. For the structure such as
16 C(C=O)..., the carbonyl group may appear in the middle of the SMILES string.
17 However, this still corresponds to a terminal carbonyl group in the molecular
18 structure, i.e., an aldehyde and not a ketone. This distinction is important: while
19 the SMILES position may suggest a non-terminal group, the actual bonding
20 context in the molecular structure confirms its identity as an aldehyde.

21 iv. Aldehyde at the end of a branch: The function examines occurrences of C=O) at
22 the end of branches. For example, the appearance of CC=O) ... pattern in the
23 SMILES counts it as an aldehyde.

24 • Conditions for counting: Specific conditions are established to ensure accurate counting
25 and avoid misidentification of aldehyde groups:

26 i. Connected to carbon atoms: The function ensures that aldehydes are counted
27 only if they are connected to carbon atoms directly. In C(C=O)... pattern, for
28 example, the c preceding the C=O indicates a valid aldehyde.

29 ii. Handling cyclic structures: The function accounts for rings, ensuring that
30 aldehydes connected to cyclic structures are counted appropriately. For
31 example, in C1(...)C=O, the presence of C1 before the aldehyde indicates a
32 connection to a ring, thus counting it as an aldehyde.



1 iii. Parentheses handling: The function effectively manages nested structures,
2 checking for relevant connections before counting. As an instance, if C=O is
3 branched from a carbon e.g., in C(...)(C=O)... pattern, it counts as an
4 aldehyde.

5 2.1.8. Carboxylic and peroxy acid groups

6 The function `carboxylic_acid_group(s)` is designed to identify the presence of carboxylic
7 acid groups within a compound. The detection is briefly broken down into the following steps:

- 8 • Carboxylic acid at the beginning: The function begins by checking if the SMILES string
9 starts with the characteristic patterns for a carboxylic acid group. Allowed prefixes are
10 O=C(O) or OC(=O). For example, in the SMILES string O=C(O)CC, the carboxyl group
11 O=C(O) at the beginning of the molecule is detected.
- 12 • Carboxylic acid as the last characters: The function checks if the SMILES string ends
13 with patterns such as C(=O)O and C(O)=O, indicating a carboxylic acid group at the end
14 of the molecule. As an instance, for the compound CCC(=O)O, the carboxylic acid group
15 at the end is correctly identified as C(=O)O.
- 16 • Carboxylic acid in the middle of the SMILES: The function searches for occurrences
17 of the carboxylic acid group in the middle of the SMILES string using the patterns like
18 C(=O)O and C(O)=O. Each time one of these patterns is found, the function increases
19 the count of carboxylic acid groups. For instance, in the SMILES
20 CC(C(=O)O)C(C(O)=O)C, the carboxylic acid groups C(=O)O and C(O)=O in the middle
21 are detected.

22 The detection of peroxy acid groups follows a similar pattern to that used for carboxylic acid
23 groups within the model. The key difference in recognizing peroxy acid groups lies in the
24 algorithm's adjustment to look for two consecutive oxygen atoms e.g., C(=O)OO and C(OO)=O
25 instead of single oxygen atoms found in carboxylic acids e.g., C(=O)O and C(O)=O. By
26 implementing this modification throughout the model, peroxy acid functional groups can be
27 counted as well.

28 2.1.9. Ketone group



1 The function `ketone_group(s)` is designed to identify and count the presence of ketone groups
2 in the SMILES string of a compound. These patterns can be located at the start, middle, or end
3 of the SMILES string, as well as within rings. The detection process is divided into several
4 steps:

- 5 • Ketone as the first characters: The function checks if the SMILES string begins with a
6 ketone pattern (e.g., O=C(C...)C in O=C(CCC)C or O=C(C...)C... in
7 O=C(c1cccc1)CCO). These patterns represent the ketone group at the beginning of the
8 SMILES string, followed by either a non-aromatic or aromatic carbon. The function
9 also handles ketones connected to rings, such as O=C1... in O=C1CCCC1 indicating a
10 carbonyl group attached to the first position of a ring. Note that although the ketone
11 group appears at the start of the SMILES string, it may not be at the beginning of the
12 molecular structure itself. The pattern is recognized based on bonding context, not
13 string position.
- 14 • Ketone as the last characters: The function checks if the SMILES string ends with =O,
15 indicating a ketone group at the end of the molecule. If the ketone is ring-connected
16 (e.g., c1cccc1=O), additional checks ensure the presence of a carbonyl group within the
17 ring. Similarly, the ketone group may appear at the end of the SMILES notation, but in
18 the molecular structure, it could be part of a cyclic or internal configuration. The
19 detection logic is based on chemical connectivity, not linear SMILES order.
- 20 • Ketone in the middle of the SMILES: The function searches for ketone groups within
21 the middle of the SMILES string using patterns such as C(=O), representing a carbonyl
22 group between two carbons. Also, the function carefully checks if the ketone is
23 branched, ensuring accurate identification of the ketone group in the middle. As an
24 example, in the SMILES string CC(=O)C, the middle ketone group C(=O) is detected
25 between two carbon atoms.
- 26 • Handling parentheses: The function accounts for complex SMILES structures that
27 contain nested parentheses. It ensures that ketone groups within branches (e.g.,
28 ...C(C(C...)C(=O)...)... in CC(C(CCC)=O)CC) are properly identified by finding the
29 matching opening and closing parentheses.

30 *2.1.10. Open-chain, alicyclic, and aromatic ether groups*



1 Three functions `open_chain_ether(s)`, `alicyclic_ether(s)`, and `aromatic_ether(s)` are
2 defined to distinguish between ethers present in the SMILES notation (open-chain, alicyclic,
3 and aromatic ethers). Below is a breakdown of the major components and logical flow within
4 this algorithm.

5 A key part of the function's operation involves detecting in-ring ethers. The function starts by
6 searching for the sequence "OC" within the SMILES string, which is the primary sequence of
7 the ether functional group. The other carbon atom bonded to the oxygen atom of the "OC"
8 sequence (e.g., `COC` or `C(OC...)`) is then searched to identify if the ether function is alicyclic or
9 open-chain. For each occurrence, the algorithm checks if the sequence is part of a non-aromatic
10 ring by comparing its position relative to the numerical ring indicators. This is done by locating
11 the positions of ring closure numbers (e.g., `C1...C1`) relative to the ether group. If the ether
12 group lies outside the ring closure points (e.g., `C1CCCC1COC`), it is classified as an open-chain
13 ether. On the other hand, if the ether group is found within the two identical ring numbers (e.g.,
14 `C1...COC...C1`), the algorithm hesitates if the ether is part of the ring or an open-chain type.
15 Therefore, the function employs a nested structure-parsing approach to clarify this issue. It uses
16 parentheses to detect branching points or nested structures within the molecule. The algorithm
17 carefully traces the boundaries of rings and other nested structures simultaneously. If the ether
18 group is embedded within parentheses and surrounded by ring numbers (e.g., pattern
19 `C1...(...COC...)...C1` in SMILES `C1CC(COC)CC1`), it is open-chain ether, while if it
20 follows patterns, such as `...C1...(...COC...C1...)...`, it is alicyclic (e.g., SMILES
21 `C1C(COCC1C)CC`). On the other hand, the function also checks if the ether is located near the
22 start or end of a ring closure, which would indicate that the ether is alicyclic (e.g., pattern
23 `C1...C(...O1)`), otherwise, it would be an open-chain ether. In compounds containing
24 multiple rings, the function iterates through each ring. It systematically searches for ether
25 groups within and around each ring, ensuring that all possible locations are checked for ether
26 group presence. If the ether group is attached to a non-aromatic ring, special handling is
27 performed to ensure accurate detection.

28 On the other hand, for aromatic ethers, the model looks for occurrences of 'Oc' within the
29 SMILES string. The logic checks different cases based on what precedes 'Oc':

- 30
- Simple alkyl group ('C' before 'Oc'): This detects linear aromatic ethers where the
31 oxygen is connected directly to an alkyl group (e.g., `COc1ccccc1`).



- 1 • Nested structures with parentheses: The model handles cases where the ether group is
2 part of a complex structure. If the ether group is surrounded by parentheses, the model
3 traces the original alkyl chain, ensuring it connects to an aromatic ring (e.g., patterns
4 C(Oc and C(...)Oc in SMILES CC(Oc1cccc1)CCC and CC(C(O)=O)Oc1cccc1).
- 5 • In case the ether group is attached to one aromatic and one non-aromatic ring, the model
6 iterates over possible ring numbers to find ether groups (e.g., 'c1CCCC1Oc1cccc1'
7 where 'c1' is part of a ring and 'Oc1' is the ether group). The same logic is applied to
8 nested structures with parentheses, ensuring that the correct connection between alkyl
9 and aromatic groups is maintained.
- 10 • The model repeats the search, but this time looking for 'OC', where the aromatic group
11 comes first (e.g., 'c1cccc1OC' where 'c1' indicates an aromatic ring and 'OC' is the
12 ether). Similar checks are performed for direct aromatic ether group connection
13 ('c1OC'), complex nested structures (e.g., c1(OC)ccc1 and c1cc(OC)cc1), and
14 parentheses-based structures (e.g., c1cc(ccc1)OC).

15 2.1.11. Ester and nitroester groups

16 The developed function, `ester_group(s)`, identifies and quantifies ester groups within a given
17 SMILES string. The ester group is characterized by the bonding of a carbonyl group (C=O) to
18 an alkoxy group (O-R). This pattern is typically (not always) represented in SMILES as
19 C(=O)O, and variations in its placement within the molecule must be accounted for. For
20 instance, when analyzing the ester group, if the alkoxy group (-OR) shows up in the SMILES
21 first and then the carbonyl group (C=O), common patterns, such as
22 ...C(...)(...)OC(=O)C... (e.g., in CC(CO)(C)OC(=O)CC with bolded characters) and
23 ...C(...)(OC(=O)C...)... (e.g., in CC(CO)(OC(=O)CC)C with bolded characters) are
24 identified. Conversely, if the SMILES notation proceeds from the acid-side carbon (i.e., carbon
25 attached directly to the carbonyl group in the ester bond -C(O)OR), ester groups may be
26 represented as ...C(...)(...)C(=O)OC... (e.g., in CC(CO)(C)C(=O)OCC with bolded
27 characters) or ...C(...)(C(=O)OC...)... (e.g., in CC(CO)(C(=O)OCC)C with bolded
28 characters). These patterns ensure that ester groups are recognized irrespective of their position
29 within a molecule.

30 The above-mentioned patterns with several other ones related to ester group may be determined
31 whether at the start, middle, or end of the SMILES string:



- 1 • If the ester group is at the beginning, the function looks for patterns such as
2 O=C(C...)OC... and O=C(OC...)C (e.g., in SMILES O=C(CCO)OCCC and
3 O=C(OCCC)CCO), where the group starts with the double-bond oxygen atom.
- 4 • For esters embedded within the middle of a SMILES string, the function searches for
5 the ester signatures such as C(=O)O and OC(=O), confirming that the carbonyl group is
6 bonded to the oxygen atom of the alkoxy group and followed by a suitable molecular
7 fragment. For instance, in a SMILES notation like CCC(=O)OCC and CCOC(=O)CC, the
8 ester group is correctly identified as part of the main chain.
- 9 • If an ester group is located at the end of the SMILES string, the function identifies
10 patterns such as)=O sequence as an indicator of the terminal carbonyl group, followed
11 by the appropriate bonding structure. This ensures that molecules like CC(OCC)=O are
12 correctly parsed with the ester group assigned to the end of the chain.
- 13 • In cases of aromatic esters, where alternating single and double bonds are common, the
14 function adapts the detection logic to properly handle aromaticity. For example, it
15 correctly identifies the ester group in a molecule like O=C(OCCC)c1ccccc1, where the
16 ester is attached to an aromatic benzene ring.
- 17 • The algorithm is designed to detect esters even in highly branched molecules. For
18 instance, in a molecule like CCC(C(=O)OCCC)(C)CC, the ester group is part of a
19 branching structure, and the function ensures it is accurately parsed by considering the
20 nested arrangement of atoms.

21 On the other hand, the model locates the starting and ending positions of the acid-side branch
22 for the ester group and searches for the presence of a nitro group (N(=O)=O) in the branch. If
23 a nitro group is detected, the compound is classified as a nitroester, and the
24 `nitroester_number` is incremented. In cases where no nitro group is present, the compound
25 is classified as a regular ester, and the `ester_number` is incremented accordingly.

26 2.1.12. C=CC=O in non-aromatic rings

27 The function `nonaromatic_CCCO(s)` is designed to quantify C=C-C=O substructures in a
28 molecular SMILES notation. This section explains how the function works.

29 Since the substructure of interest occurs within rings, the function first checks whether the
30 molecule contains any cyclic structures. The function uses `find_cycle_number(s)` to
31 determine if any rings exist in the molecular structure. If no rings are found, the function



1 terminates early. When rings are detected, the function proceeds to locate their positions within
2 the SMILES string.

3 The next step is to detect the C=C–C=O group within the identified rings. This involves
4 scanning the part of the SMILES string that represents each ring and checking for the specific
5 pattern of atoms and bonds. Several important aspects are considered during this analysis:

- 6 • Pattern search: Within the extracted portion, the function looks for patterns that match
7 the C=C–C=O group. For example, the presence of C=O at the start of a ring followed
8 by a conjugated double bond (C=C) within the ring (e.g., pattern **O=C1...C(...)=C1**
9 **bolded in SMILES O=C1CCC(OC)=C1C**), or variations where the C=C and C=O groups
10 may be spaced by additional atoms, or where they may appear in different locations
11 within the ring (e.g., pattern **...1...C(=O)C(...)=C...1...** **bolded in SMILES**
12 **O=C1CC(=O)C(CC)=CC1C**).
- 13 • Handling complex ring structures and nested rings: The function accounts for these
14 complexities by carefully navigating through the parentheses and ensuring that the
15 entire cyclic structure is examined for the C=C–C=O group (e.g., pattern
16 **...C(=C(C(...1...)=O)...)...** **bolded in SMILES CC1C(=C(C(CC1C)=O)C)CC**).

17 2.1.13. Nitrate group

18 The `nitrate_number(s)` function identifies and quantifies nitrate functional groups within the
19 provided SMILES notation. This function examines specific configurations characteristic of
20 nitrate, including the standard nitrate structure **ON(=O)=O**, the quaternary form **O[N+](=O)[O-]**,
21 the N-nitro structure **O=N(=O)O**, and variations like **O=[N+](O-)** and **[O-][N+](=O)O**,
22 which collectively reflect diverse bonding scenarios in nitrate chemistry. The function
23 calculates the count of each of these configurations and aggregates them to derive the total
24 number of nitrate groups present.

25 2.1.14. Nitro group

26 The identification of nitro groups is executed through the `nitro_group(s)` function, which
27 analyzes the SMILES notation to quantify nitro functional groups. This function employs a
28 series of search operations to locate specific nitro structures, notably **N(=O)=O**, **O=N(=O)**, and
29 various ionic forms such as **[N+](=O)[O-]**, **O=[N+](O-)**, and **[O-][N+](=O)**. Each search
30 utilizes a loop that not only finds occurrences of these structures but also ensures that adjacent



1 atoms do not disrupt the nitro configuration—specifically, it checks that there is no oxygen
2 atom directly connected to the nitro group, which would suggest an alternative bonding
3 scenario (i.e., nitrate group). As the final filtration step, the algorithm verifies whether the
4 identified nitro group is not part of an aromatic ring (i.e., benzene). If this condition is met, it
5 then checks for the presence of a hydroxyl group in the ring according to section 2.6. If no
6 hydroxyl group is found, the nitro group is counted independently. However, if a hydroxyl
7 group is present, the compound is categorized as a nitrophenol group. A detailed explanation
8 of this functional group is provided in a subsequent subsection.

9 *2.1.15. Nitrophenol*

10 The function `nitrophenol_group(s)` is designed to identify and count the number of
11 nitrophenol groups. After identifying aromatic hydroxyl group, the function then checks for
12 the presence of a nitro group in the same aromatic ring within the SMILES string:

- 13 • As the last character sequence: It inspects whether the string ends with the motif
14 N(=O)=O, which signifies a nitrophenol. The function checks if it is connected to the
15 aromatic cyclic structure by examining the characters before it. For example, if the
16 character is an aromatic ring identifier, the counter is incremented accordingly (e.g.,
17 c1cc(O)ccc1N(=O)=O).
- 18 • As the first character sequence: The function checks if the string starts with the pattern
19 O=N(=O)c1, indicating a nitrophenol positioned within the SMILES string (e.g.,
20 O=N(=O)c1cc(O)ccc1). If this pattern is detected, the counter is incremented.
- 21 • As middle character sequences: A loop is employed to search for occurrences of
22 patterns such as c(N(=O)=O), which indicates that nitrophenol is positioned within the
23 structure (e.g., c1cc(N(=O)=O)cOcc1O). Each time this pattern is found, the counter is
24 increased. The function also iterates through the previously established list of rings to
25 search for patterns of the form j(N(=O)=O) (where *j* is an aromatic carbon identifier).
26 Whenever a match is found, the counter is incremented (e.g., c1(N(=O)=O)cc(O)ccc1).

27 *2.1.16. Acylperoxynitrate group*

28 The function `carbonylperoxynitrate_group(s)` is designed to detect the presence of
29 acylperoxynitrate groups within a given SMILES string *s*. This function systematically counts
30 occurrences of several distinct structural motifs characteristic of acylperoxynitrates, including



1 OON(=O)=O, OO[N+](=O)[O-], O=N(=O)OO, O=[N+]([O-])OO, and [O-][N+](=O)OO. The
2 function aggregates the counts of these motifs into a single variable, `carbonylperoxynitrate`,
3 which represents the total number of acylperoxynitrate groups identified.

4 *2.1.17. Peroxide group*

5 The function `peroxide_group(s)` is engineered to identify and quantify peroxide groups. The
6 function first determines the number of cyclic structures within the SMILES string by calling
7 the `find_cycle_number(s)` function. If cyclic structures are found, a list of ring identifiers,
8 comprising both carbon (c) and aromatic (c) rings, is generated to facilitate later checks.
9 Subsequently, the function employs a loop to search for occurrences of the `ooc` motif, which
10 signifies the presence of peroxy groups. Each iteration of the loop calls `s.find('ooc', ...)`
11 to locate the next occurrence of the motif. Multiple conditions are assessed to ensure that the
12 identified `ooc` is correctly positioned relative to other atoms or rings:

- 13 • Adjacent carbon or aromatic carbon atoms: If the character preceding `ooc` is C, Cx or cX
14 (x=1,2,... is cyclic index), indicating that `ooc` is bonded to a carbon atom, the
15 `peroxide_number` is incremented (e.g., CCCCOCC, C1CCCC1OCC and Cc1cccc1OCC).
- 16 • Branching structures: If the character before `ooc` is a parenthesis, the function retrieves
17 the index of the last corresponding parenthesis before `ooc` and verifies that the atom
18 preceding this parenthesis is a carbon atom or an aromatic ring. If so, the counter is
19 increased. For example, when the character before `ooc` is a closing parenthesis, the
20 function checks whether the `ooc` is connected to a carbon atom in a similar manner as
21 described previously. This involves searching back to the last opening parenthesis and
22 ensuring the atom connected to that parenthesis is a carbon atom or part of a cyclic
23 structure (e.g., pattern **...c(...)ooc...** bolded in SMILES ccc(cc)ooccc).

24 These checks comprehensively ensure that only valid peroxide groups are counted, accounting
25 for the complex connectivity possible within SMILES representations. The function ultimately
26 returns the `peroxy_number`, providing a quantitative measure of peroxide groups within the
27 molecular structure.

28 *2.1.18. Aromatic amine group*



1 Following the identification of aromatic rings, the function `aromatic_amine_group(s)`
2 locates nitrogen atoms (N) within the SMILES string and determines their bonding to aromatic
3 carbons.

- 4 • Direct bond to aromatic carbon (`cN`): If nitrogen (N) is found immediately following an
5 aromatic carbon (c), it is counted as part of an aromatic amine group. To ensure the
6 nitrogen atom belongs to an aromatic amine and not a nitro group (`-NO2`), the model
7 incorporates an additional condition.
- 8 • Nitrogen in numbered rings (`c1N`, `c2N`, etc.): If nitrogen is attached to a carbon in a
9 numbered ring, such as `c1N`, the nitrogen is identified as part of an aromatic amine
10 group.
- 11 • Parenthetical structures (`c(...)N`): In cases where nitrogen is attached within
12 parentheses following a cyclic group, the function checks if the nitrogen is part of an
13 aromatic ring by verifying the bonding pattern of the cyclic carbon to the nitrogen.
14 Parentheses in SMILES represent branching, and this function ensures that any
15 branching nitrogen groups attached to aromatic carbons are also detected.

16 The function iterates through the SMILES string to ensure that all occurrences of nitrogen
17 atoms are evaluated. The bonding pattern of each nitrogen atom is assessed against the aromatic
18 rings identified in the first step. If the nitrogen atom is confirmed to be attached to an aromatic
19 ring and not part of a nitro group, it is counted as part of an aromatic amine group. The total
20 count of such groups is stored in the variable `aromatic_amine_number` and returned as the
21 output of the function.

22 *2.1.19. Primary, secondary, and tertiary amide and amine groups*

23 The model is designed to detect and count primary, secondary, and tertiary amide and amine
24 groups in a molecule represented by a SMILES string. A primary amide has the functional
25 group structure `-C(=O)NH2`, and the model identifies both simple and branched forms of this
26 group. To differentiate between primary amides and primary amines, the model specifically
27 excludes patterns without a double-bonded oxygen, i.e., `-C(...)NH2` where ‘...’ is not `=O`, thereby
28 ensuring correct identification of amide groups versus amine counterparts.

- 29 • The function first identifies primary amide groups located at the beginning of the
30 SMILES string by searching for patterns such as `O=C (N)` or `NC (=O)`, as seen in SMILES



1 O=C(N)CCCC and NC(=O)CCCC). It also detects primary amides at the end of the SMILES
2 string by tracing patterns like C(=O)N and C(N)=O, as in CCCC(=O)N and CCCC(N)=O.
3 To identify primary amine groups at the beginning or end of the strings, the function
4 checks for similar patterns but without a carbonyl group (=O). In other words, if a
5 nitrogen is bonded to a carbon that does not carry a =O, it is interpreted as an amine
6 rather than an amide. For example, C(N)CCCC, NC(C)CC, CC(C)CN, and CCCCN are
7 recognized as containing primary amine groups.
8 • The model also identifies primary amides within branches or internal positions in the
9 molecule by searching for specific patterns, such as C(=O)N and C(N)=O, as seen in
10 the SMILES strings CC(C(=O)N)CC and CC(C(N)=O)CC. To identify primary amine
11 groups in similar positions, the model checks for the absence of the double-bonded
12 oxygen (=O) on the carbon adjacent to the nitrogen. This results in patterns like CN in
13 CC(CN)CC and CC(CN)CC.

14 For secondary and tertiary amides, the model similarly searches for patterns where the nitrogen
15 atom is bonded to two or three carbon atoms, respectively. Secondary amides have the structure
16 -C(=O)NR, where R represents an alkyl group starting with a carbon atom attached to the
17 nitrogen atom, and tertiary amides have the structure -C(=O)N(R)R', where both R and R' are
18 alkyl groups starting with carbon atoms attached to the nitrogen atom. The model identifies
19 these structures by recognizing the presence of additional carbon attachments to the nitrogen
20 atom. For secondary amides, the function searches for patterns that indicate the nitrogen is
21 bonded to one additional carbon group, distinguishing them from primary amides by checking
22 for two single bonds to nitrogen, along with the carbonyl group. Similarly, for tertiary amides,
23 the function detects two alkyl groups attached to the nitrogen atom in addition to the carbonyl
24 group. Once these amide patterns are identified, the model applies the same exclusion method
25 for the double-bonded oxygen, converting these amides into their corresponding secondary and
26 tertiary amines. For secondary amines, the nitrogen is attached to two carbon atoms, and for
27 tertiary amines, the nitrogen is bonded to three carbon atoms. This method ensures that the
28 correct amide or amine group is identified and classified, whether it is primary, secondary, or
29 tertiary, based on the number of carbon attachments to the nitrogen atom.

30 Finally, the model counts the number of carbon atoms in the secondary and tertiary amides that
31 are not part of the R and R' groups in the structure. This count is considered as the number of
32 carbons on the acid side of the amides. For primary amides, since there are no additional alkyl



1 groups attached to the nitrogen atom, all carbon atoms in the structure are considered to be on
2 the acid side of the amide. This ensures accurate categorization and counting of carbon atoms
3 associated with the amide's acid side, contributing to the overall structural analysis of the
4 molecule.

5 **2.2.Saturation vapor pressure calculation**

6 After identifying the functional groups, the detection functions return integer values
7 representing the occurrence of each functional group, which are stored in an array. These values
8 are critical input parameters for the subsequent saturation vapor pressure calculation, as the
9 saturation vapor pressure is quantified as the sum of functional group contributions in the
10 SIMPOL method. Thus, each SMILES string (e.g., for compound i) is processed to retrieve the
11 saturation vapor pressures according to the SIMPOL method. The SIMPOL method defines
12 how each functional group contributes to the saturation vapor pressure—at different
13 temperatures. In the model:

- 14 i. Matrix B is read from a pre-defined text file containing coefficients of $B_{k,1}$, $B_{k,2}$, $B_{k,3}$,
15 and $B_{k,4}$ for each functional group k, according to Table 5 of (Pankow and Asher 2008).
16 A select functional group is assigned a value for contribution to saturation vapor
17 pressure in the i^{th} SMILES string (such as hydroxyl, aldehyde, and ketone groups, etc.).
- 18 ii. Then, $b_k(T)$ and $P_{L,i}^0(T)$ are calculated for any given temperature according to Equations
19 1 and 2. The total liquid (saturation) vapor pressure, $P_{L,i}^0$ (atm), is calculated as the sum
20 of all functional group contributions.

$$b_k(T) = \frac{B_{k,1}}{T} + B_{k,2} + B_{k,3} T + B_{k,4} \ln T \quad (1)$$

$$\log_{10} P_{L,i}^0(T) = \sum_k v_{k,i} b_k(T) \quad k = 0, 1, 2, \dots \quad (2)$$

21 where $v_{k,i}$ is the number of groups of type k, $b_k(T)$ is the contribution to $\log_{10} P_{L,i}^0(T)$ by
22 each group of type k, and T is the temperature. Also, 0 and L show the reference and
23 Liquid.



- 1 iii. To fit the saturation vapor pressure data to the Antoine equation (Equation 3) to enable
2 further use of the saturation vapor pressure values in different applications, the
3 saturation vapor pressure is calculated at 1000 temperature points across a wide range
4 of temperatures from 220 K to 450 K according to Equation 2. The obtained data are
5 then used in a non-linear least squares fitting procedure, which minimizes the difference
6 between the data and the saturation vapor pressure values predicted by the Antoine
7 equation. The Antoine equation parameters (i.e., A, B, and C) are then obtained for each
8 compound.

$$\log P_{sat} = A - B/(T + C) \quad (3)$$

- 9 iv. After obtaining the Antoine equation parameters, the vaporization enthalpy relationship
10 can be derived using the Clausius-Clapeyron equation:

$$\frac{d \log P_{sat}(T)}{d(\frac{1}{T})} = -\frac{\Delta H_{vap}(T)}{2.303R} \quad (4)$$

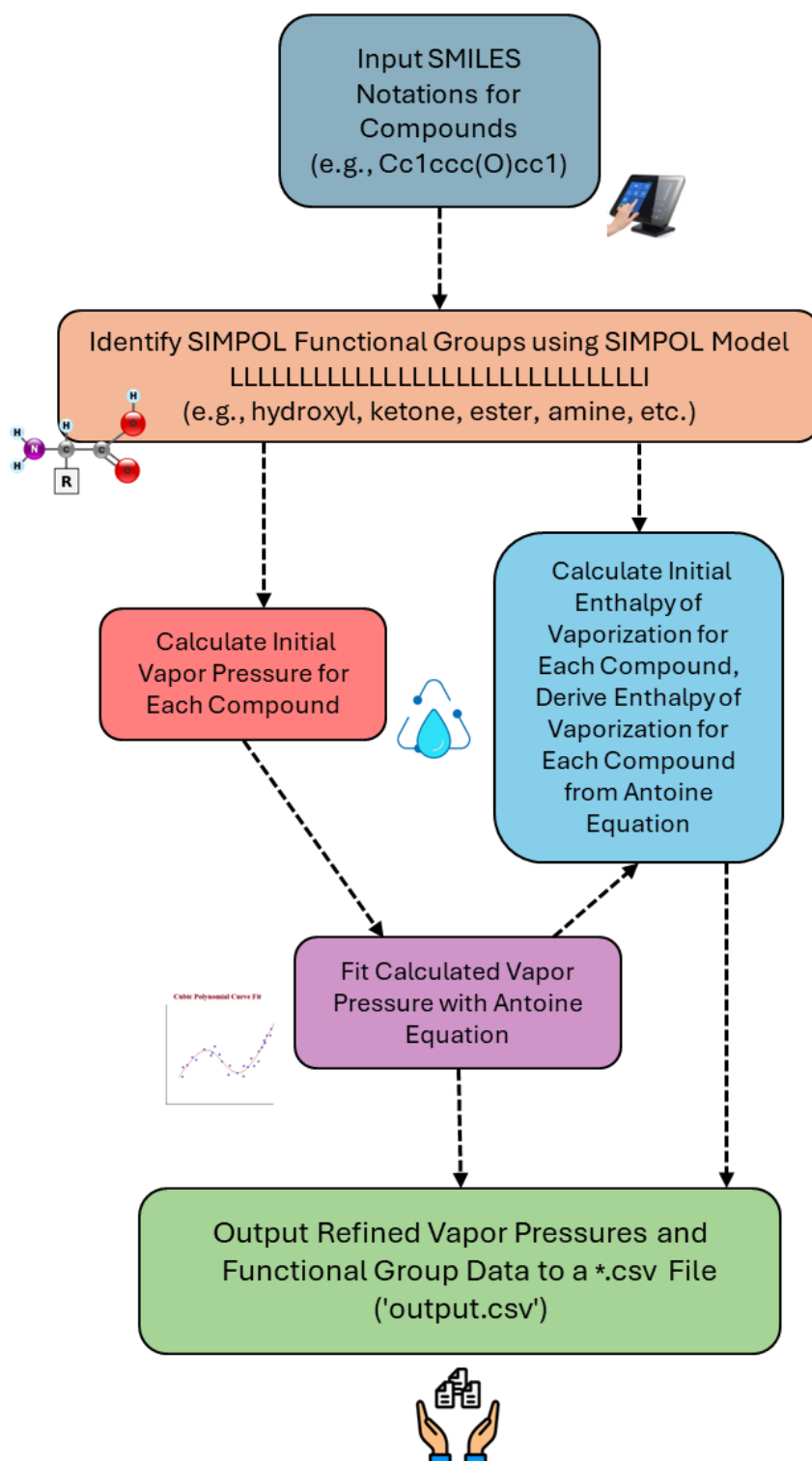
$$\Delta H_{vap}(T) = -2.303R \left(\frac{d \log P_{sat}(T)}{d(\frac{1}{T})} \right) \quad (5)$$

- 11 Here, $\Delta H_{vap}(T)$ is the temperature-dependent enthalpy of vaporization, and R is the universal
12 gas constant. This expression relates the slope of the logarithm of the saturation pressure with
13 respect to the inverse of temperature to the enthalpy of vaporization. The Antoine equation
14 provides a framework to calculate saturation vapor pressure at any given temperature, and this
15 relationship extends the utility of the model by allowing the determination of thermodynamic
16 quantities such as vaporization enthalpy.

- 17 v. The output is written to a CSV file named by the user (e.g., output.csv), where each line
18 corresponds to the i^{th} compound and its associated data, including its SMILES string,
19 the count of each functional group in its structure, its saturation vapor pressure at 300
20 K, and its fitted Antoine equation parameters.

21 Figure 3 illustrates a flowchart of the VaPOrS from input to output.

22





1 **Figure 3.** Flowchart of the automated process of VaPOrS, beginning with SMILES notation input, followed by
2 the identification of functional groups, initial saturation vapor pressure and enthalpy of vaporization
3 calculations, their fitting to the Antoine equation, and final output.

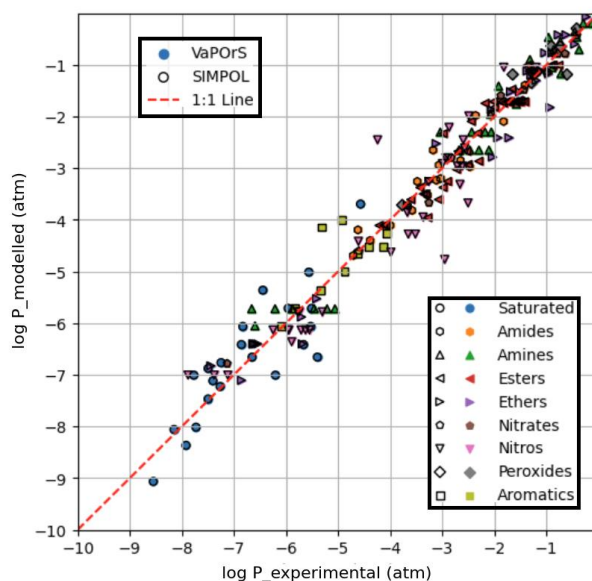
4

5 **3. Results and discussion**

6 **3.1. Saturation vapor pressure**

7 To evaluate the accuracy of the developed automated VaPOrS in predicting the saturation vapor
8 pressures of organic compounds using the SIMPOL method, the model was applied to a subset
9 of the original dataset used to develop the SIMPOL parameterization. Saturation vapor
10 pressures for 224 organic compounds were calculated using VaPOrS and compared against
11 experimental data to assess the accuracy of the implementation and ensure consistency with
12 established results.

13 Figure 4 presents a comparative analysis between the saturation vapor pressures computed by
14 the VaPOrS and those reported by the SIMPOL method and measurement at a specific
15 temperature (i.e., 333.15 K). The x-axis represents the experimental saturation vapor pressures,
16 while the y-axis represents the numerical values calculated by the VaPOrS (filled symbols with
17 no edge color) and the SIMPOL method (black-edged hollow symbols). Due to a complete
18 overlap in the data points, the black-edged symbols from the SIMPOL method can be seen over
19 the filled symbols from VaPOrS. A diagonal line is shown in the figure, indicating the ideal
20 correlation where the computed values would perfectly match the measured saturation vapor
21 pressures. Data points positioned close to this diagonal demonstrate a high level of agreement
22 between the two approaches.



1

2 **Figure 4.** Comparison of saturation vapor pressure values calculated by the VaPORs (filled symbols with no
3 edge color) and the SIMPOL method (symbols with no face color and black edges) against measured saturation
4 vapor pressures at 333.15 K. The diagonal line represents the ideal 1:1 correlation, and the proximity of data
5 points to this line indicates the accuracy of both methods in predicting saturation vapor pressures. The overlap in
6 symbols is visible due to the black-edged SIMPOL markers covering the filled VaPORs symbols.

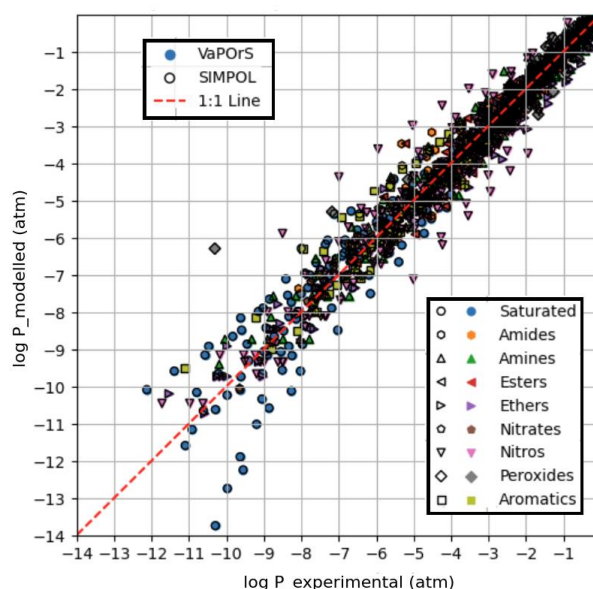
7

8 Figure 5 illustrates the saturation vapor pressure results obtained using the VaPORs model, the
9 SIMPOL method, and experimental measurements at six different temperatures: 273.15 K,
10 293.15 K, 310.15 K, 333.15 K, 353.15 K, and 373.15 K. Similar to Figure 2, the x-axis
11 represents the measured saturation vapor pressures, while the y-axis displays the values
12 calculated by VaPORs (filled symbols with no edge color) and the SIMPOL method (black-
13 edged hollow symbols). This figure includes a larger dataset, allowing for a more
14 comprehensive assessment of model performance across a range of temperatures. Many data
15 points are clustered close to this line, further confirming the effectiveness of the VaPORs model
16 in predicting saturation vapor pressures for various compounds across different temperatures
17 as well. Some deviations from the line are observed e.g., the experimental values for the nitro
18 and saturated compound saturation vapor pressures and T -dependencies have large
19 uncertainties as seen in Figures 4 and 5. These discrepancies, which were present in the original
20 dataset, do not undermine the overall trend, which demonstrates strong agreement among the



1 three methods and suggests that the VaPOrS code is reliable and robust across a wider
2 temperature range. Most importantly, we see complete agreement in the output of VaPOrS and
3 SIMPOL.

4



5

6 **Figure 5.** Comparison of saturation vapor pressures obtained from the VaPOrS model, the SIMPOL method,
7 and experimental measurements at six different temperatures (273.15 K, 293.15 K, 310.15 K, 333.15 K, 353.15
8 K, and 373.15 K). The x-axis represents the measured saturation vapor pressures, while the y-axis shows the
9 values calculated by VaPOrS (filled symbols) and the SIMPOL method (hollow symbols with black edges). The
10 diagonal line indicates the ideal correlation, with points near the line demonstrating good agreement between the
11 methods.

12

13 3.2. Enthalpy of vaporization

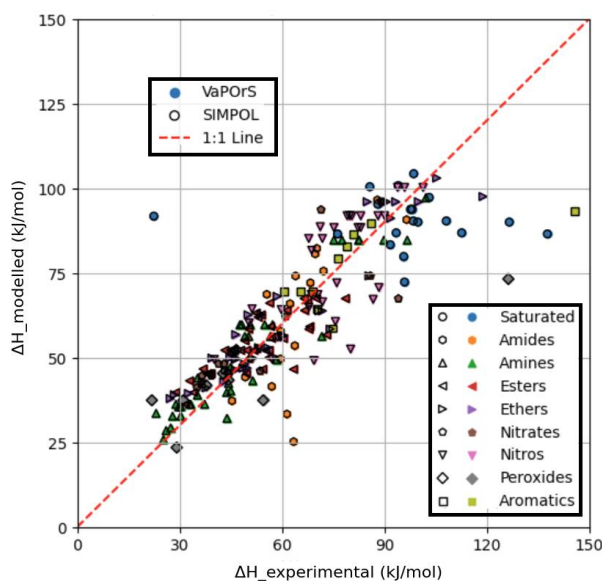
14 Figure 6 illustrates the relationship between the experimental vaporization enthalpy values and
15 those calculated by both the VaPOrS and the SIMPOL method at 333.15 K. The x-axis
16 represents the experimentally measured values, while the y-axis displays the calculated values.
17 The results from the VaPOrS are represented by filled symbols without edge color, indicating
18 a direct prediction from the present model. In contrast, the SIMPOL method results are depicted



1 as symbols with no face color and distinct black edges. The presence of overlapping symbols
2 highlights instances where the calculated values from SIMPOL cover those from VaPOrS.
3 According to Figures 4 to 6, many data points are clustered near the diagonal line,
4 demonstrating the effectiveness of the model in predicting saturation vapor pressure and
5 enthalpy of vaporization for several compounds. However, there are some outliers where the
6 calculated values deviate from the experimental counterparts. These discrepancies could be
7 attributed to the structural complexity of certain compounds making intramolecular
8 interactions important and not amenable to simple group additivity predictions, which is also a
9 known limitation in the SIMPOL method itself, or potentially, they could point out issues in
10 the original experimental measurements.

11 The ability to estimate vaporization enthalpy alongside saturation vapor pressure enables a
12 more comprehensive analysis of the compounds' behavior, especially for estimating the
13 tendency to form aerosol, in various atmospheric and environmental conditions.

14



15

16 **Figure 6.** Comparison of enthalpy of vaporization values calculated by the VaPOrS (filled symbols without
17 edge color) and the SIMPOL method (symbols with no face color and black edges) against experimental
18 measurements at 333.15 K. The diagonal line indicates the ideal 1:1 correlation, showcasing the accuracy of the

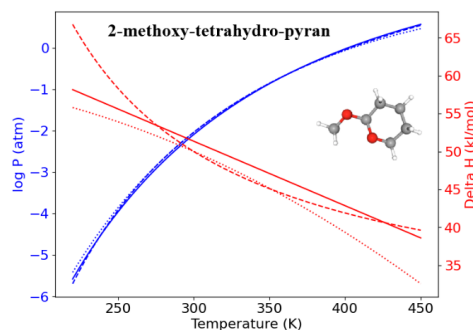
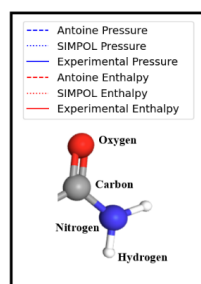


1 models. Overlapping symbols are observed, with black-edged SIMPOL markers obscuring the filled VaPOrS
2 symbols.

3 3.3. Antoine equation

4 The results of the fitting process for nine compounds, i.e., 2-methoxy-tetrahydro-pyran,
5 decanedioic acid, methyl-benzoate, phenylamine, hexanamide, phenylmethyl-nitrate, 2-
6 methyl-6-nitrobenzoic acid, diethyl-peroxide, and 2-naphthol as representatives of ethers,
7 saturated, esters, amines, amides, nitrates, nitro-compounds, peroxides, and aromatics,
8 respectively, are visualized in Figure 7. The figure illustrates the temperature-dependent
9 behavior of both pressure and enthalpy of vaporization for Antoine and SIMPOL relationships
10 generated by VaPOrS and compares them with experimental data across varying temperatures.
11 The left y-axis represents the logarithmic saturation vapor pressure in atmospheres, while the
12 right y-axis shows the enthalpy of vaporization in kJ/mol. This dual-axis representation enables
13 a direct visual comparison between pressure and enthalpy trends as the temperature increases.
14 The fitting results demonstrate a high degree of agreement between the Antoine and SIMPOL
15 curves for all compound classes, implying that the Antoine equation given by VaPOrS can be
16 applied effectively to estimate saturation vapor pressures with good accuracy across a broad
17 range of temperatures, enhancing the utility of the saturation vapor pressure data for various
18 applications.

19



20

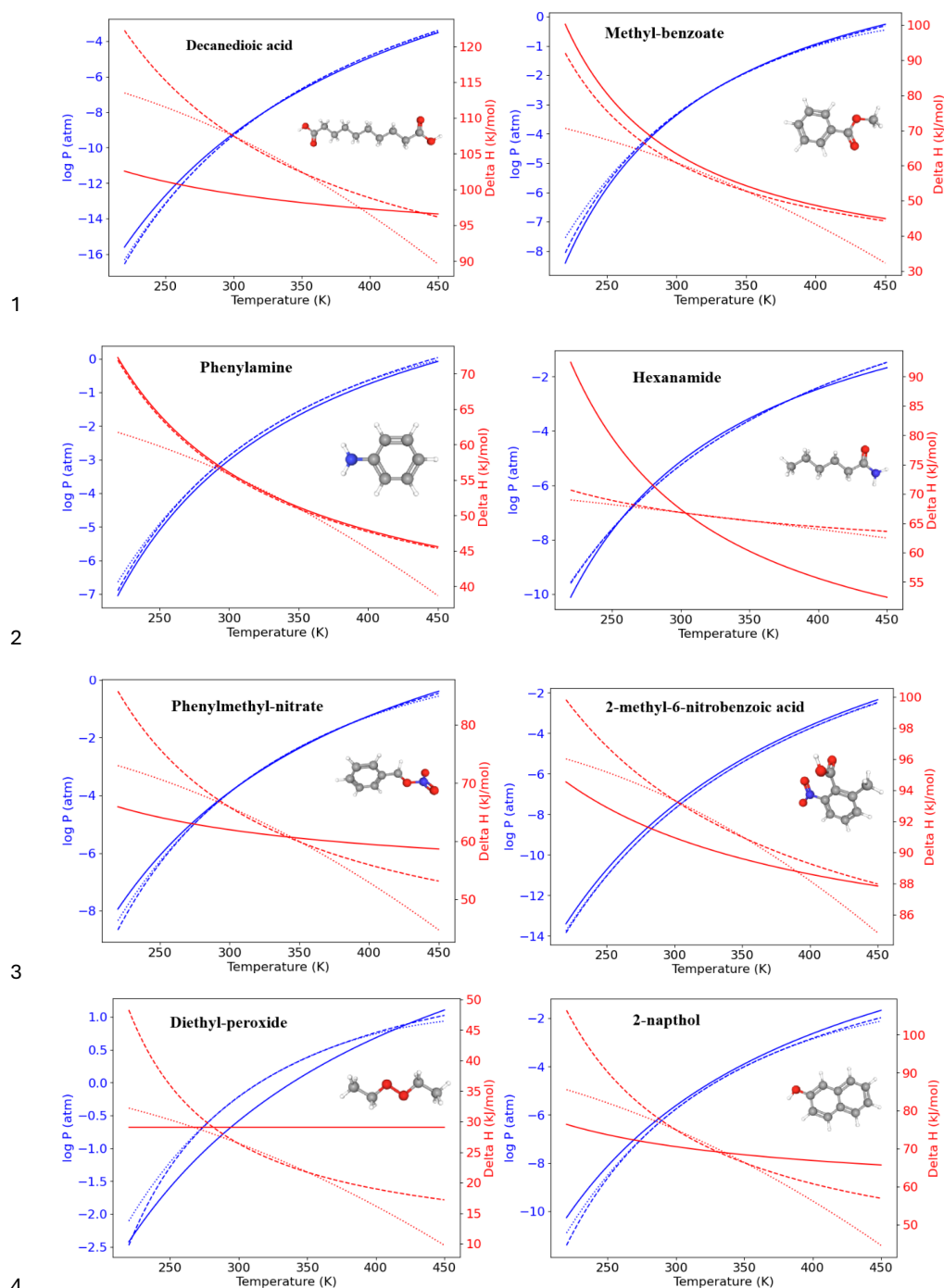


Figure 7. Temperature dependence of saturation vapor pressure and enthalpy of vaporization for nine representative organic compounds with nine distinct functional groups predicted by VaPORs using the Antoine



1 and SIMPOL equations. The left y-axis shows the logarithmic saturation vapor pressure (in atm), and the right
2 y-axis displays the enthalpy of vaporization (in kJ/mol). The results demonstrate the data generated by Antoine
3 and SIMPOL methods across the temperature range, with experimental data closely matching the theoretical
4 predictions.

5

6 **3.4.VaPOrS application**

7 Using VaPOrS, a detailed comparison was performed between the manual counting and
8 calculation of functional groups and saturation vapor pressures and the automated results
9 generated from the compounds' SMILES notations. The manual counting involved a systematic
10 review of each compound's molecular structure, visually identifying and recording the
11 functional groups, followed by calculating its saturation vapor pressure according to SIMPOL
12 group contributions. This was then cross-referenced with the automated results generated by
13 VaPOrS to ensure consistency. The procedure was performed in three steps described next.

14 **3.4.1. MCM data**

15 In the first step, a dataset of 126 primary VOCs sourced from the Master Chemical Mechanism
16 (MCM) database are evaluated. While the MCM database provides detailed chemical
17 mechanisms for atmospheric chemistry, its coverage of primary organic compounds is
18 relatively limited compared to the vast diversity of VOCs present in the atmosphere.
19 Nonetheless, it serves as a valuable resource for validation, given its detailed representation of
20 key compounds. Notably, the MCM database not only provides the structures of these
21 compounds but also includes their corresponding SMILES notation, facilitating an accurate
22 assessment of functional group presence through the VaPOrS method.

23 Table 1 presents a sample comparison between the manual and automated counts of the
24 functional groups for representative compounds from several categories in the MCM, including
25 Alcohols and Glycols, Aldehydes, Alkanes, Alkenes, Alkynes, Aromatics, Dialkenes, Esters,
26 Ethers and Glycol Ethers, Ketones, Monoterpenes and Sesquiterpenes, Organic Acids, and
27 Unclassified compounds. As shown, the results from both methods are in complete agreement,
28 with 0% discrepancy across all cases.

29



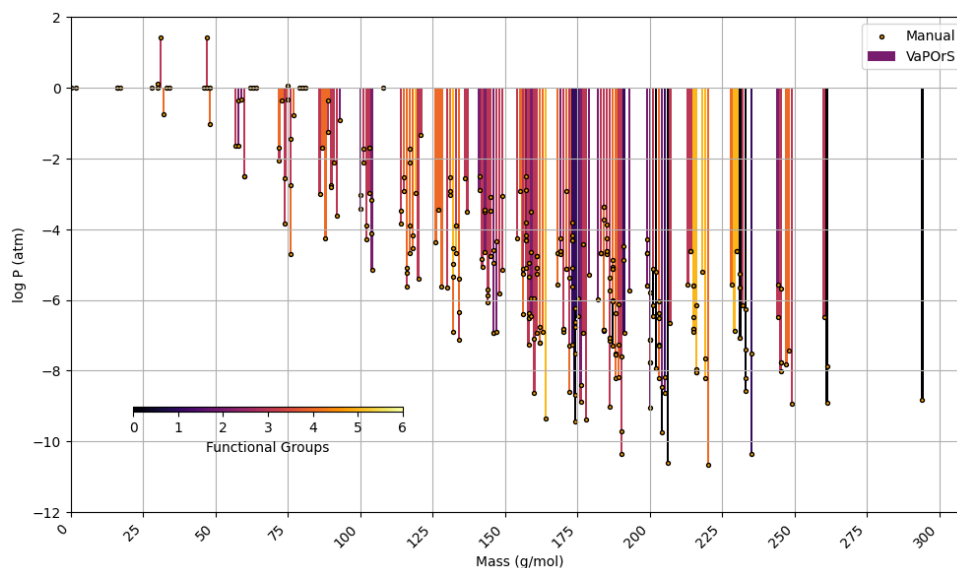
Table 1: Comparison of Manual and Automated Counts of Functional Groups for Representative Compounds across Various Categories in the MCM Database.

Category	Compounds	SMILES	Functionals	Consistency
Alcohols and Glycols	CYCLOHEXANOL	<chem>OC1CCCCC1</chem>	6 carbons, 1 nonaromatic ring, 1 hydroxyl	100%
Aldehydes	PROPENAL	<chem>C=CC=O</chem>	3 carbons, 1 C=C (non-aromatic), 1 aldehyde	100%
Alkanes	3-METHYLPENTANE	<chem>CCC(C)CC</chem>	6 carbons	100%
Alkenes	1-HEXENE	<chem>CCCCC=C</chem>	6 carbons, 1 C=C (non-aromatic)	100%
Alkynes	ETHYNE	<chem>C#C</chem>	2 carbons	100%
Aromatics	ETHYL BENZENE	<chem>CCc1ccccc1</chem>	8 carbons, 1 aromatic ring	100%
Dialkenes	1-3 BUTADIENE	<chem>C=CC=C</chem>	4 carbons, 2 C=C (non-aromatic)	100%
Esters	ETHYL ACETATE	<chem>CCOC(=O)C</chem>	4 carbons, 1 ester	100%
Ethers and Glycol Ethers	2-METHOXY ETHANOL	<chem>COCCO</chem>	3 carbons, 1 hydroxyl, 1 ether	100%
Ketones	CYCLOHEXANONE	<chem>O=C1CCCCC1</chem>	6 carbons, 1 nonaromatic ring, 1 ketone	100%
Monoterpenes and Sesquiterpenes	ALPHA-PINENE	<chem>CC1=CCC2CC1C2(C)C</chem>	10 carbons, 2 nonaromatic rings, 1 C=C (non-aromatic),	100%
Organic Acids	PROPANOIC ACID	<chem>CCC(=O)O</chem>	3 carbons, 1 carboxylic acid	100%
Unclassified	ETHYLENE OXIDE	<chem>O1CC1</chem>	2 carbons, 1 nonaromatic ring, 1 ether (alicyclic)	100%

In the second phase, alpha-pinene and benzene were selected as case studies to evaluate the species formed during their tropospheric degradation via gas-phase chemical processes, focusing on functional group occurrence and saturation vapor pressure. This analysis leveraged the detailed mechanism in the MCM to further validate the automated functional group detection system's accuracy in modeling atmospheric chemistry. For each species, the occurrences of functional groups were manually counted, and saturation vapor pressure was calculated using the SIMPOL method. Their SMILES notation was then input into the VaPOrS to automatically obtain functional group counts and saturation vapor pressures. The saturation vapor pressures obtained through both methods are compared in Figures 8 and 9 for alpha-pinene and benzene, respectively, where the y-axis represents the logarithmic saturation vapor pressure and the x-axis the molar mass of each species. Automated results are shown as color bars based on the number of detected functional groups, and their manual counterparts are displayed as points. The perfect alignment of points atop bars for each species indicates excellent agreement between both approaches. It is worth mentioning that C₆H₆N₂O₁₁ and CH₂O (i.e., NNCATECOOH and HCHO in the MCM) were recognized as the least and most volatile species in the benzene oxidation process. On the other hand, C₉H₁₆O₆ and CH₃O



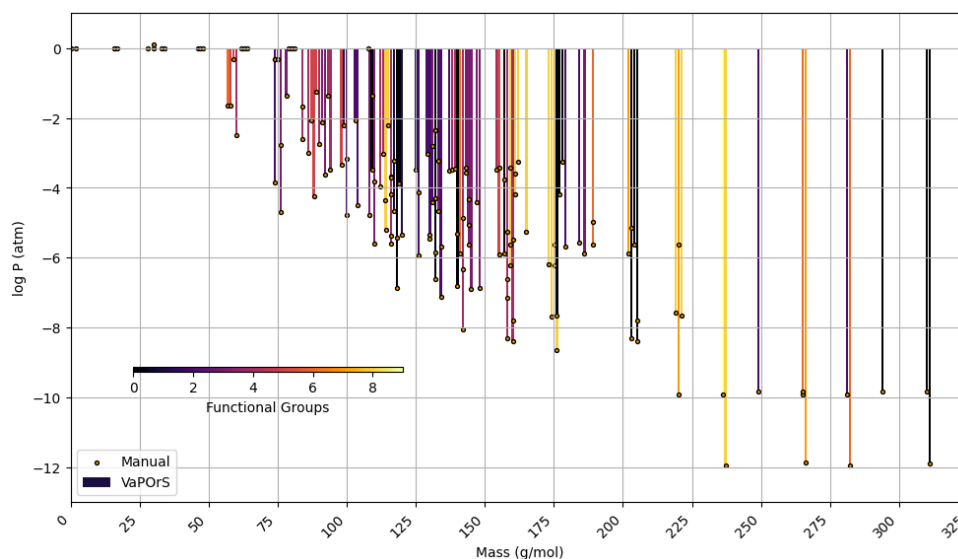
- 1 (i.e., C922OOH and CH3O in the MCM) were the least and most volatile species in the alpha-
- 2 pinene oxidation process.



- 3
- 4 **Figure 8.** Comparison of saturation vapor pressure results for tropospheric degradation species of alpha-pinene.
- 5 The figure presents the log saturation vapor pressure versus molar mass for species formed from the
- 6 tropospheric oxidation of alpha-pinene according to MCM. Bars show results from the automated VaPOrS code,
- 7 with colors based on detected functional group number involved in the chemical structure of species, and points
- 8 reflect manually calculated values. The alignment of points atop bars demonstrates the perfect consistency
- 9 between automated and manual calculations.

10

11



1

2 **Figure 9.** Comparison of saturation vapor pressure results for tropospheric degradation species of benzene. The
 3 figure illustrates the log saturation vapor pressure versus molar mass for species derived from benzene
 4 degradation according to MCM. The bars represent saturation vapor pressures calculated by the automated
 5 VaPOrS code, with colors based on detected functional group number involved in the chemical structure of
 6 species, while points indicate manually obtained values. The close alignment of points with bars highlights the
 7 accuracy of the automated method relative to manual calculations.

8

9 **3.4.2. autoAPRAM-fw data**

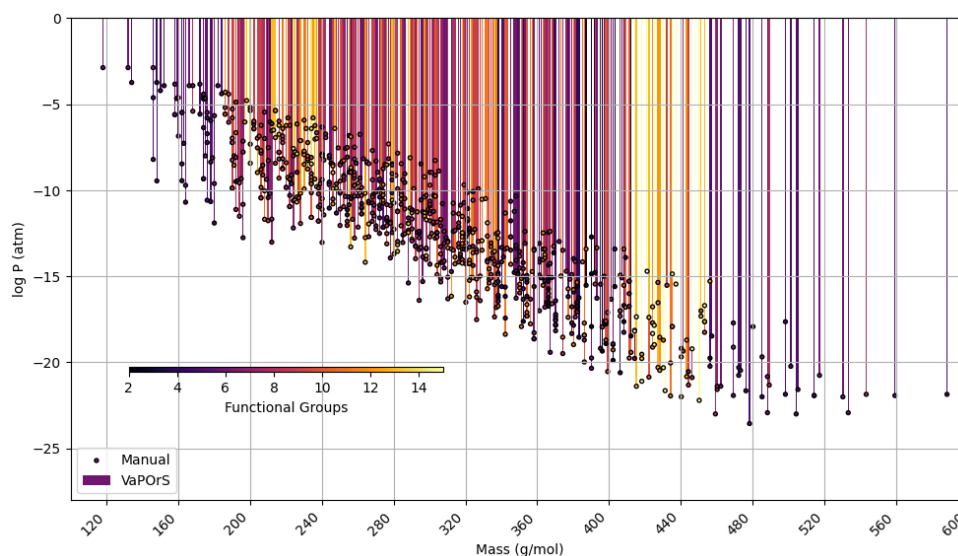
10 In the final stage of this analysis, the VaPOrS model was utilized to determine the saturation
 11 pressures of over 850 potential chemical species produced through the autoxidation of alkoxy
 12 and peroxy radicals, which emerge during benzene degradation. The initial radicals were
 13 defined by the MCM, with their respective saturation vapor pressures detailed in Figure 9.
 14 Conversely, the products of autoxidation were generated using the autoAPPRAM-fw tool
 15 (Pichelstorfer et al. 2024), with their potential structures represented by SMILES notation.
 16 Demonstrating its efficiency, VaPOrS analyzed all SMILES entries within a single second,
 17 accurately counting the required functional groups and calculating corresponding saturation
 18 vapor pressures. The results of these predictions are illustrated in Figure 10.

19 Figure 10 demonstrates the saturation vapor pressures achieved through the VaPOrS and
 20 compares them with their manually calculated counterparts. Automated results are shown as



1 colorful bars based on the number of detected functional groups reaching as high as 15 for
2 some autoAPRAMfw products, and manual results are depicted as points. The compatibility
3 of points atop bars for each species indicates excellent agreement between both approaches.

4



5

6 **Figure 10.** Comparison of saturation vapor pressure results for autoxidation species of benzene. The figure
7 illustrates the log saturation vapor pressure versus molar mass for species derived from autoxidation of initial
8 alkoxy and peroxy radicals of benzene degradation. The bars represent saturation vapor pressures calculated by
9 the automated VaPOrS code, with colors based on detected functional group number involved in the chemical
10 structure of species, while points indicate manually obtained values. The close alignment of points with bars
11 highlights the accuracy of the automated method relative to manual calculations.

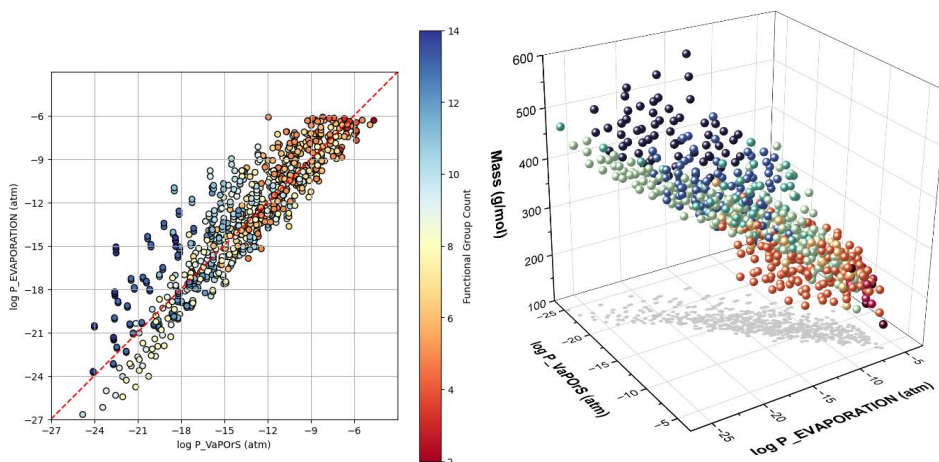
12

13 Figure 11 presents 2D and 3D scatter plots comparing the logarithmic saturation vapor
14 pressures obtained from VaPOrS ($\log P_{\text{VaPOrS}}$) with those predicted by the
15 EVAPORATION ($\log P_{\text{EVAPORATION}}$), Myrdal-Yalkowsky ($\log P_{\text{Myrdal_Yalkowsky}}$),
16 and Nanoolal ($\log P_{\text{Nanoolal}}$) methods for autoxidation products. Each marker represents a
17 compound, with color indicating the number of functional groups in its structure. The red
18 dashed line in the 2D plots signifies the theoretical 1:1 relationship, where the saturation vapor
19 pressures predicted by VaPOrS and the other models would be equivalent. Moreover, the 3D
20 plots include the molar mass of species to give more details of the achieved results.



1 Further analysis shows that the current VaPORs predictions based on SIMPOL
2 parameterization align closely with those from the EVAPORATION method in the higher
3 saturation vapor pressure range (approximately -6 to -15 on the logarithmic scale), with the
4 EVAPORATION method tending to overestimate saturation vapor pressures for species with
5 greater functional complexity. In contrast, VaPORs demonstrates good agreement with the
6 Nanoolal method in the lower saturation vapor pressure range (approximately -15 to -24),
7 particularly for molecules with a high functional group count. The Myrdal-Yalkowsky method,
8 however, totally overestimates saturation vapor pressures across the board compared to
9 VaPORs, with deviations increasing as functional group complexity rises. It is important to note
10 that the current values are essentially SIMPOL-based predictions, so while these comparisons
11 are informative, the general trends have been discussed in previous studies. However, as
12 highlighted by (Isaacman-VanWertz and Aumont 2021), a combination of existing methods,
13 potentially an average of them, has been suggested to yield the most reliable saturation vapor
14 pressure estimates. Acknowledging this, a future refinement of the current approach could
15 involve assessing whether incorporating such a hybrid method improves agreement with
16 experimental data.

17 The VaPORs model thus distinguishes itself by providing structure-based estimations for
18 saturation vapor pressure predictions, offering a valuable tool for the assessment of complex
19 organic compounds in atmospheric chemistry. This first application of the code utilized
20 SIMPOL group contribution method for saturation vapor pressure and enthalpy of vaporization
21 estimation, yet the code can be extended to work with any structure based thermodynamic
22 property estimator, thereby streamlining work like secondary aerosol modelling considerably.



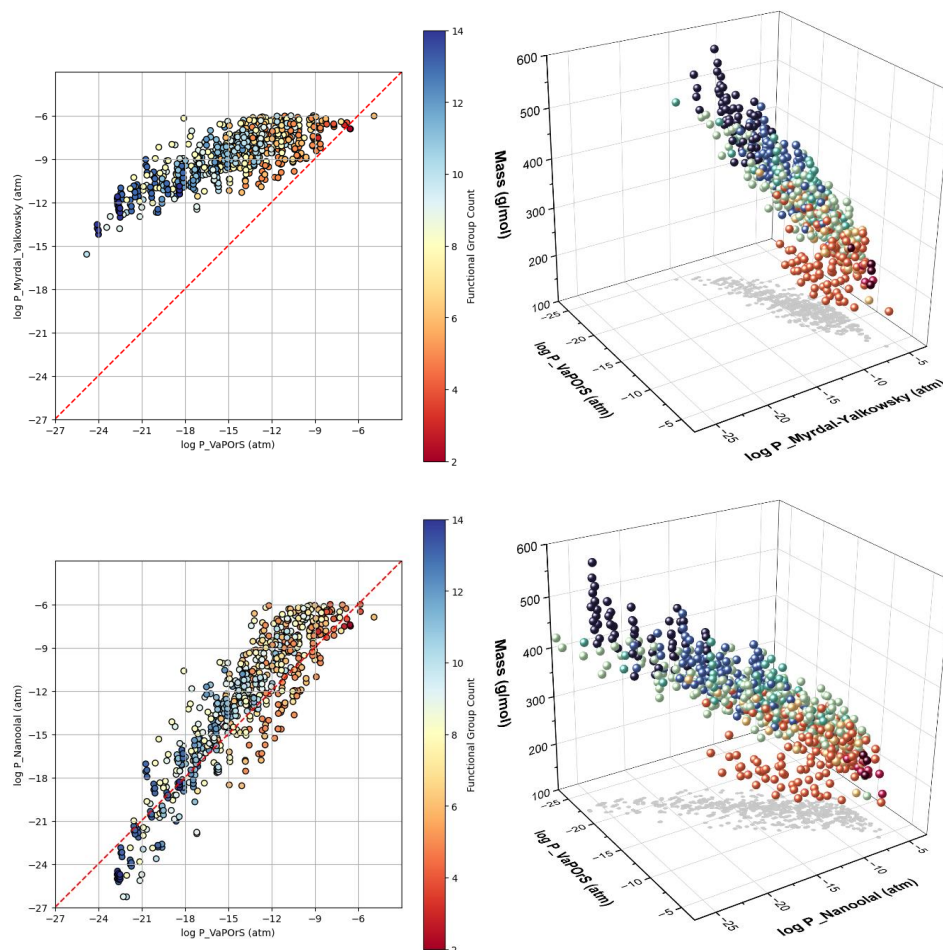


Figure 11. Comparison of logarithmic saturation vapor pressures predicted by the VaPOrS model ($\log P_{\text{VaPOrS}}$) with those from the EVAPORATION ($\log P_{\text{EVAPORATION}}$), Nanoolal ($\log P_{\text{Nanoolal}}$), and Myrdal-Yalkowsky ($\log P_{\text{Myrdal_Yalkowsky}}$) methods for autoxidation products derived from benzene degradation. Each data point represents a compound, with color indicating the functional group count. The red dashed line in the 2D plots represents the theoretical 1:1 relationship, where predictions from VaPOrS and other methods would be equivalent. The 3D plots include molar mass variation as well to give a comprehensive view of relationships between the parameters.

The versatility of VaPOrS lies in its potential for seamless integration into a range of widely used atmospheric models, enhancing their predictive capabilities. As seen, the Master Chemical Mechanism (MCM), with its comprehensive SMILES database of VOCs, can directly leverage VaPOrS to automate the generation of temperature-dependent saturation



1 vapor pressure equations. This automation ensures uniformity in saturation vapor pressure
2 predictions, which is critical for large-scale atmospheric models to simulate chemical reactions
3 and transport processes accurately. By providing a consistent and efficient means of handling
4 VOC property data, VaPOrS enhances MCM's ability to support atmospheric chemistry
5 research on gas-to-particle partitioning and secondary organic aerosol formation (Saunders et
6 al. 2003; M. E. Jenkin et al. 2003).

7 In MEGAN (A. Guenther et al. 2006), VaPOrS can improve biogenic VOC emission estimates
8 by rapidly providing saturation vapor pressure and enthalpy of vaporization data. Similarly,
9 LOTOS-EUROS (Schaap et al. 2008) can benefit from more accurate secondary organic
10 aerosol formation and gas-to-particle partitioning predictions. Global-scale models like GEOS-
11 Chem (Bey et al. 2001) and WRF-Chem (Grell et al. 2005) are supported by VaPOrS's ability
12 to process large datasets of VOCs efficiently, enabling more reliable simulations of
13 atmospheric chemical processes. Regional models such as CMAQ (Byun and Schere 2006) are
14 also enhanced, with VaPOrS contributing to improved particulate matter formation predictions.

15 Additionally, in SAPRC (Carter, n.d.) mechanisms, VaPOrS can automate the calculation of
16 key VOC properties, streamlining computational workflows. Beyond these applications,
17 specialized models like ADCHAM and ADCHEM (Roldin et al. 2014; 2011) can leverage
18 VaPOrS to refine aerosol growth rate predictions and radiative forcing estimations, thus
19 advancing studies on aerosol-cloud interactions and climate feedback mechanisms. These
20 integrations highlight VaPOrS's broad applicability and its role in improving the accuracy and
21 efficiency of atmospheric models, ultimately contributing to more informed strategies for air
22 quality management and climate change mitigation.

23 It is worth mentioning that while group contribution methods, such as SIMPOL and VaPOrS,
24 offer a reliable approach for estimating saturation vapor pressure by summing the contributions
25 of individual functional groups, their accuracy tends to decrease as molecular complexity
26 increases. This decline in predictive capability is particularly noticeable in highly
27 functionalized compounds, where interactions between multiple functional groups may deviate
28 from the assumed additive behavior. Previous studies have demonstrated that the presence of
29 numerous functional groups does not always lead to a proportional reduction in saturation
30 vapor pressure, as steric effects, intramolecular hydrogen bonding, and other molecular
31 interactions can alter the expected contributions. This limitation highlights the need for caution
32 when applying these models to large multifunctional molecules and suggests potential areas



1 for future refinement, such as incorporating correction factors or machine-learning approaches
2 to account for non-additive effects.

3 **4. Conclusion**

4 This study presents a comprehensive analysis of the performance of the VaPOrS model,
5 developed for identifying and quantifying functional groups in organic compounds with a
6 specific focus on saturation vapor pressure calculations based on the SIMPOL method. The
7 research highlights three critical aspects: the identification and counting of 30 structural groups
8 recognized as influential on saturation vapor pressure, the subsequent calculation of saturation
9 vapor pressure and enthalpy of vaporization for a range of organic compounds, and
10 introduction of a temperature-dependent relationship for the saturation vapor pressure and
11 enthalpy of vaporization. This study primarily utilizes the SIMPOL method for
12 parameterization; however, the VaPOrS framework is inherently flexible and can be adapted
13 to various other structure-based parameterization methods. These include approaches such as
14 group additivity and volatility basis set (VBS) models. Furthermore, VaPOrS can be expanded
15 to estimate a range of thermodynamic properties that depend on structural group-based
16 parameterization, broadening its applicability beyond saturation vapor pressure alone. The
17 validation process involved a meticulous comparison of manual and automated counts and
18 calculations for over 1,000 organic compounds sourced from the Master Chemical Mechanism
19 (MCM) database and recently introduced autoAPRAMfw autooxidation code. The perfect
20 agreement demonstrates the accuracy for functional group identification and counting, with
21 subsequent semi-automated saturation vapor pressure determination by SIMPOL through
22 VaPOrS. This new methodology will be predominantly useful for researchers analyzing
23 organic compounds, particularly in fields related to atmospheric chemistry and especially
24 relating to aerosol formation. In conclusion, the findings from this research validate VaPOrS
25 as a robust computational tool for estimating saturation vapor pressures while providing a
26 systematic approach to functional group analysis. Its accuracy in functional group
27 quantification and saturation vapor pressure prediction demonstrates substantial potential for
28 advancing our understanding of organic compound behaviors in atmospheric and
29 environmental applications. Future work will focus on enhancing the model's capabilities,
30 exploring additional functional groups, and refining the saturation vapor pressure estimation
31 model to further improve its applicability and precision in real-world scenarios.

32 **Data availability**



1 The raw data supporting the figures in the manuscript are openly available on Zenodo at:
2 <https://doi.org/10.5281/zenodo.15688105>.

3 **Code availability**

4 The VaPOrS code used in this study is publicly available on Zenodo (Mojtaba Bezaatpour
5 2025). This archive includes a Jupyter notebook (VaPOrS.ipynb), a standalone Python script
6 (VaPOrS.py), an input file containing SMILES strings (SMILES.txt), and example output files
7 in both .txt and .csv formats. The repository is licensed under the MIT License and is fully
8 open for use and redistribution under the conditions specified therein.

9 **Author contributions**

10 M.B. conceptualized the study and developed the Python code; M.B. and M.R. conducted the
11 functional group analysis, validated the tool against existing data and contributed to data
12 visualization; M.B. prepared the manuscript, and M.D.M. and M.R. reviewed and approved the
13 final version of the manuscript.

14 **Competing interests**

15 The authors declare that they have no conflict of interest.

16 **Acknowledgements**

17 This project has received funding from the European Research Council under the European
18 Union's Horizon 2020 research and innovation programme under Grant No. 101002728 (ERC
19 Consolidator grant ADAPT) and 101096133 (PAREMPI). This work is also funded by the
20 Research Council of Finland (Grant Nos.: 331207, 336531, 346373, 353836). The AI-based
21 tools were used for language editing to improve the readability of the manuscript.

22 **References**

- 23 Atkinson, R. 2000. "Atmospheric Chemistry of VOCs and NOx." *Atmospheric Environment* 34 (12–
24 14): 2063–2101. [https://doi.org/10.1016/S1352-2310\(99\)00460-4](https://doi.org/10.1016/S1352-2310(99)00460-4).
25 Berndt, Torsten, Stefanie Richters, Tuija Jokinen, Noora Hyttinen, Theo Kurtén, Rasmus V. Otkjær,
26 Henrik G. Kjaergaard, et al. 2016. "Hydroxyl Radical-Induced Formation of Highly Oxidized
27 Organic Compounds." *Nature Communications* 7 (1): 13677.
28 <https://doi.org/10.1038/ncomms13677>.
29 Bey, Isabelle, Daniel J. Jacob, Robert M. Yantosca, Jennifer A. Logan, Brendan D. Field, Arlene M.
30 Fiore, Qinbin Li, Hongyue Y. Liu, Loretta J. Mickley, and Martin G. Schultz. 2001. "Global



- 1 Modeling of Tropospheric Chemistry with Assimilated Meteorology: Model Description and
2 Evaluation." *Journal of Geophysical Research: Atmospheres* 106 (D19): 23073–95.
3 <https://doi.org/10.1029/2001JD000807>.
- 4 Bianchi, Federico, Theo Kurtén, Matthieu Riva, Claudia Mohr, Matti P. Rissanen, Pontus Roldin,
5 Torsten Berndt, et al. 2019. "Highly Oxygenated Organic Molecules (HOM) from Gas-Phase
6 Autoxidation Involving Peroxy Radicals: A Key Contributor to Atmospheric Aerosol."
7 *Chemical Reviews* 119 (6): 3472–3509. <https://doi.org/10.1021/acs.chemrev.8b00395>.
- 8 Byun, Daewon, and Kenneth L. Schere. 2006. "Review of the Governing Equations, Computational
9 Algorithms, and Other Components of the Models-3 Community Multiscale Air Quality
10 (CMAQ) Modeling System." *Applied Mechanics Reviews* 59 (2): 51–77.
11 <https://doi.org/10.1115/1.2128636>.
- 12 Carter, William P L. n.d. "Documentation of the SAPRC-99 Chemical Mechanism for VOC
13 Reactivity Assessment, Report to the California Air Resources Board". CA: University of
14 California, Riverside, USA, 2000, <https://intra.engr.ucr.edu/~carter/pubs/s99doc.pdf>.
- 15 Compennolle, S., K. Ceulemans, and J.-F. Müller. 2011. "EVAPORATION: A New Vapour Pressure
16 Estimation Method for Organic Molecules Including Non-Additivity and Intramolecular
17 Interactions." *Atmospheric Chemistry and Physics* 11 (18): 9431–50.
18 <https://doi.org/10.5194/acp-11-9431-2011>.
- 19 Crounse, John D., Lasse B. Nielsen, Solvejg Jørgensen, Henrik G. Kjaergaard, and Paul O. Wennberg.
20 2013. "Autoxidation of Organic Compounds in the Atmosphere." *The Journal of Physical
21 Chemistry Letters* 4 (20): 3513–20. <https://doi.org/10.1021/jz4019207>.
- 22 Donahue, N. M., A. L. Robinson, C. O. Stanier, and S. N. Pandis. 2006. "Coupled Partitioning,
23 Dilution, and Chemical Aging of Semivolatile Organics." *Environmental Science &
24 Technology* 40 (8): 2635–43. <https://doi.org/10.1021/es052297c>.
- 25 Eckert, Frank, and Andreas Klamt. 2002. "Fast Solvent Screening via Quantum Chemistry: COSMO-
26 RS Approach." *AIChE Journal* 48 (2): 369–85. <https://doi.org/10.1002/aic.690480220>.
- 27 Ehn, Mikael, Joel A. Thornton, Einhard Kleist, Mikko Sipilä, Heikki Junninen, Iida Pullinen, Monika
28 Springer, et al. 2014. "A Large Source of Low-Volatility Secondary Organic Aerosol." *Nature*
29 506 (7489): 476–79. <https://doi.org/10.1038/nature13032>.
- 30 Epping, Ruben, and Matthias Koch. 2023. "On-Site Detection of Volatile Organic Compounds
31 (VOCs)." *Molecules* 28 (4): 1598. <https://doi.org/10.3390/molecules28041598>.
- 32 Goldstein, Allen H., and Ian E. Galbally. 2007. "Known and Unexplored Organic Constituents in the
33 Earth's Atmosphere." *Environmental Science & Technology* 41 (5): 1514–21.
34 <https://doi.org/10.1021/es072476p>.
- 35 Grell, Georg A., Steven E. Peckham, Rainer Schmitz, Stuart A. McKeen, Gregory Frost, William C.
36 Skamarock, and Brian Eder. 2005. "Fully Coupled 'Online' Chemistry within the WRF
37 Model." *Atmospheric Environment* 39 (37): 6957–75.
38 <https://doi.org/10.1016/j.atmosenv.2005.04.027>.
- 39 Guenther, A., T. Karl, P. Harley, C. Wiedinmyer, P. I. Palmer, and C. Geron. 2006. "Estimates of
40 Global Terrestrial Isoprene Emissions Using MEGAN (Model of Emissions of Gases and
41 Aerosols from Nature)." *Atmospheric Chemistry and Physics* 6 (11): 3181–3210.
42 <https://doi.org/10.5194/acp-6-3181-2006>.
- 43 Guenther, Alex, C. Nicholas Hewitt, David Erickson, Ray Fall, Chris Geron, Tom Graedel, Peter
44 Harley, et al. 1995. "A Global Model of Natural Volatile Organic Compound Emissions."
45 *Journal of Geophysical Research: Atmospheres* 100 (D5): 8873–92.
46 <https://doi.org/10.1029/94JD02950>.
- 47 Isaacman-VanWertz, Gabriel, and Bernard Aumont. 2021. "Impact of Organic Molecular Structure on
48 the Estimation of Atmospherically Relevant Physicochemical Parameters." *Atmospheric
49 Chemistry and Physics* 21 (8): 6541–63. <https://doi.org/10.5194/acp-21-6541-2021>.
- 50 Iyer, Siddharth. 2023. "Molecular Rearrangement of Bicyclic Peroxy Radicals: Key Route to Aerosol
51 from Aromatics," February. <https://doi.org/10.5281/ZENODO.8214481>.
- 52 Iyer, Siddharth, Matti P. Rissanen, Rashid Valiev, Shawon Barua, Jordan E. Krechmer, Joel Thornton,
53 Mikael Ehn, and Theo Kurtén. 2021a. "Molecular Mechanism for Rapid Autoxidation in α -
54 Pinene Ozonolysis." *Nature Communications* 12 (1): 878. <https://doi.org/10.1038/s41467-021-21172-w>.



- 1 ———. 2021b. “Molecular Mechanism for Rapid Autoxidation in α -Pinene Ozonolysis.” *Nature*
2 *Communications* 12 (1): 878. <https://doi.org/10.1038/s41467-021-21172-w>.
- 3 Jenkin, M. E., S. M. Saunders, V. Wagner, and M. J. Pilling. 2003. “Protocol for the Development of
4 the Master Chemical Mechanism, MCM v3 (Part B): Tropospheric Degradation of Aromatic
5 Volatile Organic Compounds.” *Atmospheric Chemistry and Physics* 3 (1): 181–93.
6 <https://doi.org/10.5194/acp-3-181-2003>.
- 7 Jenkin, Michael E., Richard Valorso, Bernard Aumont, and Andrew R. Rickard. 2019. “Estimation of
8 Rate Coefficients and Branching Ratios for Reactions of Organic Peroxy Radicals for Use in
9 Automated Mechanism Construction.” *Atmospheric Chemistry and Physics* 19 (11): 7691–
10 7717. <https://doi.org/10.5194/acp-19-7691-2019>.
- 11 Jimenez, J. L., M. R. Canagaratna, N. M. Donahue, A. S. H. Prevot, Q. Zhang, J. H. Kroll, P. F.
12 DeCarlo, et al. 2009. “Evolution of Organic Aerosols in the Atmosphere.” *Science* 326
13 (5959): 1525–29. <https://doi.org/10.1126/science.1180353>.
- 14 Joback, K.G., and R.C. Reid. 1987. “ESTIMATION OF PURE-COMPONENT PROPERTIES FROM
15 GROUP-CONTRIBUTIONS.” *Chemical Engineering Communications* 57 (1–6): 233–43.
16 <https://doi.org/10.1080/00986448708960487>.
- 17 Klamt, Andreas. 1995. “Conductor-like Screening Model for Real Solvents: A New Approach to the
18 Quantitative Calculation of Solvation Phenomena.” *The Journal of Physical Chemistry* 99 (7):
19 2224–35. <https://doi.org/10.1021/j100007a062>.
- 20 Klamt, Andreas, Volker Jonas, Thorsten Bürger, and John C. W. Lohrenz. 1998. “Refinement and
21 Parametrization of COSMO-RS.” *The Journal of Physical Chemistry A* 102 (26): 5074–85.
22 <https://doi.org/10.1021/jp980017s>.
- 23 Luo, Hao, Luc Vereecken, Hongru Shen, Sungah Kang, Iida Pullinen, Mattias Hallquist, Hendrik
24 Fuchs, et al. 2023. “Formation of Highly Oxygenated Organic Molecules from the Oxidation
25 of Limonene by OH Radical: Significant Contribution of H-Abstraction Pathway.”
26 *Atmospheric Chemistry and Physics* 23 (13): 7297–7319. <https://doi.org/10.5194/acp-23-7297-2023>.
- 27
- 28 Mattila, Timo, Markku Kulmala, and Timo Vesala. 1997. “On the Condensational Growth of a
29 Multicomponent Droplet.” *Journal of Aerosol Science* 28 (4): 553–64.
30 [https://doi.org/10.1016/S0021-8502\(96\)00458-2](https://doi.org/10.1016/S0021-8502(96)00458-2).
- 31 McDonald, Brian C., Joost A. De Gouw, Jessica B. Gilman, Shantanu H. Jathar, Ali Akherati,
32 Christopher D. Cappa, Jose L. Jimenez, et al. 2018. “Volatile Chemical Products Emerging as
33 Largest Petrochemical Source of Urban Organic Emissions.” *Science* 359 (6377): 760–64.
34 <https://doi.org/10.1126/science.aag0524>.
- 35 Mellouki, A., T. J. Wallington, and J. Chen. 2015. “Atmospheric Chemistry of Oxygenated Volatile
36 Organic Compounds: Impacts on Air Quality and Climate.” *Chemical Reviews* 115 (10):
37 3984–4014. <https://doi.org/10.1021/cr500549n>.
- 38 Mojtaba Bezaatpour. 2025. “Mojtababzp/VaPOrS: VaPOrS v1.0.1.” Zenodo.
39 <https://doi.org/10.5281/ZENODO.15222175>.
- 40 Myrdal, Paul B., and Samuel H. Yalkowsky. 1997. “Estimating Pure Component Vapor Pressures of
41 Complex Organic Molecules.” *Industrial & Engineering Chemistry Research* 36 (6): 2494–
42 99. <https://doi.org/10.1021/ie950242l>.
- 43 Nannoolal, Yash, Jürgen Rarey, and Deresh Ramjugernath. 2008. “Estimation of Pure Component
44 Properties.” *Fluid Phase Equilibria* 269 (1–2): 117–33.
45 <https://doi.org/10.1016/j.fluid.2008.04.020>.
- 46 Pankow, J. F., and W. E. Asher. 2008. “SIMPOL.1: A Simple Group Contribution Method for
47 Predicting Vapor Pressures and Enthalpies of Vaporization of Multifunctional Organic
48 Compounds.” *Atmospheric Chemistry and Physics* 8 (10): 2773–96.
49 <https://doi.org/10.5194/acp-8-2773-2008>.
- 50 Pichelstorfer, Lukas, Pontus Roldin, Matti Rissanen, Noora Hyttinen, Olga Garmash, Carlton Xavier,
51 Putian Zhou, et al. 2024. “Towards Automated Inclusion of Autoxidation Chemistry in
52 Models: From Precursors to Atmospheric Implications.” *Environmental Science: Atmospheres*
53 4 (8): 879–96. <https://doi.org/10.1039/D4EA00054D>.
- 54 Rissanen, Matti P., Theo Kurtén, Mikko Sipilä, Joel A. Thornton, Juha Kangasluoma, Nina Sarnela,
55 Heikki Junninen, et al. 2014. “The Formation of Highly Oxidized Multifunctional Products in



- 1 the Ozonolysis of Cyclohexene.” *Journal of the American Chemical Society* 136 (44): 15596–
2 606. <https://doi.org/10.1021/ja507146s>.
- 3 Roldin, P., A. C. Eriksson, E. Z. Nordin, E. Hermansson, D. Mogensen, A. Rusanen, M. Boy, et al.
4 2014. “Modelling Non-Equilibrium Secondary Organic Aerosol Formation and Evaporation
5 with the Aerosol Dynamics, Gas- and Particle-Phase Chemistry Kinetic Multilayer Model
6 ADCHEM.” *Atmospheric Chemistry and Physics* 14 (15): 7953–93.
7 <https://doi.org/10.5194/acp-14-7953-2014>.
- 8 Roldin, P., E. Swietlicki, G. Schurgers, A. Arneth, K. E. J. Lehtinen, M. Boy, and M. Kulmala. 2011.
9 “Development and Evaluation of the Aerosol Dynamics and Gas Phase Chemistry Model
10 ADCHEM.” *Atmospheric Chemistry and Physics* 11 (12): 5867–96.
11 <https://doi.org/10.5194/acp-11-5867-2011>.
- 12 Saunders, S. M., M. E. Jenkin, R. G. Derwent, and M. J. Pilling. 2003. “Protocol for the Development
13 of the Master Chemical Mechanism, MCM v3 (Part A): Tropospheric Degradation of Non-
14 Aromatic Volatile Organic Compounds.” *Atmospheric Chemistry and Physics* 3 (1): 161–80.
15 <https://doi.org/10.5194/acp-3-161-2003>.
- 16 Schaap, Martijn, Renske M.A. Timmermans, Michiel Roemer, G.A.C. Boersen, Peter J.H. Builtjes,
17 Ferd J. Sauter, Guus J.M. Velders, and Jeanette P. Beck. 2008. “The LOTOS EUROS Model:
18 Description, Validation and Latest Developments.” *International Journal of Environment and*
19 *Pollution* 32 (2): 270. <https://doi.org/10.1504/IJEP.2008.017106>.
- 20 Topping, David, Mark Barley, Michael K. Bane, Nicholas Higham, Bernard Aumont, Nicholas
21 Dingle, and Gordon McFiggans. 2016. “UMansysProp v1.0: An Online and Open-Source
22 Facility for Molecular Property Prediction and Atmospheric Aerosol Calculations.”
23 *Geoscientific Model Development* 9 (2): 899–914. <https://doi.org/10.5194/gmd-9-899-2016>.
- 24 Vereecken, L., B. Aumont, I. Barnes, J.W. Bozzelli, M.J. Goldman, W.H. Green, S. Madronich, et al.
25 2018. “Perspective on Mechanism Development and Structure-Activity Relationships for
26 Gas-Phase Atmospheric Chemistry.” *International Journal of Chemical Kinetics* 50 (6): 435–
27 69. <https://doi.org/10.1002/kin.21172>.
- 28 WHO. 2005. “Air Quality Guidelines: Global Update 2005. Geneva: World Health Organization.”
- 29 Zhao, Defeng, Iida Pullinen, Hendrik Fuchs, Stephanie Schrade, Rongrong Wu, Ismail-Hakki Acir,
30 Ralf Tillmann, et al. 2021. “Highly Oxygenated Organic Molecule (HOM) Formation in the
31 Isoprene Oxidation by NO₃ Radical.” *Atmospheric Chemistry and*
32 *Physics* 21 (12): 9681–9704. <https://doi.org/10.5194/acp-21-9681-2021>.