

Electronegativity Equalization in Molecular Mechanics

Magnus Ringholm

Trondheim, June 11 2009

Declaration

I declare that the work presented in this thesis has been carried out independently and in agreement with “Reglement for sivilarkitekt- og sivilingeniøreksamen”.

Trondheim, June 11 2009

Magnus Ringholm

Preface

This work has been a rewarding and challenging task. In the 20 weeks allotted for the production of this thesis, several aspects encountered in computational chemistry have been visited and brought together for the results presented in this document. The largest single task undertaken has been the implementation from scratch of a genetic algorithm for the obtainment of parameters for the molecular mechanics model featured in this work. Time was also spent on the implementation of this model during the first weeks. Another large task was the obtainment of reference values for the molecules used in parametrization. For this purpose, quantum chemical calculations have been carried out for several hundred molecular systems with sizes ranging from a few atoms to upwards of sixty. Finally, the obtainment of results and their interpretation has naturally taken its share of time. Much knowledge has been gained in this time, primarily about the subject matter itself, but also about the general ways in which research is carried out.

I would sincerely like to thank my supervisor, Per-Olof Åstrand, for his valuable assistance, guidance and support during the entire course of this work. Special thanks are also given to Hans Sverre Smalø. His general help and assistance, especially in the implementation of the model and in the interpretation and discussion of results, has been valuable and is much appreciated.

Thanks are also given to my friends at the institute, including, but not limited to, Sondre Schnell Kvalvåg, Miriam Mekki, Ingrid Aaen, Mari Voldsund, Øivind Wilhelmsen, Ragnhild Skorpa, Marit Takla, Odne Burheim, Einar Ryeng and Anders Lervik for their support and high spirits, brightening the mood and making the years of study enjoyable. My flatmates Gøran Berntsen and Martin Thuve Hovden are thanked for their valuable friendship and putting up with me on a daily basis. Others not named are also thanked for their friendship.

Finally, I would like to thank my family for their kind support and motivation during my upbringing and education, for the latter particularly during the work on this thesis.

Abstract

The obtainment of atom-type parameters for a new molecular mechanics model for the calculation of the molecular dipole moment and molecular polarizability has been carried out using a genetic algorithm followed by local minimization. For the calculation of the molecular polarizability, parameters have been obtained for the elements hydrogen, carbon, oxygen, fluorine and chlorine. For the calculation of the molecular dipole moment, parameters have been obtained for hydrogen, carbon, fluorine and chlorine. The genetic algorithm is a suitable choice for this purpose. The model is able to reproduce molecular polarizabilities in good agreement with reference values for systems containing all elements for which parameters have been found. It also shows good agreement with reference values for the molecular dipole moment for halogenated alkane systems and some aromatic systems, but does not provide satisfactory agreement for other aromatic systems and alkene chains with alternating double bonds.

Contents

1	Introduction	1
2	Theory	5
2.1	The extended EEM/PDI combination model	5
	Introduction and summary of the models being combined	5
	The combined EEM/PDI model	5
	Modifying the combined EEM/PDI model	11
2.2	Genetic algorithms	18
	Introduction	18
	The individual	19
	Initial settings, selection, recombination and mutation	21
	Measuring GA performance and convergence	25
	Shortcomings and challenges	26
	Additional features	27
2.3	Finding an experimental method	28
	Obtaining reference values	29
	Setting up the optimization	30
	Convergence and ending the search	34
	Summary of the experimental procedure	34
3	Implementation and experimental details	36
3.1	The obtainment of reference values	36
3.2	The genetic algorithm and local minimization. Experimental procedure	36
4	Results and discussion	39
4.1	Polarizability	39
4.2	Dipole moment	50
4.3	Performance of the parametrization method	56
4.4	Suggested directions for further work	57
5	Conclusion	59

1 Introduction

The development of methods for the computation of electromagnetic molecular properties is a vast field of study, and a lot of effort has been made to attain higher degrees of accuracy or decrease calculation times. Today, while certainly not depleted as a subject of further study, for many properties there exist accurate and efficient methods of calculation in the quantum mechanical (QM) or *ab initio* framework. Several programs have been created for this purpose, each capable of reliable calculations of such properties both in single- and multiprocessor setups. In the development of these programs, a lot of care is taken in ensuring that algorithms show favorable scaling under ever increasing degrees of parallelization. This effort to parallelize has been essential for these programs, as the level of complexity encountered at QM levels of theory demand that, in order to be feasible in terms of wall-time, multiple computational nodes are employed for all but the simplest molecular systems or levels of theory. The high-order scaling of the quantum chemical methods (N^3 or higher, where N is the relevant size measure of the molecular system) means that their application, even in highly parallelized setups, are restricted to systems of moderate sizes [1]. Furthermore, even for systems of modest sizes, the high level of sophistication encountered in *ab initio* methods leads to a large prefactor before the scaling-determining N , limiting the number of molecular systems that can be studied in one project.

The situation for a research project may be one where the computational scenario is such that it falls under the restrictions presented above. For example, there may be only a few systems, but these systems could be so large as to render the quantum chemical computation intractable, or it may be desired to perform a screening of several thousand smaller molecules to identify those that possess a desired property or properties. In the former scenario, a qualitative or approximate result for the property under study may be deemed satisfactory. In the latter, the researcher may be willing to sacrifice some accuracy in order to obtain some approximate results or qualitative answers to identify “promising leads” for further study. These “leads” would perhaps in later work be subjects of more extensive investigation at a high level of theory, being few in number compared to the original set of candidates. As the computational resources available for each research project is not always abundant, the availability of such approximative methods is important for enabling such projects to be undertaken.

Providing such quick methods of calculation while retaining as much precision as possible could be said to be the central objective of force-field (FF) or molecular mechanics (MM) methods. The mechanics employed in these methods are generally Newtonian in nature, although hybrid semi-empirical methods seek to combine the precision of quantum mechanics with the simplicity of classical theories in what can be characterised as an intermediate level of theory. In any case, MM/FF methods, when successful, can provide results approaching quantum mechanical or experimental values while yielding drastic reductions in computational complexity, and thus, computing time. Generally, a good force field method will give good precision, possibly reducing the system size scaling to a lower order (typically no larger than N^3) and the prefactor by several orders of magnitude. Clearly, such methods could have a range of application far beyond the limits of quantum mechanical treatment.

The essential challenge when designing force-field methods are identifying the key mechanisms from which the properties of interest arise, and then describing them in such a way that the cost of their calculation are small. This will generally mean developing a simplified model of the system appropriate for the situation at hand, trying to imitate effects shown by higher levels of theory. Sometimes, the first attempt at making a model will be too simple, and steps must be taken to add the right amount of sophistication to the right places. Developing a successful method in this way may be a difficult task, but

aside from the obvious value of greatly simplified computation, there is also considerable physical understanding to be gained from identifying the important contributions to the desired property.

Ubiquitous in force-field methods are parameters [1]. Depending on the way in which the model regards or partitions the system, these numbers may be a combination of any or either of atom-type parameters with or without regard to hybridization modes, bond-type parameters, bond angle or dihedral angle parameters, global parameters, parameters for certain functional groups such as methyl or carboxyl groups, or even parameters governing secondary or tertiary structural features. The dependence on parameters could in broad terms be said to come from the fact that force-field models are essentially simplifications of QM theory. This means that some features, which would explicitly or implicitly be calculated fully by a QM model, are condensed into parameters. The parameters could then largely be said to represent features that one is willing to regard as invariant over the range for which one wishes the model to apply.

Often, the number of parameters used in the model will be intimately connected to the number of features contributing to the quantity of interest that the model endeavors to encompass. Furthermore, a force-field model describing most essential features of the property at hand would commonly be deemed elegant or, more strictly, worthwhile, if it is able to reproduce effects from QM levels of theory by using only a few moderately simple principles for its features. From this it can be argued that the merit of a model, compared to other models describing the same quantity, keeps some relation to the number of parameters it employs, where using fewer parameters to obtain largely the same result is better. Conversely, combining in a model a large number of effects, some of which may even reasonably be deemed entirely negligible, and thereby using a large number of parameters, would perhaps provide a small or moderate gain in accuracy, but it would also look bloated when compared to other, more simple methods. The use of such a model, while still likely to be efficient when compared to QM methods, may be much more computationally expensive than simpler approaches. The case is then made for keeping the model, and thus its number of parameters, as simple as could reasonably be done.

The model featured in this article is a force-field model principally intended for the description of molecular dipole moments and molecular dipole-dipole polarizabilities. This model [2], developed by Per-Olof Åstrand, Hans Sverre Smalø, and Lasse Jensen, combines the attractive features of the electronegativity equalization model (EEM)[3]-[8] with the point dipole interaction (PDI) model [9]-[13]. The model seeks to reproduce dipole moments and polarizabilities with high precision at a greatly reduced cost compared to QM approaches. The model as presented in the original article employs, with one exception, atom-type parameters for the description of key properties. The exception is one global parameter used in one of the terms governing charge transfer from one atom to a third atom through a second atom. It is hoped that the model will be able to describe atomic charges, dipole moments and polarizabilities for a wide range of compounds comprised of several kinds of elements.

In order to accomplish this, the parameters used in the model must be determined in such a way that the model is capable of describing these properties for molecules which may be quite distinct in nature. As they are mainly atom-type parameters, this means that their values should be representative of as many as possible of the situations in which atoms of a particular element may appear in molecular systems. The development of a method for reliably determining these parameters to make the model function to the best of its ability, possibly identifying where the model falls short, is the principal task of this work.

The task of determining these parameters is best regarded as a global optimization prob-

lem. The problem becomes one of finding, in the space spanned by the parameter values, a set of parameters from which the model can determine molecular dipole moments and polarizabilities sufficiently close to the correct values for a wide enough range of molecular systems. While this is a problem which is easily stated, the task in fact contains several challenges, some more subtle than others. These challenges will be presented in further detail in sections to come, but a brief summary is provided here.

The first problem is determining what actually constitutes a good set of parameters, that is, choosing a suitable function to be optimized. As favorable sets of parameters are those which produce dipole moments or polarizabilities close to the correct values, such a function must be one that is intimately linked to the difference between the values calculated by the model for a given set of parameters, and the correct values of these properties. Several choices exist for this purpose and will be discussed later, but a more important question is: What are the “correct” values? It is natural to suggest experimental values for this purpose, but especially for polarizabilities, such results are scarce, and the experimental accuracy may also be undesirably high. Furthermore, restricting oneself to a limited set of results, perhaps obtained from a variety of sources, carry two large disadvantages. The first one is the risk that some of these values may be erroneous, with errors lending themselves to differences in experimental method, setup or apparatus. By using such values, one is left vulnerable to the quality of the experimental results. The second disadvantage is that the limited set of “correct” answers from which to choose may result in an unwanted narrowing of the range of molecules of which one may judge the model to be a good description. By using QM calculated values as the “answer key”, the departure from which being the basis for the function to be optimized, one would only be limited by time and computational resources in the number of molecules available for study. Furthermore, by using the same method for all molecules thus calculated upon (assuming that the method does not employ random or pseudorandom numbers), one could disregard errors on the level of reproducibility. Systematic errors may still be present due to limitations of the level of theory employed in the QM calculations, but this can possibly be mitigated by using a high level of theory. All in all, using QM results as reference values provides a much higher level of confidence than would be attained from experimental values, and the flexibility in the number of systems available for study is another big advantage.

We have now briefly touched upon another problem: Selecting a set of molecules large enough to span a sufficiently high range of applicability. As the model uses atom-type parameters, this means providing reference values for a set of molecules large enough so that each element for which one wishes to find parameters is sufficiently described, in the sense that every “role” that the element can play in molecular systems one wishes the model to apply to, is covered by the set of reference molecules.

Next, the method used to determine these parameters should be able to locate a global minimum, or at least span a large part of parameter space in the search. For some parameters, one may have a certain degree of *a priori* knowledge, but the likely values of other parameters would be unknown. A local minimization procedure would then be more vulnerable to poor initial parameter choices than would a global procedure. In global optimization procedures, one may also need to specify ranges in which the values of the parameters may vary. With these ranges, one is allowed some degree of precision in specifying one’s uncertainty about each parameter by allowing a wide or narrow range.

Finally, the model itself may be insufficient or inaccurate for part of the range of molecules one would wish or believe it to be accurate. This may be challenging to identify in an optimization, as such apparent errors may stem from several other factors, most notably an imbalance in the set of reference molecules used in the optimization. For an example of the latter, one may wish to optimize parameters for oxygen when parameters for carbon and hydrogen are already obtained. Naively, a reference set of

mostly alcohols and a few aldehydes is chosen. The minimization may then result in excellent results for alcohols and poor results for aldehydes, when in fact the model could have given good results for both kinds of systems had a more balanced reference set been used.

An optimization method believed to be promising for this kind of parametrization is the genetic algorithm. This algorithm, to be presented in detail later, is capable of reaching a global optimum and span a large part of parameter space [14]. It is required to specify ranges for the parameters being optimized, allowing for application of *a priori* knowledge. One is free to choose from several features to include in the algorithm, based on what one wishes to emphasize and the computational resources at hand. The implementation time for such an algorithm is moderate, and makes no use of derivatives, the analytical determination of which could be a hard task.

Before moving on to presenting the theory used in this work, a brief overview of the sections to come are given here. We start with a description of the force-field model for the calculation of molecular dipole moments and polarizabilities. Then, a presentation of the genetic algorithm with key and optional features is given. Next, the procedure surrounding the optimization task is presented. Following this, implementation details are given. Next, the results will be presented, accompanied by pertinent discussion. We will end with suggestions for further work and some concluding remarks.

2 Theory

2.1 The extended EEM/PDI combination model

Introduction and summary of the models being combined

This section presents the extended EEM/PDI combination model developed by Smalø, Åstrand and Jensen [2]. This model will be referred to as “the model”, or “the new model” when compared to other models. The model will here be presented in a way closely resembling that of the original article [2]. After this presentation, we will move on to present some modifications to this model, also given in the original article. Finally, another modification to the model will be presented, the idea of which has arisen during the course of the present work.

The model is one combining, reformulating and extending the aspects of the electronegativity equalization model (EEM) and the point dipole interaction (PDI)[9]-[13] model. As its name may reveal, the EEM seeks to create a representation of the molecular electronegativity by assigning to each atom in the molecule an atomic electronegativity and chemical hardness, regarded as atom-type parameters. In broad terms, the atomic electronegativity could be said to be the amount of electronegativity that each atom “brings to the table”, that is, its initial contribution to the molecular electronegativity before equalization takes place. Its chemical hardness is then a measure of its “reluctance” to give away some of this electronegativity to other parts of the molecule. By allowing charge to flow between the atoms in the molecule, an state of equilibrium, indeed, an equalization, is attained. From this equalization, properties such as atomic charges and molecular dipole moments and polarizabilities can be found. The model seeks to be in analogy with density functional theory, in which concepts such as charge density, electronegativity and chemical hardness are regularly and readily addressed [20].

The main limitation of the EEM is that it is essentially a metallic model. This is because in the EEM, charge is allowed to flow freely between the atoms in the molecule. This is a simplification which may yield inexact results for molecules that don’t follow this behavior, and will lead to clearly erroneous results for nonmetallic molecules with long extensions in one or more directions, of which the most important example is carbon chains. For these molecules, the polarizability resulting from EEM equalization will approach infinity as the chains lengthen, due to the false condition that there are no barriers to the flow of charge between the atoms of these molecules. This problem has been addressed in various ways [21]-[25], and as will be presented later, the new model adds an energy-cost term for the transport of charge through the molecule. In this way, metallic and nonmetallic systems alike may be treated by varying this cost.

The PDI model is parametrized by isotropic atom-type polarizabilities. These polarizabilities are coupled to each other by the interactions of atomic induced dipole moments from an external electric field. From the coupled atomic polarizabilities, the molecular polarizability is obtained. A large selection of molecular systems has been studied using this model. The new model uses charge distributions[15]-[19] rather than the point particles employed in the original PDI model.

The combined EEM/PDI model

In this section, we will combine the EEM and PDI models to form a set of equations for the determination of atomic charges and molecular dipole moments and polarizabilities. We will in the following consider an arbitrary molecular system of neutral charge comprised of N atoms. The parameters of the model will be superscripted with an asterisk for ease of identification. Einstein summation is adopted throughout. The model will

first be stated in a form which is equivalent to a combined EEM/PDI model. Following this, several modifications and considerations will be made, specifying the model in its final form.

The EEM model and PDI model is linked to the charge-charge and dipole-dipole interaction energy of the molecule, respectively. When combining these models, it is also necessary to consider charge-dipole interaction energy, as this will be the coupling between them. Therefore, neglecting higher-order interactions, the molecular energy V is stated to consist of the charge-charge, the charge-dipole and the dipole-dipole interaction energies V^{qq} , $V^{q\mu}$, and $V^{\mu\mu}$ as [5], [26]

$$V = V^{qq} + V^{q\mu} + V^{\mu\mu}. \quad (1)$$

This energy will later be subjected to minimization for the obtainment of optimal atom charge transfer terms and dipole moments. The charge-charge interaction energy is given by an additive contribution from each atom I , in which the unperturbed and uncharged energy V_I^0 and its perturbed/charged counterpart is separated according to

$$V^{qq} = \sum_I^N \left(V_I^0 + (\chi_I^* + \varphi_I^{\text{ext}})q_I + \frac{1}{2}\eta_I^*q_I^2 + \frac{1}{2} \sum_{J \neq I}^N q_I T_{IJ}^{(0)} q_J \right), \quad (2)$$

where χ_I^* and η_I^* are atomic parameters for electronegativity and chemical hardness, respectively, φ_I^{ext} denotes interaction with an external electrostatic potential, q_I is atomic charge resulting from charge transfer from other atoms. The atom-pair summation term contains Coulomb electrostatic interaction, which by classical electrostatics implies that $T_{IJ}^{(0)} = \frac{1}{R_{IJ}}$, where R_{IJ} is the Euclidian distance between atoms I and J . Ignoring the constant term V_I^0 and taking $T_{II}^{(0)} = \eta_I^*$, the tidier

$$V^{qq} = \sum_I^N \left((\chi_I^* + \varphi_I^{\text{ext}})q_I + \frac{1}{2} \sum_J^N q_I T_{IJ}^{(0)} q_J \right) \quad (3)$$

is obtained. The EEM/PDI coupling term and charge-dipole interaction energy is stated by combining atomic charges and dipoles as

$$V^{q\mu} = \sum_{I,J}^N q_I T_{IJ,\alpha}^{(1)} \mu_{J,\alpha}, \quad (4)$$

where α denotes a Cartesian axis, $\mu_{J,\alpha}$ is the corresponding dipole moment of atom J , and $T_{IJ,\alpha}^{(1)}$ is the charge-dipole interaction, which is the Cartesian gradient of $T_{IJ}^{(0)}$. Finally, the dipole-dipole interaction energy combines atomic dipoles by

$$V^{\mu\mu} = \frac{1}{2} \sum_J^N \mu_{J,\alpha} \alpha_{J,\alpha\beta}^{-1} \mu_{J,\beta} - \frac{1}{2} \sum_J^N \sum_{K \neq J}^N \mu_{J,\alpha} T_{JK,\alpha\beta}^{(2)} \mu_{K,\beta} - \sum_J^N E_{J,\alpha}^{\text{ext}} \mu_{J,\alpha}, \quad (5)$$

where α and β are Cartesian axes, $T_{JK,\alpha\beta}^{(2)}$ is a dipole-dipole interaction given by the gradient of $T_{IJ,\alpha}^{(1)}$, $E_{J,\alpha}^{\text{ext}}$ is an external electric field at atom J , and $\alpha_{J,\alpha\beta}$ is an atomic polarizability. In the above expression, this polarizability could in principle be anisotropic, although in the present work an isotropic atomic polarizability α_J^* is used. This choice will generally sacrifice some accuracy in the anisotropic parts of the representation of the molecular polarizability, but is nevertheless adopted for simplicity.

This is a convenient point to rephrase the model to more easily address the notion of charge transfer [21]. The starting point for this is to identify the atomic charge q_I as arising from a sum of terms of the type q_{IJ} , which is the charge transferred to atom I from some atom J in the system. This is done by noting that

$$q_I = \sum_J^N L_{IJ} q_{IJ}, \quad (6)$$

where L_{IJ} is topology matrix determining if charge transfer is allowed between atoms I and J , where $L_{IJ} = 1$ if charge transfer is allowed and 0 if it is not [21]. This matrix serves the purpose of reducing the number of equations considered in the expressions to be presented later: As long as there for any pair of atoms exists a ‘‘charge transfer pathway’’; a pathway through other atoms for the transfer of charge for which the corresponding topology matrix elements are equal to 1, the same results are obtainable as for a model where a more liberal permission of pairwise charge transfer were assigned. The charge transfer terms will replace the atomic charges both in the expressions to come, and they will also take the role as variables under consideration rather than the atomic charges. This does not hamper the calculation of the atomic charges in any way, as they by (6) are readily determined from the appropriate charge transfer terms.

At this point we make a remark concerning the conservation of charge. In the beginning of this section it was stated that the system considered is a molecule of neutral charge. In both this case and one where the molecule carries a non-neutral charge, the conservation of charge can be dealt with by adding a Lagrange multiplier term of the type $-\lambda(q^{mol} - \sum_I^N q_I)$ to eqns. (2) and (3), where λ is identified as the molecular chemical potential. However, for neutral molecules, this conservation will automatically be included by requiring $q_{JI} = -q_{IJ}$ and $q_{II} = 0$, and so we adopt this and take the Lagrange multiplication out of consideration.

As was stated right after eqn. (1), the molecular energy V will be subjected to minimization in order to obtain the optimal atomic charge transfer values q_{SP} (for two atoms S and P) and dipole moments $\mu_{I,\alpha}$. The differentiation of V with respect to these quantities yields the expressions

$$\frac{\partial V}{\partial q_{SP}} = \frac{\partial V^{qq}}{\partial q_{SP}} + \frac{\partial V^{q\mu}}{\partial q_{SP}} = 0 \quad (7)$$

and

$$\frac{\partial V}{\partial \mu_{I,\alpha}} = \frac{\partial V^{q\mu}}{\partial \mu_{I,\alpha}} + \frac{\partial V^{\mu\mu}}{\partial \mu_{I,\alpha}} = 0, \quad (8)$$

where $S > P$, both set to zero for optimization. In the following, we will rephrase the energy contributions to conform to the the charge-transfer viewpoint, make some other notational simplifications, and finally set up some matrix equations for the determination of the molecular charges, dipole moments and polarizabilities.

The charge-charge interaction energy is rewritten in the charge transfer viewpoint by combining (3) and (6) into

$$V^{qq} = \sum_{I,J}^N (\chi_I^* + \varphi_I^{\text{ext}}) L_{IJ} q_{IJ} + \frac{1}{2} \sum_{I,J,K,M}^N L_{IK} q_{IK} T_{IJ}^{(0)} L_{JM} q_{JM}. \quad (9)$$

We note that

$$\sum_{I,J}^N (\chi_I^* + \varphi_I^{\text{ext}}) L_{IJ} q_{IJ} = \sum_{I,J}^N (\chi_J^* + \varphi_J^{\text{ext}}) L_{JI} q_{JI}, \quad (10)$$

and introduce $\chi_{IJ}^* = \chi_I^* - \chi_J^*$ and $\varphi_{IJ}^{\text{ext}} = \varphi_I^{\text{ext}} - \varphi_J^{\text{ext}}$, which, when combined with the molecular charge conservation properties $q_{JI} = -q_{IJ}$ and $q_{II} = 0$ and the topology matrix element property $L_{IJ} = L_{JI}$ yields

$$\begin{aligned} \sum_{I,J}^N (\chi_I^* + \varphi_I^{\text{ext}}) L_{IJ} q_{IJ} &= \frac{1}{2} \sum_{I,J}^N ((\chi_I^* + \varphi_I^{\text{ext}}) L_{IJ} q_{IJ} + (\chi_J^* + \varphi_J^{\text{ext}}) L_{JI} q_{JI}) \\ &= \frac{1}{2} \sum_{I,J}^N (\chi_{IJ}^* + \varphi_{IJ}^{\text{ext}}) L_{IJ} q_{IJ} = \sum_{I,J>I}^N (\chi_{IJ}^* + \varphi_{IJ}^{\text{ext}}) L_{IJ} q_{IJ}. \end{aligned} \quad (11)$$

The original expression for V^{qq} then becomes

$$V^{qq} = \sum_{I,J>I}^N (\chi_{IJ}^* + \varphi_{IJ}^{\text{ext}}) L_{IJ} q_{IJ} + \frac{1}{2} \sum_{I,J,K,M}^N L_{IK} q_{IK} T_{IJ}^{(0)} L_{JM} q_{JM}. \quad (12)$$

The expression below gives charge transfer differentiation of V^{qq} for (7), and uses the properties $\frac{\partial V^{qq}}{\partial q_{SP}} = -\frac{\partial V^{qq}}{\partial q_{PS}}$ and $T_{IJ}^{(0)} = T_{JI}^{(0)}$, where $S > P$, to yield

$$\frac{\partial V^{qq}}{\partial q_{SP}} = L_{SP} \left(\chi_{SP}^* + \varphi_{SP}^{\text{ext}} + \sum_{J,M}^N (T_{SJ}^{(0)} - T_{PJ}^{(0)}) L_{JM} q_{JM} \right). \quad (13)$$

Defining

$$T_{SP,JM}^{(0)} = (T_{SJ}^{(0)} - T_{PJ}^{(0)}) - (T_{SM}^{(0)} - T_{PM}^{(0)}) \quad (14)$$

and noting that $\frac{\partial V^{qq}}{\partial q_{SP}}$ can be nonzero only when $L_{SP} \neq 0$, leaving in consideration only atom pairs between which charge transfer is allowed (so that $L_{SP} = 1$), eqn. (13) is rewritten as

$$\begin{aligned} \frac{\partial V^{qq}}{\partial q_{SP}} &= \chi_{SP}^* + \varphi_{SP}^{\text{ext}} + \sum_{J,M}^N (T_{SJ}^{(0)} - T_{PJ}^{(0)}) L_{JM} q_{JM} \\ &= \chi_{SP}^* + \varphi_{SP}^{\text{ext}} + \sum_{J,M>J}^N (T_{SJ}^{(0)} - T_{PJ}^{(0)}) - (T_{SM}^{(0)} - T_{PM}^{(0)}) L_{JM} q_{JM} \\ &= \chi_{SP}^* + \varphi_{SP}^{\text{ext}} + \sum_{J>M}^N T_{SP,JM}^{(0)} L_{JM} q_{JM}, \end{aligned} \quad (15)$$

yielding, by the final expression, a compact form. It is noted that $T_{SP,JM}^{(0)} = -T_{PS,JM}^{(0)} = -T_{SP,MJ}^{(0)}$, from which it can be seen that the resulting charge fulfills $q_{JM} = -q_{MJ}$, leaving only the cases where $S < P$ and $J < M$ to be considered.

The appropriate charge-transfer terms replace the atom charges in eqn. (4) to give

$$V^{q\mu} = \sum_{I,K,J}^N L_{IK} q_{IK} T_{I,J,\alpha}^{(1)} \mu_{J,\alpha} = \sum_{I,K>I,J}^N L_{IK} q_{IK} (T_{I,J,\alpha}^{(1)} - T_{K,J,\alpha}^{(1)}) \mu_{J,\alpha}, \quad (16)$$

Introducing

$$T_{SP,J,\alpha}^{(1)} = T_{SJ,\alpha}^{(1)} - T_{PJ,\alpha}^{(1)}, \quad (17)$$

and minimizing eqn. (16) with respect to q_{SP} , gives

$$\begin{aligned} \frac{\partial V^{q\mu}}{\partial q_{SP}} &= L_{SP} \sum_J^N (T_{SJ,\alpha}^{(1)} - T_{PJ,\alpha}^{(1)}) \mu_{J,\alpha} \\ &= L_{SP} \sum_J^N T_{SP,J,\alpha}^{(1)} \mu_{J,\alpha}, \end{aligned} \quad (18)$$

again yielding a compact form. Using the symmetry $T_{JS,\alpha}^{(1)} = -T_{SJ,\alpha}^{(1)}$, introducing

$$T_{J,SP,\alpha}^{(1)} = T_{SJ,\alpha}^{(1)} - T_{PJ,\alpha}^{(1)} = -\left(T_{JS,\alpha}^{(1)} - T_{JP,\alpha}^{(1)}\right) = \left(T_{SP,J,\alpha}^{(1)}\right)^T, \quad (19)$$

in analogy with eqn. (18), another compact form is obtained as

$$\begin{aligned} \frac{\partial V^{q\mu}}{\partial \mu_{J,\alpha}} &= \sum_{S,P>S}^N (T_{SJ,\alpha}^{(1)} - T_{PJ,\alpha}^{(1)}) L_{SP} q_{SP} \\ &= -\sum_{S,P>S}^N (T_{SJ,\alpha}^{(1)} - T_{PJ,\alpha}^{(1)}) L_{SP} q_{SP} \\ &= -\sum_{S,P>S}^N T_{J,SP,\alpha}^{(1)} L_{SP} q_{SP}. \end{aligned} \quad (20)$$

The remaining differentiation is given by

$$\frac{\partial V^{\mu\mu}}{\partial \mu_{J,\alpha}} = \alpha_{J,\alpha\beta}^{-1} \mu_{J,\beta} - \sum_{K \neq J}^N T_{JK,\alpha\beta}^{(2)} \mu_{K,\beta} - E_{J,\alpha}^{\text{ext}}. \quad (21)$$

Finally, defining $T_{II,\alpha\beta}^{(2)} = (\alpha_I^*)^{-1} \delta_{\alpha\beta}$, where α_I^* is an atom-type isotropic polarizability parameter and $\delta_{\alpha\beta}$ is the Kronecker delta, eqns. (15), (18), (20) and (21) are inserted into eqns. (7) and (8) to obtain the matrix equation

$$\begin{bmatrix} T_{SP,JM}^{(0)} & T_{SP,J,\alpha}^{(1)} \\ T_{I,JM,\alpha}^{(1)} & T_{IJ,\alpha\beta}^{(2)} \end{bmatrix} \begin{bmatrix} q_{JM} \\ \mu_{J,\beta} \end{bmatrix} = \begin{bmatrix} -\chi_{SP}^* - \varphi_{SP}^{\text{ext}} \\ E_{I,\alpha}^{\text{ext}} \end{bmatrix}, \quad (22)$$

which brings together all the results so far. For the parts of the matrix governed by atom pairs, only pairs S, P and J, M where charge transfer is allowed are included, since otherwise the corresponding matrix elements would be zero as L_{SP} or L_{JM} , respectively, would be zero. Hence and henceforth, all references to and summations over atom pairs will be implicitly understood to only include atom pairs between which charge transfer is allowed, so that $\sum_{I,J}$ means a sum over all pairs where $L_{IJ} = 1$ and $\sum_{I,J>I}$ means

a sum over all pairs so that $J > I$ and $L_{IJ} = 1$. Such pairs will also be referred to as “active pairs”.

Let P_a denote the number of active pairs in the molecular system. Recalling that N is the number of atoms in the system, the matrix in eqn. (22) is then of size $(P_a + 3N \times P_a + 3N)$. The number P_a will be equal to the number of bonds in the system if the topology matrix elements are set to unity only for atoms that are bonded, and zero otherwise. Adopting this topology matrix convention, the size of the matrix will be some low multiple of N , yielding a satisfactory scaling with system size.

Moving on to develop expressions for the calculation of the molecular dipole moment and polarizability, eqn. (22) is restated in field-independent and field-dependent forms. Let q_{JM}^0 and $\mu_{J,\alpha}^0$ denote a field-independent atom pair charge transfer term and atomic dipole moment, respectively. Disregarding any external electric interaction by setting $\varphi_{SP}^{\text{ext}}$ and $E_{I,\alpha}^{\text{ext}}$ equal to zero, the expression

$$\begin{bmatrix} T_{SP,JM}^{(0)} & T_{SP,J,\alpha}^{(1)} \\ T_{I,JM,\alpha}^{(1)} & T_{IJ,\alpha\beta}^{(2)} \end{bmatrix} \begin{bmatrix} q_{JM}^0 \\ \mu_{J,\alpha}^0 \end{bmatrix} = \begin{bmatrix} -\chi_{SP}^* \\ 0 \end{bmatrix}, \quad (23)$$

solvable for these quantities, is found. Once obtained, these field-independent quantities are readily inserted into the expression for the molecular dipole moment μ_α^{mol} as

$$\begin{aligned} \mu_\alpha^{\text{mol}} &= \sum_I (R_{I,\alpha} q_I^0 + \mu_{I,\alpha}^0) = \sum_{I,M} R_{I,\alpha} q_{IM}^0 + \sum_I \mu_{I,\alpha}^0 \\ &= \sum_{I,M>I} R_{IM,\alpha} q_{IM}^0 + \sum_I \mu_{I,\alpha}^0. \end{aligned} \quad (24)$$

In the above expression, $R_{I,\alpha}$ and $R_{IM,\alpha}$ denote Cartesian coordinate of atom I and interatomic distance between atoms I and M , respectively. For the field-dependent quantities $\partial q_{JM}/\partial E_\gamma^{\text{ext}}$ and $\partial \mu_{J,\alpha}/\partial E_\gamma^{\text{ext}}$, eqn. (22) is differentiated with respect to the electric field. It is assumed that the electric field is homogeneous so that $\varphi_I^{\text{ext}} = -R_{I,\alpha} E_\alpha^{\text{ext}} + C$, where C is a constant vanishing under differentiation. The resulting expression is

$$\begin{bmatrix} T_{SP,JM}^{(0)} & T_{SP,J,\alpha}^{(1)} \\ T_{I,JM,\alpha}^{(1)} & T_{IJ,\alpha\beta}^{(2)} \end{bmatrix} \begin{bmatrix} \partial q_{JM}/\partial E_\gamma^{\text{ext}} \\ \partial \mu_{J,\alpha}/\partial E_\gamma^{\text{ext}} \end{bmatrix} = \begin{bmatrix} R_{SP,\gamma} \\ \delta_{\beta\gamma} \end{bmatrix}. \quad (25)$$

The solution to eqn. (25) may be used to find the molecular polarizability $\alpha_{\alpha\beta}^{\text{mol}}$ from

$$\begin{aligned} \alpha_{\alpha\beta}^{\text{mol}} &= \frac{\partial \mu_\alpha^{\text{ind}}}{\partial E_\beta^{\text{ext}}} = \sum_I R_{I,\alpha} \frac{\partial q_I}{\partial E_\beta^{\text{ext}}} + \frac{\partial \mu_{I,\alpha}}{\partial E_\beta^{\text{ext}}} \\ &= \sum_{I,M>I} R_{IM,\alpha} \frac{\partial q_{IM}}{\partial E_\beta^{\text{ext}}} + \sum_I \frac{\partial \mu_{I,\alpha}}{\partial E_\beta^{\text{ext}}}. \end{aligned} \quad (26)$$

From eqns. (25) and (26) it is clear that the molecular polarizability is independent of the atomic electronegativity. This completes the presentation of the combined EEM/PDI model. In the following section, we will look at some modifications marking the departure from where the model is solely a combination of the EEM and PDI models.

Modifying the combined EEM/PDI model

Partly based on earlier work by the authors, several modifications to the model thus far presented have been introduced to deal with various effects.

We start with a topic mentioned in the introduction to the model: Using atomic charge distributions instead of point charges[15]-[19]. Using point charges is a simplification which may not be close enough to the true distribution of charge around an atom. By using some function to represent the atomic charge distribution, a more realistic representation can likely be attained. The atomic charge distribution considered for the model is based on a Gaussian function [18], which by the distance r_I from the nucleus is given as

$$q_I \left(\frac{\Phi_I^*}{\pi} \right)^{\frac{3}{2}} e^{-\Phi_I^* r_I^2}, \quad (27)$$

where Φ_I^* is an atom-type parameter describing the atomic propensity of radial charge density decline. The norming factor preceding the exponential term ensures that the total atomic charge is still q_I . By integration, the electrostatic interaction energy V_{IJ} between two such charge distributions is found as[27]

$$V_{IJ} = q_I q_J \frac{\text{erf}(\sqrt{a_{IJ}} R_{IJ})}{R_{IJ}}, \quad (28)$$

where

$$a_{IJ} = \frac{\Phi_I^* \Phi_J^*}{\Phi_I^* + \Phi_J^*}. \quad (29)$$

To more closely resemble the familiar electrostatic potential energy between two point charges, and for a convenient introduction of charge distributions into the EEM/PDI combination model, the notion of a scaled distance R_{IJ}^s is introduced. This scaled distance is given as

$$R_{IJ}^s = \frac{R_{IJ}}{\text{erf}(\sqrt{a_{IJ}} R_{IJ})}, \quad (30)$$

and may be inserted into (28) to give

$$V_{IJ} = \frac{q_I q_J}{R_{IJ}^s}. \quad (31)$$

Determining the error function may be computationally intensive and dependent on the accurate implementation of a procedure for its calculation. Therefore, an approximate scaled distance may be useful. Such an approximation should reproduce the error function with some accuracy and in addition show the same behavior at the limits $R_{IJ} \rightarrow 0$ and $R_{IJ} \rightarrow \infty$. This is accomplished by[18]

$$R_{IJ}^s = \sqrt{R_{IJ}^2 + \frac{\pi}{4a_{IJ}}}, \quad (32)$$

which is adopted hereafter. In the relay tensors $T_{IJ}^{(0)}$, $T_{IJ}^{(1)}$ and $T_{IJ}^{(2)}$ introduced in the previous section, the interatomic distances and their differentiation employs the scaled distance. Another concept considered is exchange interaction, for which a natural starting point for discussion is the Hartree-Fock separation of energy terms, where the energy V is written as a sum of the one-electron term H_I , the self-Coulomb and interparticle

Coulomb terms J_{II} and J_{IJ} and the self-exchange and interparticle exchange terms K_{II} and K_{IJ} , so that[1]

$$V = H_I + J_{II} + J_{IJ} + K_{II} + K_{IJ}, \quad (33)$$

where the self-exchange term cancels exactly half of the self-Coulomb term. From this, the self-exchange term is accounted for by the modification[28]

$$\frac{1}{2}\eta_I^* q_I^2 \leftarrow \frac{1}{4}\eta_I^* q_I^2 \quad (34)$$

in eqn. (2). However, this could just as easily have been done by doubling the value of η_I^* unless η_I^* is also involved in some other energy term. By requiring the Coulomb term in the limit of $R_{IJ} \rightarrow 0$ to be equal to the self-Coulomb term[29], the relation

$$\eta_I^* = \lim_{R_{IJ} \rightarrow 0} \frac{1}{R_{IJ}^s} = \sqrt{\frac{2\Phi_I^*}{\pi}} \quad (35)$$

achieves the desired coupling. Now, the interparticle exchange energy V_{IJ}^{exch} could in quantum terminology be approximatively regarded as proportional to the square of the overlap of the particle wavefunctions ψ , so that [30]

$$V_{IJ}^{\text{exch}} = C \langle \psi_I | \psi_J \rangle, \quad (36)$$

where C is an as of yet undetermined prefactor. An approximation of this expression suitable for use in this model is the overlap of the corresponding electronic charge distributions ρ , given by [31], [32]

$$V_{IJ}^{\text{exch}} = C \int \rho_I \rho_J d\tau, \quad (37)$$

which is also adopted. The electronic charge distribution can be regarded as the distribution of the n_I^{val} valence electrons of an atom I , which by analogy to (27) is given by

$$\rho_I = n_I^{\text{val}} \left(\frac{\Phi_I^*}{\pi} \right)^{\frac{3}{2}} e^{-\Phi_I^* r_I^2}. \quad (38)$$

Inserting this expression into (37) and carrying out the integration yields

$$V_{IJ}^{\text{exch}} = C n_I^{\text{val}} n_J^{\text{val}} \left(\frac{a_{IJ}}{\pi} \right)^{\frac{3}{2}} e^{-a_{IJ} r_{IJ}^2}. \quad (39)$$

The prefactor C is determined by another appeal to limiting behavior. By using that the interparticle exchange term approaches the self-exchange term at short distances, i.e. taking

$$-\frac{1}{2}\eta_I^* = -\sqrt{\frac{\Phi_I^*}{2\pi}}, \quad (40)$$

eqn. (39) can be rewritten without the prefactor as

$$V_{IJ}^{\text{exch}} = -n_I^{\text{val}} n_J^{\text{val}} \sqrt{\frac{a_{IJ}}{\pi}} e^{-a_{IJ} r_{IJ}^2}. \quad (41)$$

At this point, the charge transfer viewpoint is adopted so that the atomic charge, given by n_I^{val} and the effective nuclear charge from the nucleus and core electrons, Z_{eff} , is rewritten as

$$q_I = Z_{\text{eff}} + n_I^{\text{val}} = Z_{\text{eff}} + n_I^{(0)} + \sum_K q_{IK}, \quad (42)$$

where $n_I^{(0)} = -Z_{\text{eff}}$ is the number of valence electrons in the unperturbed atom. From this expression, n_I^{val} can be understood to be $n_I^{(0)}$ plus a perturbation $\sum_K q_{IK}$ representing the electrons given by charge transfer. Inserting this for n_I^{val} in eqn. (41) yields the final expression for V_{IJ}^{exch} as

$$\begin{aligned} V_{IJ}^{\text{exch}} &= - \sum_{I,J,K,M} \sqrt{\frac{a_{IJ}}{\pi}} e^{-a_{IJ}r_{IJ}^2} q_{IK} q_{JM} - 2 \sum_{I,J,M} \sqrt{\frac{a_{IJ}}{\pi}} e^{-a_{IJ}r_{IJ}^2} n_I^{(0)} q_{JM} \\ &= - \sum_{I,J} \sqrt{\frac{a_{IJ}}{\pi}} e^{-a_{IJ}r_{IJ}^2} n_I^{(0)} n_J^{(0)}, \end{aligned} \quad (43)$$

where the last term can be disregarded due to its independence of a charge transfer term, thereby vanishing in the differentiation (7). In the model, the exchange energy is taken into account by comparing with eqn. (9). From this it is seen that χ_J^* and $T_{IJ}^{(0)}$ is modified as

$$\chi_J^* \leftarrow \chi_J^* - 2 \sum_I \sqrt{\frac{a_{IJ}}{\pi}} e^{-a_{IJ}r_{IJ}^2} \quad (44)$$

and

$$T_{IJ}^{(0)} \leftarrow T_{IJ}^{(0)} - \sqrt{\frac{a_{IJ}}{\pi}} e^{-a_{IJ}r_{IJ}^2}. \quad (45)$$

In the introduction to this section, the limitation of the EEM was mentioned. It was stated that this behavior arises from the fact that the EEM allows charge to flow freely in the molecule, thereby closely resembling the behavior of metallic system. As this will introduce a source of error in nonmetallic systems, there is a potential increase in accuracy to be gained by addressing this limitation. The metallic behavior can be stated as the tendency for charge transfer to occur between two particles at infinite separation. In the EEM, such charge transfer will take place if the particles are different, thereby showing a difference in electronegativity, or if an external electrostatic potential difference is present, such as that generated by a homogeneous electric field [33]. This behavior remains the same for infinitely long molecular chains as charge is still allowed to flow freely, thereby resulting in polarizabilities approaching infinity as the chain length approaches infinity.

This behavior can also be demonstrated in in equation form by applying equation (7) to a system consisting of two atoms 1 and 2. Disregarding for the moment the topology matrix, this application yields one equation

$$(\eta_1^* + \eta_2^* - 2T_{12}^{(0)})q_{12} = -(\chi_{12}^* + \varphi_{12}^{\text{ext}}). \quad (46)$$

As the relation $q_{12} = q_1 = -q_2$ must hold for this system, eqn. (46) can be rearranged to give

$$q_1 = -q_2 = \frac{-\chi_{12}^*}{\eta_1^* + \eta_2^* - 2T_{12}^{(0)}} + \frac{-\varphi_{12}^{\text{ext}}}{\eta_1^* + \eta_2^* - 2T_{12}^{(0)}}, \quad (47)$$

separated into a permanent atomic charge term and a contribution from a homogeneous external electric field, respectively. In the long distance limit the relay tensor element will tend to zero, giving

$$q_1 = -q_2 = \frac{-\chi_{12}^*}{\eta_1^* + \eta_2^*} + \frac{-\varphi_{12}^{\text{ext}}}{\eta_1^* + \eta_2^*}, \quad (48)$$

from which the metallic effect is demonstrated. In the case of a nonzero external electric field, the second term will contribute to the polarizability. A molecular chain with a hetero atom at one end under the same conditions (but under the restrictions imposed by a topology matrix) will also give a nonzero charge transfer between the end atoms.

As was mentioned in the introduction, this behavior is addressed by adding a charge transfer energy cost or “penalty” term in order to describe both metallic and non-metallic systems. In the search for adequate functional terms to accommodate this, the quantum mechanical Mulliken approach [34] provides a convenient starting point. Here, the charge transfer term q_{IJ} between two different atoms I and J is given by regarding the sets of basis functions B_I and B_J restricted to atoms I and J , respectively, by the expression

$$q_{IJ} = \sum_{i \in B_I} \sum_{j \in B_J} D_{ij} S_{ij}, \quad (49)$$

where D_{ij} and S_{ij} are density and overlap matrix elements, respectively. In this context, the important information in this equation is the overlap S_{ij} . This function declines exponentially as the distance R_{IJ} increases, thereby bringing about a decline in the charge transfer term as well. This kind of behavior could introduce nonmetallicity to the EEM. By regarding eqn. (48) and assuming the use of a function declining exponentially with interatomic distance, this can be approached by modifying the chemical hardness so that [24], [25]

$$(\eta_1^* + \eta_2^*) \leftarrow (\eta_1^* + \eta_2^*) S_{12}^{-1}, \quad (50)$$

thereby adding an exponentially increasing energy cost to the charge transfer as the interatomic separation increases. Recalling the chemical hardness to be understood as a measure of the atomic “reluctance” to give away its original electronegativity, this modification could be said to retain the essential nature of this parameter. For a discussion of a general modification to the chemical hardness, it is convenient to take as a starting point an energy term of the form

$$\frac{1}{2} \eta_I^* q_{IK} q_{IM}, \quad (51)$$

an atomic charge equivalent of which is found in eqn. (2) and appearing in “disguise” in charge transfer form in eqn. (9) by the convention adopted for $T_{II}^{(0)}$. A candidate for introducing a general modification of the chemical hardness is modifying eqn. (51) as

$$\frac{1}{2} \eta_I^* q_{IK} q_{IM} \leftarrow \frac{1}{2} \eta_I^* S_{IK}^{\frac{1}{2}} S_{IM}^{\frac{1}{2}} q_{IK} q_{IM}. \quad (52)$$

There remains to determine a functional form for the overlap term. In keeping with the Gaussian “convention” adopted for charge distribution and exchange terms, a candidate for such a form may be

$$S_{IJ} = e^{-a_{IJ}R_{IJ}^2}, \quad (53)$$

This will give the desired behavior for the two-atom system considered above and has been proposed elsewhere [23]. Now, turning to the other system under consideration, a molecular chain of length approaching infinity, it is seen that the modification by eqn. (52) does not suffice for eliminating the problem of the occurrence of charge transfer. Two atoms at each end of the chain will be connected by a series of atoms for which a “pathway” of nonzero topology matrix elements exists. If consequent atoms in the chain are equidistant, their overlap will be equal in this “pathway” between the end atoms, and the energy cost desired to happen could be nullified simply by scaling η_I^* . The consequence is that charge transfer between the end atoms will still take place, which is an unacceptable result. It is therefore necessary to look at another form of an energy cost term modifying the chemical hardness. Such a form is the expression

$$\frac{1}{2}\eta_I^*q_{IK}q_{KM} \leftarrow \frac{1}{2}\epsilon\eta_I^*q_{IK}q_{KM}, \quad (54)$$

retaining the form of the energy term used in (51). In this expression, ϵ is real and positive. Now, for determining a desirable form of ϵ , a three-particle system consisting of atoms I , K and M is considered. In this system, atoms K and M both connected to the central atom I but are not connected to each other (so that $L_{KM} = 0$). In this system, charge may be transferred between atoms K and I and atoms I and M . Now, if this system is part of a large molecular chain, the transport cost modified ϵ will be applied repeatedly as charge attempts to flow from one end atom to another. By finding a method of tuning ϵ so that metallic and nonmetallic systems both can be accommodated, a suitable modification can be obtained.

The chemical hardness contribution to the energy, V^η , in the charge transfer phrasing (the “disguised” term of eqn. (9)) corresponding to the modification of eqn. (54), is given by

$$V^\eta = \frac{1}{2}\eta_K^*q_{KI}^2 + \frac{1}{2}\eta_I^*(q_{KI}^2 - 2(1-\epsilon)q_{KI}q_{IM} + q_{IM}^2) + \frac{1}{2}\eta_M^*q_{IM}^2, \quad (55)$$

which, differentiated with respect to q_{KI} and q_{IM} yields

$$\frac{\partial V^\eta}{\partial q_{KI}} = (\eta_K^* + \eta_I^*)q_{KI} - (1-\epsilon)\eta_I^*q_{IM} \quad (56)$$

and

$$\frac{\partial V^\eta}{\partial q_{IM}} = (\eta_M^* + \eta_I^*)q_{IM} - (1-\epsilon)\eta_I^*q_{KI}, \quad (57)$$

respectively. From the last terms of the above two expressions, it is seen that ϵ can not be larger than unity, as this will result in an effectively negative chemical hardness for atom I, restricting ϵ to the interval (0, 1]. By comparing this interval to eqn. (54), it can be stated that nonmetallic behavior should be observed when ϵ is close to 1, and conversely, metallic behavior should be consistent with ϵ being close to 0.

An important type of system for which the model should produce accurate results is carbon chains. Depending on the hybridization state of the atoms comprising such a chain, the interatomic bonds will have a varying degree of σ and π bonding, and this has significant effect on the ease with which charge is transferred through the chain. Now, carbon-carbon bonds may largely be said to be single, double or triple, disregarding any coordination perturbations from other atoms. As these bonds will have different lengths,

albeit only a small difference, a function dependent on the interatomic distances could be able to discriminate between them, thus also being able to scale the “metallicity” by linking such a function to ϵ . In the original article, such a function is introduced as $g_{I,KM}(R_{IK}, R_{IM})$, for which the subscript should be taken to mean “concerning atom I , involving atoms K and M ”. In this article, the functional form of $g_{I,KM}$ is designed for the discrimination between carbon single and double bonds, thereby leaving the description of triple bonds out of consideration. This function is introduced together with the overlap functions by a final modification along the lines of eqn. (52) as

$$\frac{1}{2}\eta_I^*q_{IK}q_{IM} \leftarrow \frac{1}{2}\eta_I^*S_{IK}^{\frac{1}{2}}S_{IM}^{\frac{1}{2}}g_{I,KM}q_{IK}q_{IM}. \quad (58)$$

Finally, this modification is linked to ϵ by requiring that $\epsilon = 1 - g_{I,KM}$. This means that $g_{I,KM}$ should approach 1 for metallic systems, and, conversely, approach 0 for nonmetallic systems. As doubly bonded carbon atoms, showing a significant amount of π bonding, will be more conducive to charge transport, $g_{I,KM}$ for which at least one of the bonds is a carbon-carbon double bond should be closer to 1 than a corresponding system where this carbon-carbon bond is single. A declining exponential function could be able to produce values varying sensitively within the interval $(0, 1]$ based on bond order distinction provided that a suitable exponent is supplied. In the original article, a comparison of bond distances in conjugated alkene molecules of varying sizes was used to demonstrate that the variation between carbon single and double bond distances is high compared to the sum of two adjacent interatomic distances in those carbon chains. As a conjugated alkene system by experience should transport charge much more readily than a corresponding alkane system, the exponent of $g_{I,KM}$ should still provide the metallic behavior even when one of the carbon-carbon bonds under consideration is single if the other carbon-carbon bond is double, as will be the case for all such terms when all atoms I , K and M are carbon atoms. For this reason, the functional form

$$g_{I,KM} = e^{-C(R_{IK}+R_{IM}-2R_I^*-R_K^*-R_M^*)^2}, \quad (59)$$

is proposed, where R_I^* is an atom-type parameter for atom I and C is a global parameter. The sum $R_I^* + R_J^*$ is intended to be approximately equal to the notion of an “equilibrium” bond length between atoms I and J , while C is intended to govern the ability of $g_{I,KM}$ to switch rapidly between 1 and 0 when the the sum of the actual length of two adjacent bonds KI and IM deviates from the sum of the equilibrium bond lengths.

The final modification given in the original article uses the newly adopted parameters R_I^* to model the interatomic overlap, reminiscent of eqn. (53) as

$$S_{IJ} = e^{-a_{IJ}(R_{IJ}-R_I^*-R_J^*)^2} \quad (60)$$

if $R_{IJ} > (R_I^* + R_J^*)$, and $S_{IJ} = 1$ otherwise. This is because the overlap will be close to unity when the actual bond length is close to the equilibrium value, while decreasing exponentially with further interatomic departure.

The final modification given in eqn. (58) is incorporated into the model by modifying $T_{SP,JM}^{(0)}$ of eqn. (14). This tensor is only modified when one or more of the relations $S = J$, $P = M$, $K = J$ or $K = M$ holds. Assuming all indices J , M , S and P are different, the modified tensor is given by

$$\begin{aligned}
T_{SP,SM}^{(0)} &= T_{SS}^{(0)} S_{SP}^{-1} S_{SM}^{-1} g_{S,PM} - T_{PS}^{(0)} - T_{SM}^{(0)} - T_{PM}^{(0)} \\
T_{SP,JS}^{(0)} &= T_{SJ}^{(0)} - T_{PJ}^{(0)} - T_{SS}^{(0)} S_{SP}^{-1} S_{JS}^{-1} g_{S,PJ} + T_{PS}^{(0)} \\
T_{SP,PM}^{(0)} &= T_{SP}^{(0)} - T_{PP}^{(0)} S_{PS}^{-1} S_{PM}^{-1} g_{P,SM} - T_{SM}^{(0)} + T_{PM}^{(0)} \\
T_{SP,JP}^{(0)} &= T_{SJ}^{(0)} - T_{PJ}^{(0)} - T_{SP}^{(0)} + T_{PP}^{(0)} S_{PS}^{-1} S_{PJ}^{-1} g_{P,SJ} \\
T_{SP,SP}^{(0)} &= T_{SS}^{(0)} S_{SP}^{-1} S_{SP}^{-1} - T_{PS}^{(0)} - T_{SP}^{(0)} + T_{PP}^{(0)} S_{SP}^{-1} S_{SP}^{-1} \\
T_{SP,PS}^{(0)} &= T_{PS}^{(0)} - T_{SS}^{(0)} S_{SP}^{-1} S_{SP}^{-1} - T_{PP}^{(0)} S_{SP}^{-1} S_{SP}^{-1} + T_{SP}^{(0)} \\
T_{SP,JM}^{(0)} &= T_{SJ}^{(0)} - T_{PJ}^{(0)} - T_{SM}^{(0)} + T_{PM}^{(0)}.
\end{aligned} \tag{61}$$

This marks the end of the presentation of the model as it is given in the original article [2]. During the course of the present work, a proposed modification has been introduced. This modification is based on the idea that the propensity of charge transfer between atoms of different elements may vary. As the parameter C is presented as a global parameter, it will to some extent dictate the ability of every bond to stop charge from being transferred, treating bonds between any two elements in the same way. To illustrate why this may introduce errors, consider the end of an arbitrarily long molecular chain with a heteroatom substituent at one end. Suppose that the heteroatom actually behaves in such a way that charge transfer is effectively stopped between the heteroatom and the first atom of the chain. Now, let the chain be comprised of carbon atoms, as will commonly be the case. Consider the two cases where the first carbon-carbon chain bond, adjacent to the carbon-hetero atom bond, is either a single or double bond. Depending on the value of R_C^* , where the subscripted C denotes carbon, the parentheses term of the exponent in $g_{I,KM}$ could by circumstance approach zero for one of these cases (creating the opposite effect for the other case), thereby increasing the tendency of $g_{I,KM}$ to take a value close to unity and impose a “locally metallic” behavior for that case, when in fact this behavior would run contrary to the supposed behavior of the heteroatom. As the global parameter C may, for sake of producing satisfactory results for molecules in which this heteroatom behavior is not observed, may take a small value, its effect on the exponent may be too weak to produce the desired reduction of $g_{I,KM}$ to a value near zero. To counteract this, a modified form of eqn. (59) is proposed, in which

$$g_{I,KM} = e^{-(C_I^*)^2 C_K^* C_M^* (R_{IK} + R_{IM} - 2R_I^* - R_K^* - R_M^*)^2}, \tag{62}$$

where C_I^* is another atom-type parameter. In this way, elements behaving in the manner supposed for the hetero atom in the previous paragraph will have the ability to locally block charge transfer more effectively, while effects for other elements could still be accommodated by the tuning of their corresponding parameters. This modification is adopted in the present work. A risk of introducing this new atom-type parameter is that there will be a strong coupling between the parameters R_I^* and C_I^* , so that in an optimization, one may become unreasonably large while the other becomes unreasonably small.

To recapitulate the parametrization, the original model uses five atom-type parameters: The electronegativity χ_I^* , the chemical hardness η_I^* , the (isotropic) polarizability α_I^* , the charge distribution radial decline term Φ_I^* and the parameter R_I^* used for the overlap and $g_{I,KM}$. In addition, there is one global parameter C , used in $g_{I,KM}$, which in the present work is substituted for a set of atom-type parameters C_I^* . The molecular polarizability is independent of the electronegativity parameter, and can therefore be taken out of consideration in that regard. This marks the end of the theory concerning the model and its features. In the next section we will look at the genetic algorithm.

2.2 Genetic algorithms

Introduction

In this section, a presentation of the genetic algorithm is given. The theory concerning genetic algorithms presented here is given in a review chapter [14]. The notation used in the previous section detailing the extended EEM/PDI combination model is disregarded in this section. We start off with a brief primer on the subject of optimization, wherein the genetic algorithm is introduced, and move on to walking through the steps taken in a genetic algorithm. Some of the problems facing the genetic algorithm will be discussed, and various extensions that could mitigate some of these weaknesses and enhance the overall performance will be described, whereby a genetic algorithm may be tailored to some level of specialization to suit the situation under study.

A genetic algorithm (GA) is an optimization method, that is, a method for locating parameters which are optimal with respect to some situation of interest. Genetic algorithms are named for their resemblance to biological evolution in their method of functioning, and the majority of concepts from which GAs are constructed are in close analogy to some biological counterpart. A genetic algorithm is regarded as potentially powerful tool for searching parameter space.

Before proceeding, let's make some general observations on optimization. The canonical scenario is that a system described by a set of parameters is given, and we would like to find the parameter configuration for which the system is in an optimal state with respect to some method of evaluation. Let the set of parameters be denoted by x , and let the evaluation of these parameters, i.e. the value of the quantity of interest, be represented by the function f , so that $f = f(x)$. Finally, let x^* be an optimal set of parameters for some f .

In order to treat the concepts introduced and their visualization in a consistent manner, we adopt the convention that optimal points in parameter space are characterized as minima of the function to be optimized. It is clear that any optimization of a function can always be seen from this viewpoint: Should the desired optimum be a maximum of some function f , the problem is trivially restated as locating the minimum of the negative $-f$. Now, more formally put, a local minimizer x^* of f is a set of parameters so that $f(x^*) \leq f(x)$ for all $x \in B(x^*; \epsilon)$, where $B(x^*; \epsilon)$ denotes a neighborhood of size $\epsilon > 0$ around x^* . We say that $f(x^*)$ is then a local minimum [35]. Now, there may be several local minima for one f , but only one global minimum. A global minimum is a local minimum such that no other local minima are of a lower value. A set of parameters corresponding to the global minimum is called a global minimizer, denoted by x^{**} . To avoid confusion, we remark that f may have several global minimizers, but only one global minimum. This is trivially demonstrated by the function $f(x) = \sin x$, which will have a global minimum of -1 , which is obtained at all global minimizers $x^{**} = \pi + 2k\pi, k \in \mathbb{Z}$.

While a large number of optimization paradigms are limited to the descent into a local minimum, a properly configured GA could locate the global minimum for the system at hand. In local minimization methods, the identification of a global minimum may be difficult if at all possible. For such a localization to succeed by means other than mere fortune, it is usually required to apply prior knowledge about the situation at hand, starting the method at a point believed to be close to the global minimum and hoping that the method will descend to the proper location. As many such searches may be needed, this may be difficult and time-consuming, to the limit of rendering such a search unfeasible. A GA's potential of locating the global minimum, on the other hand, is not dependent on some initial set of parameter values supplied by the user. However, it would be incorrect to assume that a GA requires no knowledge at all of the situation.

As will be detailed later, the algorithm is dependent on the specification of intervals in which each parameter may vary. Neglectful specification of these ranges may lead to sampling of absurd or “impossible” parts of parameter space.

Let’s move on to a brief description of the basic principles of genetic algorithms before treating these principles in detail. At the core of the GA is the population. The population is made up of a set of individuals. An individual corresponds to a particular arrangement of parameters, thereby corresponding to a particular configuration of the system under study. Now, the goal of the GA is to continually evolve the population, that is, the collection of tuples of parameters, into another population which, by some metric, is closer to the minimum of the evaluation function taking the parameters of the individuals as arguments. A particular population is called a generation, and the GA proceeds by creating another generation from the current one. The mechanisms by which this process takes place is essentially modelled on biological evolution.

These mechanisms can be divided into three main steps. The first step is the selection. In this process, the individuals in the current population are first evaluated according to a function. This function, called the fitness function, is the same function f that we wish to minimize. The function value associated with each individual, in other words the function value resulting from the parameters making up each individual, is called its fitness. When this evaluation is finished, a subset of the individuals in the population are chosen according to some set of selection rules. These rules should favor the most fit individuals, that is, the individuals whose parameters give the lowest value of f (recall that f can be and is stated as a function to be minimized - therefore, the most fit individuals would be the ones where the resulting value of f is lowest). This subset is called the breeding pool - the set of genetic material from which the next generation is created.

Next, individuals in the breeding pool are utilized to create offspring - the next generation of individuals. This stage is called recombination. Here, several strategies exist for producing offspring. Being covered in later sections and only mentioned here, most of these strategies are variations on combining the genetic material of one individual with that of another to produce another individual, in the same way that the genetic material of offspring in nature is inherited from its parents.

The last main step of the GA is mutation. In this stage, an alteration of the genetic material of one or more individuals is randomly introduced through some probabilistic rule. This step also finds a clear counterpart in biological evolution, wherein for most purposes, genetic mutation can be regarded as a random event. However, the rate of mutation introduced into a genetic algorithm is typically several orders of magnitude higher than that found in nature. The reason for including mutation is that it allows for a much vaster search of parameter space. In fact, mutation by its introduction of randomness ensures that, given enough time, the entire parameter space will be searched. It is this mechanism that gives genetic algorithms the potential to locate the global minimum.

The individual

The first step in presenting the genetic algorithm is taking a closer look at the way an individual is represented. As mentioned in the previous section, the individual consists of a set of parameters corresponding to a particular configuration of the system under study. This set of parameters is called the chromosome, taking the shape of a continuous string or tuple of individual parameter values. These individual values are called genes.

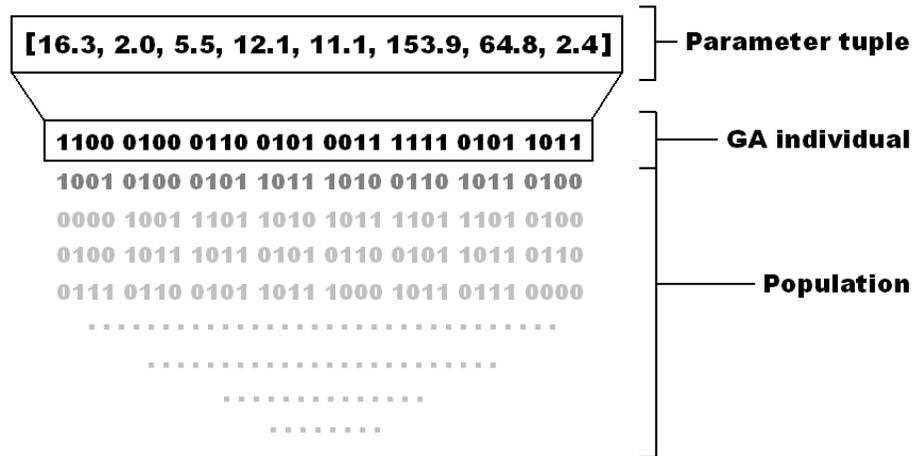


Figure 2.1: *The individual and the population in the genetic algorithm.*

The illustration above should give a clear impression of how the concepts discussed until now are related. From the illustration, it can be seen that individual genes are represented as binary strings. This is how parameters were first represented in genetic algorithms, and it has been adopted as a convention. As will be detailed later, there is also the possibility of representing parameters as real numbers, in effect using their actual values, but for the purpose of presenting the main features of the GA, it is convenient to retain the binary encoding convention. This binary encoding ensures that for each parameter, the values it may take can vary on a suitable interval of predetermined range and degree of discreteness.

This encoding begins by selecting an interval $[x_{min}, x_{max}]$ in which the parameter is allowed to vary. Next, the number of bits B to be used to represent each parameter is chosen. Including the start and end points x_{min} and x_{max} , the 2^B different values that a binary number of B bits can take are mapped onto 2^B equidistant points on the interval. Let I denote the collection of these points. Hence, by using a large number of bits for the representation, a high-resolution partition of the interval is achieved, and vice versa by using a small number of bits. The resolution can thus be scaled to suit one's own wishes. As the computational cost of the rest of the GA is usually thoroughly dwarfed by the cost of evaluating the fitness function, the number of bits chosen will have little direct impact on the computational expense, apart from the fact that convergence by some measures to be presented later may be slowed by using a large number of bits.

Taking d_i as the decimal value of a binary string of length B representing point x_i on the interval given above, the mapping from real interval to binary is given by the expression

$$d_i = (2^B - 1) \frac{x_i - x_{min}}{x_{max} - x_{min}} \quad (63)$$

and the corresponding reverse mapping is given by

$$x_i = x_{min} + (x_{max} - x_{min}) \frac{d_i}{2^n - 1} \quad (64)$$

Encoding chromosomes using the standard bit notation, the “powers-of-two” notation, may at first glance seem innocuous, but it carries some implications to the procedures for mutation. A typical mutation procedure, to be presented later, entails switching one randomly chosen bit of the chromosome. As such a procedure normally would not

distinguish between the significance of the bits used for each parameter, the mutation could give a drastic change in parameter values. This will usually be an unwanted consequence as the intention of mutation is introducing small local variations, and by this, the need arises for an encoding scheme in which such a mutation procedure wouldn't give such drastic consequences.

Table 2.1: Gray coding.

Gray code	Decimal	Standard binary
000	1	000
001	2	001
011	3	010
010	4	011
110	5	100
111	6	101
101	7	110
101	8	111

Such an encoding scheme is the so-called Gray coding [36]. This scheme is demonstrated above. Its main consequence is that the switching of one bit, in any position, will only shift the corresponding parameter value by one distance resolution unit as dictated by $[x_{min}, x_{max}]$ and B .

Initial settings, selection, recombination and mutation

Now that the parameter encoding scheme is taken care of, it is time to move on to the main mechanisms of the genetic algorithm. The first step is selecting the number of individuals N in the initial population. In most flavors of GA, this population size will remain constant during the course of the algorithm, but this is by no means required. The recommended population size is at least $20 \times P$, where P is the number of parameters, and at least 100 individuals in any case, should $P < 5$. Next, the genes of the initial population must be selected. As it is desired to perform a broad search of parameter space, the obvious choice for this is generating a population drawn from a uniform random distribution on I . Other choices of initial genetic makeup are of course allowed, for instance drawing from nonuniform distributions, but such choices would likely be too tailored to a specific situation to merit a mention in this general text.

Following this, the iterative part of genetic algorithm can commence. Recalling that one iteration of the GA is equivalent to the creation of a new generation of individuals, the terms "iteration" and "generation" are largely interchangeable. The first part of progressing through a generation is evaluating the fitness function for each individual in the population. As mentioned earlier, this will likely be the most time-consuming step as the fitness function in itself is usually sophisticated, and, furthermore, it must be calculated for every individual. When the fitness has been found for each individual, it is time to perform the three steps determining the genetic makeup of the next generation: Selection, recombination and mutation.

In selection, we recall that individuals deemed fit for building the next generation is put into a so-called breeding pool by some method or combination of methods. A typical size of the breeding pool is half the population size. There exist many selection methods, of which some are simple while others carry some degree of sophistication. The size of the breeding pool is generally not equal to the size of the population. Generally, it is desired to select the fittest individuals for breeding, but it may be that some individuals with low fitness in fact carry a desirable configuration in parts of their chromosome. As will be shown, some selection procedures take this possibility into consideration, while others

disregard it. Another important point, especially concerning selection, is that as the GA runs its course, it is generally desired that a balance is struck between maintaining genetic diversity and progressing towards some kind of convergence, the determination of which will be the subject of the next section. Therefore, care should be taken not to adopt procedures strongly inclined towards one of these extremes.

In the following, assume for convenience of explanation that the individuals have been put into a sorted list according to their fitness, where the most fit individuals are at the top of the list. An individual's placement in this sorted list is called its rank. A simple selection method is the so-called cutoff selection. In this method, some fraction of the individuals in the sorted list are simply put directly into the breeding pool.

Another simple method is tournament selection. Here, two individuals are chosen at random, and their fitness is compared. The individual with the highest rank is then put into the breeding pool. This ensures that on average, the most fit individuals will be selected, but individuals with low fitness also stand a chance of being selected from winning an "underdog fight", which will take place if the random selection happens to give two individuals both with low fitness. This procedure is continued until the desired breeding pool size is reached.

A third method is roulette selection. This method is in analogy to the spinning of a roulette wheel. Here, a random number is selected uniformly on some selection interval, typically $[0, 1]$ as this is the canonical interval for the drawing of a uniform random number. In this interval, each individual is represented as a subinterval of the selection interval. The size of each individual's subinterval is determined according to a proportioning rule. The random number selected will be in one of the subintervals thus defined, and the individual corresponding to this subinterval will be the one selected for the breeding pool. If the selection interval is conceived to be mapped onto a circle, each subinterval can then be thought of as a sector of this circle, thereby illustrating the analogy to a roulette wheel.

There are two well-known approaches to roulette selection. The first approach is the so-called fitness-based roulette selection. In this approach, the size of each subinterval is determined according to each individual's fitness, where individuals with higher fitness will receive a larger subinterval. A problem with fitness-based roulette selection is that in some populations, there may be just a few individuals with a very high fitness. As these individuals will receive a very large subinterval compared to other individuals and thus stand a very high likelihood of being selected for breeding, the genetic diversity of the population may quickly be reduced dramatically, thereby hampering the progression of the algorithm and bringing about premature convergence. In conclusion, using fitness-based roulette selection alone can prove to be a poor choice.

An alternative to fitness-based roulette selection is rank-based roulette selection. This method will size each subinterval according to the corresponding individual's fitness rank, where higher-ranking individuals will receive a larger subinterval. Selection by rank does away with the main problem of selecting purely by fitness, as the size of the subintervals will not be subject to such drastic differences as could be experienced in the latter. Comparing rank-based with fitness-based roulette selection, one would expect that the former should retain genetic diversity for longer.

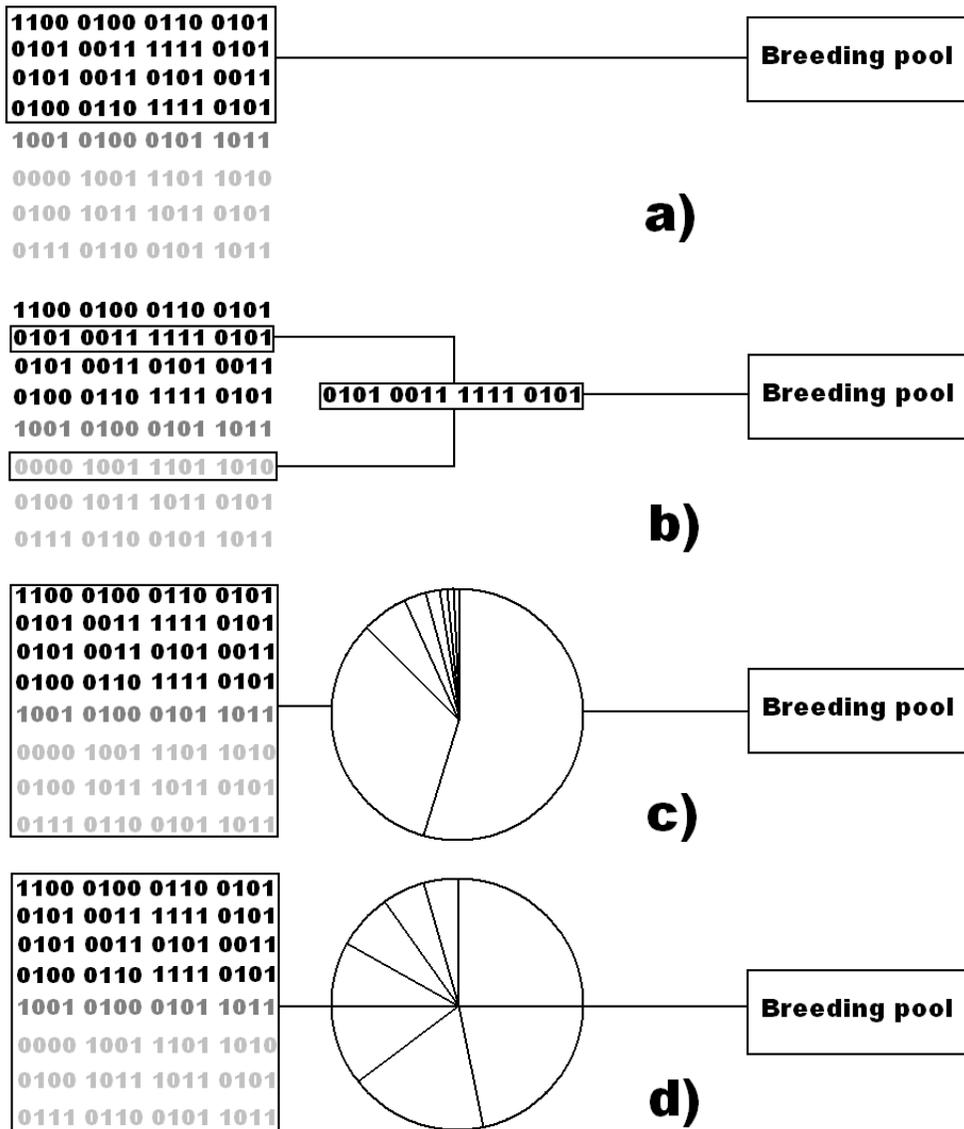


Figure 2.2: Selection procedures in a genetic algorithm: a) Cutoff selection. b) Tournament selection. c) Fitness-based roulette selection. d) Rank-based roulette selection.

The illustration in figure 2 above provides a summary of the selection procedures. After selection is completed it is time for recombination. In this stage, the genetic material of the individuals in the breeding pool is combined in some way to produce new individuals.

The simplest recombination procedure is single-point crossover. In this procedure, two individuals are chosen at random from the breeding pool. Next, a crossover point is chosen as a number N_C between 1 and the $L_C - 1$, where L_C is the bitwise length of the chromosome. Then, all bits up to and including position N_C from the first individual are combined with all bits succeeding position N_C from the second individual to form a new individual. The two individuals thus combined may be called the parents, and the resulting, recombined individual the child.

The crossover concept can be extended to include any number of crossover points. For instance, in 2-point crossover, two crossover points N_{C_1} and N_{C_2} are chosen so that

$N_{C_2} > N_{C_1}$. Then, all bits up to and including N_{C_1} are taken from parent 1, the succeeding bits up to and including N_{C_2} are taken from parent 2, and the final bits are taken from parent 1 to form a child chromosome. This general procedure is called n-point crossover.

A generalization of the crossover scheme is uniform crossover. In this procedure, a string of bits the length of the chromosome is created. In this string, dubbed the mask, bits are randomly assigned as 0 or 1. The new individual then takes the bit values corresponding to the positions where the mask bit is 0 from parent 1, and the values corresponding to the positions where the mask bit is 1 from parent 2, putting each bit value in the position from which it was originally taken. The recombination procedures presented here are illustrated in figure 3 below.

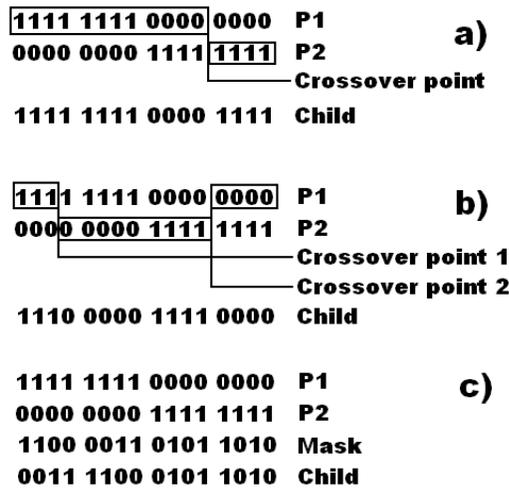


Figure 2.3: *Recombination procedures in a genetic algorithm. a) Single-point crossover. b) 2-point crossover. c) Uniform crossover.*

Finally, a feature is mentioned here that, although not recombination in strict terms, is still convenient to introduce at this point. This feature is called elitism and comes from the desire to make sure that the most fit individual or individuals are carried over to the next generation. In elitism, highly fit individuals are simply put in the next generation without being altered in any way. By this it is ensured that the best “knowledge” of the situation being optimized is not gained only to be lost in later generations.

The last stage of the GA iteration is mutation. Several options exist for introducing mutation. The most common mutation method is bit flipping. Mentioned earlier, it is recalled that this involves flipping a randomly chosen bit to its opposite value - from 1 to 0 or from 0 to 1. Other methods are inversion, for which part of or all of the chromosome is inverted to its opposite direction, or reordering, for which two substrings of arbitrary length are swapped in the chromosome. Individuals carried over by elitism are not subjected to mutation. The probability of mutation for a single bit is typically set between 0.001 and 0.1.

When mutation is applied, the creation of the next generation of individuals is complete. This new generation is in turn subjected to the same steps as the preceding one until the GA is ended.

Measuring GA performance and convergence

An interesting problem with the genetic algorithm is determining when it could reasonably be said to converge, and along the same lines, determining the performance of the algorithm as it progresses. In other minimization procedures, determining convergence is sometimes straightforward, an unambiguous metric for this purpose being readily available. In the GA, a clear-cut measure of convergence does not exist, neither qualitatively nor quantitatively.

A naive first approach to determine the performance of the GA is looking at the average fitness $P_{\text{on-line}}$ of all calls to the fitness function during the entire GA run. This is called the on-line performance, and is given by [37]

$$P_{\text{on-line}} = \frac{1}{C} \sum_{i=1}^C f(i). \quad (65)$$

In this expression, $f(i)$ is the value of function call i and C is the total number of such calls. While this quantity will usually tend to increase during the course of the GA, it may oscillate during the course of the GA run, and it may also be markedly increased by the presence of a few very unfit individuals not representative of the actual improvement of the population. This means that the on-line performance is not really an acceptable estimate.

A better choice for this purpose is the so-called off-line performance. This metric resembles the on-line performance, but is more adjusted towards looking at the best individuals. Assuming that individuals are not arranged in any particular order of fitness, the off-line performance is given by

$$P_{\text{off-line}} = \frac{1}{C} \sum_{i=1}^C f^*(i), \quad (66)$$

where $f^*(i)$ is the best fitness encountered up to and including function evaluation i when performing the summation. The off-line performance is more closely linked to the best individuals while still retaining some information about the rest of the population, and will also decrease monotonically during the GA run, making this a better choice than the on-line performance. A related but simpler way of measuring performance is just taking the fitness of the best individual so far. This will generally give a more discontinuous result than the off-line performance, but is still an acceptable choice.

Measuring convergence is related to measuring performance, but it is generally of interest to include more information about the entire population rather than just the best individuals. To explain the distinction: Performance is related to the goodness of the results obtained thus far, while convergence is related to the likelihood that further improvements will occur. A common measure of convergence is looking at the number of converged bit positions in the chromosome. A bit is said to be converged if it takes the same value in a fraction of individuals greater than some threshold, typically 0.9. Convergence may then be declared when the fraction of bits converged, compared to L_C , is greater than some other threshold, typically 0.8. Another measure of convergence is simply looking at the change in on-line or off-line performance from one generation to the next. Convergence may be declared if this number stays low or equal to zero for a threshold number of generations. This line of thought may also be applied to the simple performance measure using just the fitness of the best individual.

Shortcomings and challenges

As a perfect optimization routine does not exist, every such algorithm carries some weaknesses. The main weaknesses of the simple genetic algorithm are presented in this section. Some of these shortcomings is inherent to the GA framework, while some are particular to some implementation. Stating these weaknesses in this section provides a starting point for the discussion of their potential mitigation and is useful for further reference.

A genetic algorithm may be prone to premature convergence. The presence of fit individuals will mean that the presence of fit substrings of their genes will increase rapidly from these individuals being favored for breeding. While this will be good if the resulting individuals are sufficiently fit (so that the results are good enough for termination of the GA), the risk remains that they aren't, and once the diversity is decreased, it will have a hard time growing again. This risk may be mitigated by employing methods for which the tendency of rapidly decreased diversity is lessened. An example of this is basing roulette selection on rank rather than fitness, as detailed in the previous section. Another obvious way of providing longer lasting diversity is increasing the population size. As will sometimes be the case, the number of available fitness function calls are limited due to time or computational resource constraints. As the number of such calls used during the course of the GA will be equal to $N_i \times N_g$, where N_i is the number of individuals in the population (assuming a population of fixed size) and N_g is the number of generations for which the algorithm is run, such a limitation would mean that increasing the number of individuals would mean decreasing the number of generations for which the algorithm can be run. Should such a limitation be present, a balance must be struck between retaining diversity and degree of coverage of parameter space, and allowing for progression of the population into a fit state. Another method of postponing convergence is increasing the mutation rate. By introducing a larger number of small variations into the genetic bulk, a higher degree of "local diversity" is likely to be observed. Of course, as mutation is intended to produce only such local variations, this will not do much to strengthen the preservation of radically different chromosomes.

Another problem, sampling of absurd parts of parameter space, was mentioned earlier. This kind of sampling must be distinguished from mere sampling of a bad, low-fitness part of parameter space: Absurd sampling means that a combination of parameters is taken that cannot possibly be a configuration for the system under study. As will be discussed in the outline of the experimental method section, parameter ranges can, in the present work, be set up so that this risk is absent, and therefore this will not be a problem for the purposes of this work.

A third feature of the genetic algorithm is that it uses a relatively high amount of evaluations of the fitness function during its course. As this is not strictly a problem concerning the actual functioning of the algorithm, it will not be treated in detail, but it is remarked that this would give an added incentive to produce an optimized procedure for the evaluation of the fitness function.

Next, crossover methods may be implemented in such a way that long substrings providing good fitness for the individual of which they are part are disrupted, thereby hampering the algorithm's progression. To counteract this, parameters believed to be related or coupled in some way may be placed close together on the chromosome, so that the chance of their separation is lessened, or the crossover scheme may be revised, for example by decreasing the number of crossover points used in an n -point crossover method.

The final problem mentioned is that of the GA's tendency to perform poorly in local search. The genetic algorithm is often stated to be well suited for coarse searches of parameter space, and not as a local minimization tool. However, an optimization procedure

principally based on the genetic algorithm can be equipped to tackle local optimization problems. In the following sections some of these will be detailed.

Additional features

In this section, a few variations and extensions on the features presented above will be given. As there is a vast number of potential such extensions, this section will be limited to features considered for implementation in the present work.

The first variation presented concerns the way in which individuals are represented. A variation on the binary representation detailed previously is representing individuals by their true parameter values, that is, not performing the transformation into binary. By using this “real representation”, some of the methods presented with respect to the binary representation will be changed slightly, but their working principles and intended functioning will remain largely the same. This is mainly due to the fact that the smallest unit of the chromosome by real value representation will be the parameter, without further subdivision into bits. Apart from the obvious aspect of not having to transform between the binary representation of chromosomes to real values, the calculation of the fitness function is unaffected. This is also true for selection methods as they are independent of the representation. Recombination methods will essentially remain the same, save for the fact that crossover points will be located between entire parameter values, without the possibility of location “inside parameter borders” as could be done in the binary representation by taking crossover points inside the substring representing a single parameter.

The mutation methods described above, however, are intimately linked to the binary representation and do not readily find a counterpart in the real representation, apart from the notion of mutation rate. If a real representation is to be used, a different procedure must be devised for performing mutation. Such a method should keep with the main idea of creating a local variation in the genetic material as opposed to radically changing individuals subjected to mutation. A method suited for this task is Gaussian mutation. Suppose that some parameter x_i in a chromosome has been selected for mutation. Then, Gaussian mutation will introduce a local variation by assigning

$$x_i \leftarrow x_i + M_i, \tag{67}$$

where M_i is a random variable drawn from a normal distribution $N(0, S_i)$, where S_i is a variance metaparameter for x_i . One disadvantage of this method is the possibility that such a mutation may step outside the predefined parameter bounds. This can be counteracted by discarding such a mutation and performing another one until a valid mutation is given. This will ensure that parameter values stay inside bounds, but it will also to some extent bias the mutation by an “edge unfavorability effect”. This is because if the value of x_i is close to one of the bounds, mutations moving away from this bound will have a higher probability of being accepted. The probability is not strictly equal to 1 as such a mutation may be large enough to bring the parameter beyond the other bound. However, as a sensible choice for S_i will be much smaller than the present bounds, as the intended effect of mutation is local variation, such large mutations will under the normal distribution be negligibly unlikely for practical purposes. In addition to this, the edge unfavorability effect will be less pronounced, thereby rendering Gaussian mutation as a feasible alternative for most applications.

Finally, by using real parameters it is not possible to measure bit convergence as detailed in section 2.2, but in conclusion, using real parameters instead of a binary representation does not yield any significant disadvantages, and the behavior of both implementations will largely be the same. As the implementation time will likely be lower for the former,

it is an attractive alternative for situations where time is limited and the algorithm is implemented from scratch.

The other topic of this section is local optimization. As was mentioned in the previous section, genetic algorithms show poor performance in local searches as they are mainly designed for broad coverage of parameter space. However, local optimization may be introduced to a genetic algorithm. This is called Lamarckian GA after Lamarck's hypothesis of inherited characteristics. Suppose that some local optimization procedure is available. In a typical Lamarckian implementation, individuals will be subjected to local minimization before their fitness is evaluated. The original genes of the chromosome will then be replaced by the optimized parameters, and the corresponding fitness will be the function value at this minimum. If binary representation is used, this will be wrapped in a decoding/re-encoding stage from binary to actual and back. Decoding/re-encoding is obviously not needed if real representation is used. A variation on this scheme is discarding the minimized parameter values, keeping only the function value to be used as the fitness. The rationale for this is, as parameters in a local minimum may be very specialized, combining them with other parameters, as is done in the recombination stage of the GA, could create a very unfit offspring. This effect may not be so dramatic for parameter values emerging from the non-Lamarckian workings of the GA, and so, keeping the unminimized parameter values could prove beneficial for subsequent generations.

The main disadvantage of a Lamarckian GA is the fact that it may use a very large amount of function calls. It is clear that local optimization may create highly fit individuals more rapidly, but as the population will usually contain a large proportion of unfit individuals, the function calls spent on optimizing these "lost causes" will essentially be wasted. This may be counteracted by only introducing local optimization at a late point in the GA, that is, after several generations, or alternatively, doing it only intermittently, for instance at every tenth generation. Even then, the number of function calls used in local optimization may quickly grow to rival the number used for all the non-Lamarckian generations. Often, these calls are better spent performing more of the coarse, standard GA searching. In any case, parameters taken from a completed non-Lamarckian GA run should be subjected to local minimization as this will often provide a significant improvement on what that GA can provide. This may be all the local minimization that is actually needed.

2.3 Finding an experimental method

Now that the model and the genetic algorithm has been presented, it is time to put the preceding two sections together to create an experimental procedure by which the values of the parameters used in the model can be determined. As was mentioned in the introduction, the key idea for determining these parameters is using some kind of optimization scheme wherein the function optimized is one that is closely linked to the accuracy of the properties calculated by the model, taking a proposed set of parameters as its arguments. An optimal point of such a function will be a point for which the corresponding parameters hopefully could be said to enable the model to produce accurate results. As will be discussed later, just the task of finding such a function is nontrivial as the notion of accuracy is ambiguous.

The theory presented in this section will be limited to cover the obtainment of parameters for prediction of the molecular dipole moment and the static molecular polarizability tensor, that is, a molecular polarizability tensor for which the electric field invoking the polarizability is static. The article[2] presenting the model has carried out such a parametrization for a limited selection of systems using a local optimization method. In that work, two different sets of parameters were obtained; one set for the prediction of the dipole moment and another set for the polarizability. The argument for not seeking

to find a unified set of parameters for the prediction of both properties is that the mechanisms or effects from which each of these properties arise may be said to be quite different. For the polarizability, being a response property for an external electric field, band-gap considerations are important, while for the dipole moment, the formation of chemical bonds and local effects play a large role.

Certain aspects of the procedure presented here, principally concerning the use of a genetic algorithm, could conceivably be employed for the obtainment of parameters for the calculation of other properties and may thus have applicability beyond those just mentioned, but the exploration of this potential applicability is beyond the scope of the present work and will therefore not be covered.

Obtaining reference values

An obvious prerequisite for any scheme seeking to determine a good set of parameters is having some means of comparing the values obtained by a prospective set of parameters to a set of reference values in which there must be a satisfactory level of confidence. In other words, the first task to be completed is obtaining a collection of such values for use as “answer keys”. It was mentioned in the introduction that for this purpose, *ab initio* quantum mechanical methods is considered the best choice available. In earlier work by the authors of the model [40]-[43], Hartree-Fock (HF) calculations were used to provide these reference values. For the properties considered in this work, the results given by density functional theory (DFT) are usually regarded as superior to those obtained from HF theory, but a drawback to DFT has been problematic behavior for large systems, and improvements have been suggested [44]-[52]. This would be an argument against using DFT for the reference values in the present work, as one point of interest is studying how the properties change with increasing molecular size, for which increasing chain length is the most tangible measure. However, the introduction of current-DFT [53] provides improved description of the molecular polarizability of large systems as compared to other flavors of DFT, rendering this an attractive choice for our purposes.

Assuming that a molecule or collection of molecules for which reference values is to be found has been decided upon, the first step of calculation is the optimization of the molecular geometry. The usual scenario is that molecules are “built” using a graphical interface, in which the geometry is primitively optimized, typically using some low-level force field theory. An optimization of these approximate geometries, consistent with the level of theory to be employed in the calculation of the properties of interest, is clearly called for. The appropriate properties can then be calculated from the optimized geometries. Finally, these geometries along with the calculated properties must be put into a format available for use by the procedure employed for the optimization of the parameters in the model.

Some considerations concerning the choice of exchange/correlation functional and basis set must also be made. In the model article, the BLYP functional and a TZP basis set has been used for the geometry optimization and calculation of molecular dipole moments. For the calculation of polarizabilities, a current-DFT implementation and an augmented TZP basis set has been used. The current-DFT scheme available for use in this work, namely the one implemented in the DFT program ADF [54]-[56], uses a linear density functional, which is regarded as a simplified model, meaning that while being a good choice for the calculation of the molecular polarizability due to its improved description of large systems, this scheme may not be such a favorable choice for geometry optimization and the calculation of the molecular dipole moment when more sophisticated functionals are available for this purpose.

Setting up the optimization

In this and succeeding sections, it will be assumed that accurate reference values for the properties under study are readily obtainable. When that aspect is covered, it is time to look at how one may develop a procedure for the optimization of the model parameters.

It will come as no surprise that an optimization method believed to be a suitable choice for the present work is the genetic algorithm. The original model article employed a local optimization method. As such methods are dependent on the choice of an initial guess for all the parameters subjected to optimization, the employment of such a method would require a quite high degree of initial knowledge about the parameter values. If such knowledge is not available, it would likely be necessary to perform a large number of local optimizations so that a broad coverage of the parameters for which the knowledge is not good can be attained. As there might be several such parameters, this kind of search can be tedious to the level of unfeasibility. Of course, such a series of local optimizations may also be performed automatically through the use of a script, with a low degree of or even no human interaction, but the number of local searches required may still prove to be prohibitively large.

A global optimization method has the potential to do away with this problem. Given that such a method can be made to cover a large part of parameter space to a sufficiently large level of detail, there will be no need to systematically perform local minimization, save perhaps for a few select promising individuals found by a global, broad search method. As the genetic algorithm is a method designed for this kind of search, its employment could show to be an improvement over local methods employed in earlier parametrization efforts.

Using a genetic algorithm, the only strict requirement of initial parameter knowledge is the specification of the range in which each parameter is allowed to vary. This range may be tailored to each parameter to reflect the degree of initial knowledge about that parameter. By this, a parameter for whose value the knowledge is near certain may be given a very narrow range. By using such a narrow range, assuming that the knowledge employed in this narrow range specification is accurate, individuals in the GA population will tend to show better fitness compared to a situation where a wider range had been given. In a sense, it could be said that individuals by this will not be "wasted" on a broad search of a parameter range for which the feasible value range is actually narrow. On the other hand, parameters for whose values the knowledge is poor may be given a large parameter range to provide a broad coverage. For these parameters, the GA computational effort associated with searching such a relatively broad range could be said to be "time well spent". All in all, a sensible specification of parameter ranges will ensure that the GA is made to search with a suitable level of broadness for each parameter.

Before moving on to the specification of the fitness function for use in the GA, it is necessary to make some observations. There are three important considerations to be taken here - the first is the potential limitations of the model for which the parameters are to be determined, and the second one is choosing an appropriate size and content for the set of reference molecules. Finally, if parameters for several elements are sought, it is necessary to determine if and how the optimization can be partitioned in a sensible and computationally economical way.

The first consideration should really be regarded as obvious, but it may be easy to forget when attempting to perform a successful parametrization: If the model is fundamentally weak in some respect, no amount of parameter optimization will overcome this weakness. The reason for mentioning this is that such a weakness may lead to significant trouble in deciding whether a parametrization could be deemed to be successful, meaning to

have provided the model with a set of parameters so that it will function to the best of its ability. The potential trouble is due to the fact that it may be difficult to decide whether some significant prediction error from a set of parameters should be ascribed to shortcomings in the optimization procedure or an inherent weakness in the model. There is no clear-cut way of deciding this, but it will often be that model weaknesses will be more inclined towards producing larger errors for some particular type of system rather than showing a more uniform distribution of errors. This may however be confounded by effects arising from the choice of reference sets. This is the next consideration.

The selection of a good reference set may be a difficult procedure. As the model under study employs atom-type parameters, this discussion will be from this perspective. There are two important factors which must be taken into consideration when selecting a reference set. The first factor is ensuring that all roles that an element may take in a molecule are found somewhere in the reference set. If limitations to the model are known, the roles falling under these limitations should naturally be excluded. An example of this is the fact that carbon-carbon triple bonds are not covered by the model featured in this work. Now, the ambition to cover all roles is of course impossible to fulfill completely, as this would mean using every possible molecule. However, by using a fairly large reference set, one may achieve a certain level of confidence in the width of applicability of the parameters thus obtained. Again, there is no certain way of knowing if a reference set is large enough. For this purpose it is necessary to use chemical intuition or experience. The second factor concerns the balancing of the reference set. In the introduction, an example involving alcohols and aldehydes was mentioned, wherein it was suggested that using an imbalanced reference set could lead to correspondingly imbalanced accuracy, wherein the precision for one type of system is sacrificed for the benefit of another type. The reliance on reference set balancing may to some extent be relaxed by the choice of fitness function, as will be discussed later.

The final consideration to be presented before discussing the fitness function is the partitioning of the optimization scheme. When using atom-type parameters, it may be beneficial to partition the optimization to cover a few elements at a time, fixing the parameter values already optimized when extending the optimization to cover other elements. The way in which this is done also a nontrivial task, as possible weak or strong dependencies between parameters may mean that parameters suited for a particular collection of elements may be completely unsuitable for the incorporation of other elements. This may be because during an optimization, some effects may be wrongly described by parameters for one element, when in fact these effects should have been ascribed to other elements or some intermediate balance. Then, when another element is introduced, there may be no way to achieve satisfactory results as the parameters for the other elements are fixed. This may to some extent be mitigated by including molecules consisting of only a subset of all the elements, but this may not always be convenient or even possible.

Another point to be made on this topic concerns the number of elements to take into consideration for a particular optimization run. As was seen in section 2.2, the population size of a genetic algorithm should increase with the number of parameters, but if more elements are to be included in a single run, the reference set must also increase, thereby resulting in roughly squared scaling with the number of elements for which parameters are to be determined in a single run. Limiting the number of parameters in each run will reduce the scaling with roughly an order of magnitude when compared to taking more elements (and thus determining more parameters) in a single run, as the single-run scaling will be decreased by a squared factor, but the number of runs needing to be performed will add roughly another order of magnitude, decreasing the total run time.

As organic systems generally are of high interest, the first run will usually have to

include both carbon and hydrogen. The risk of using only these two elements is the one mentioned above - that the parameters of one of them will take the “burden” for the description of the entire effect. This risk will likely tend to decrease if more elements are included in this run, but as this will increase the computation time it is generally not desired to use too many if the run-time tends to approach prohibitive levels. A compromise solution adopted in this work is taking carbon, hydrogen and one other element in the first run. By this it is hoped that a balance will be struck between the problems described above.

The considerations above have provided a background for the discussion of the choice of fitness functions, and this will be covered presently. It is obvious that the fitness function will have to be some measure of the discrepancy between the values of the properties as predicted by the model, and the QM reference values, but there are several choices each of which may be said to carry some validity. In the following, the phrases “fitness function” and “error function” should be understood to be equivalent. The term “error function” is not to be confused with its namesake from the integration of a Gaussian function.

The first property under consideration is the molecular dipole moment. This quantity will be a three-dimensional vector describing the dipole moment along each axis. Let $\mu_{j, \text{ref}}^i$ be the reference value for Cartesian axis j for this quantity for molecule i in the test set, and let $\mu_{j, \text{calc}}^i$ be the corresponding value as predicted by the model for some collection of parameters. The Euclidian dipole moment error $E_{\mu, \text{euc}}$ will then be given by

$$E_{\mu, \text{euc}} = \sum_{i=1}^M \left(\sum_{j=1}^3 (\mu_{j, \text{calc}}^i - \mu_{j, \text{ref}}^i)^2 \right)^{\frac{1}{2}}, \quad (68)$$

where M is the number of molecules in the reference set. This error measure corresponds to the total length of the error vector, and will therefore be dominated by the largest component of this vector. Another measure, $E_{\mu, \text{txc}}$, is in analogy to the so-called “taxicab” metric, and is given by

$$E_{\mu, \text{txc}} = \sum_{i=1}^M \sum_{j=1}^3 |\mu_{j, \text{calc}}^i - \mu_{j, \text{ref}}^i|. \quad (69)$$

This “taxicab” error measure will measure the sum of the discrepancy of each component, and could be said to pay more equal attention to each component of the error vector than will its Euclidean counterpart. Both of these error measures can be used in an implementation and will usually only give slightly different results.

The polarizability is a 3×3 tensor. For molecule i , let $\alpha_{jk, \text{ref}}^i$ and $\alpha_{jk, \text{calc}}^i$ be the reference and model-calculated jk -component of the polarizability tensor, respectively. In analogy with the error measures for the dipolemoment, the Euclidian and taxicab measures $E_{\alpha, \text{euc}}$ and $E_{\alpha, \text{txc}}$ for the polarizability is then given as

$$E_{\alpha, \text{euc}} = \sum_{i=1}^M \left(\sum_{j=1}^3 \sum_{k=1}^3 (\alpha_{jk, \text{calc}}^i - \alpha_{jk, \text{ref}}^i)^2 \right)^{\frac{1}{2}} \quad (70)$$

and

$$E_{\alpha, \text{txc}} = \sum_{i=1}^M \sum_{j=1}^3 \sum_{k=1}^3 |\alpha_{jk, \text{calc}}^i - \alpha_{jk, \text{ref}}^i|, \quad (71)$$

respectively. An intermediate summation in which rows or columns of the polarizability tensor are summed and then combined in some way does not carry any physical meaning and should therefore be left out of consideration as an error measure.

Now, the error measures detailed above can all be used as a fitness function for the minimization of the total error. However, all of them will be sensitive to the balancing of the reference set. If one of these error measures are employed for an optimization run where the reference set is markedly imbalanced toward one or more types of molecular systems, the optimization may tend to minimize the errors for the systems frequently represented at the cost of increasing the errors for those that are not, as this will provide an overall decrease in the total error. While this may be satisfactory for some applications, one should at least consider finding an error function for which the sensitivity to reference set balance is not so high. One such function is the ‘‘maximal absolute error function’’ $E_{max}^{abs}[E(\cdot)_i]$, given as

$$E_{max}^{abs} = \max[E(\cdot)_i], \quad (72)$$

where i is the molecular system of the test set giving the largest error as calculated from one of eqns. (68)-(71), denoted by $E(\cdot)_i$. Using this error measure, the algorithm will seek to minimize the largest error in the test set, thereby eliminating the sensitivity to test set balancing. However, as the reference values of the properties under study will often be subject to large variations from one reference molecule to the next, as is a very marked tendency for polarizability values and also to some extent for dipole moments, this error measure will usually prove to be a poor choice. That is because a system for which the reference property is large will usually also produce a large error vector compared to a system for which the reference property takes a small value. This means that even if such an error is relatively large, it may actually be a quite satisfactory result. Furthermore, the algorithm may be able to bring this absolute error for large reference value systems down to a level comparable to the errors encountered for systems for which the reference value is small. This will mean that while the accuracy achieved for large reference value systems may be excellent, the ‘‘small-value’’ accuracy may be abysmal. This is obviously an unacceptable consequence. The resolution of this problem comes from switching from maximum to relative error. The relative-error equivalent of eqns. (68)-(71) are

$$E_{\mu, \text{euc}}^{\text{rel}} = \sum_{i=1}^M \frac{\left(\sum_{j=1}^3 (\mu_{j, \text{calc}}^i - \mu_{j, \text{ref}}^i)^2 \right)^{\frac{1}{2}}}{\left(\sum_{j=1}^3 (\mu_{j, \text{ref}}^i)^2 \right)^{\frac{1}{2}}}, \quad (73)$$

$$E_{\mu, \text{txc}}^{\text{rel}} = \sum_{i=1}^M \frac{\sum_{j=1}^3 |\mu_{j, \text{calc}}^i - \mu_{j, \text{ref}}^i|}{\sum_{j=1}^3 |\mu_{j, \text{ref}}^i|}, \quad (74)$$

$$E_{\alpha, \text{euc}}^{\text{rel}} = \sum_{i=1}^M \frac{\left(\sum_{j=1}^3 \sum_{k=1}^3 (\alpha_{jk, \text{calc}}^i - \alpha_{jk, \text{ref}}^i)^2 \right)^{\frac{1}{2}}}{\left(\sum_{j=1}^3 \sum_{k=1}^3 (\alpha_{jk, \text{ref}}^i)^2 \right)^{\frac{1}{2}}}, \quad (75)$$

and

$$E_{\alpha, \text{euc}}^{\text{rel}} = \sum_{i=1}^M \frac{\sum_{j=1}^3 \sum_{k=1}^3 |\alpha_{jk, \text{calc}}^i - \alpha_{jk, \text{ref}}^i|}{\sum_{j=1}^3 \sum_{k=1}^3 |\alpha_{jk, \text{ref}}^i|}. \quad (76)$$

These relative error sums will not carry the same dependence on the size of the reference values as their absolute-error counterparts and are also viable candidates for use as fitness functions. The maximum relative error function may now be introduced as

$$E_{\max}^{\text{rel}} = \max[E(\cdot)_i^{\text{rel}}], \quad (77)$$

where $E(\cdot)_i^{\text{rel}}$ is one of the error measures (73)-(76).

Another possible choice for an error function may be found by combining a summative and a maximal error function. As these error functions may take widely different values, such a combination should use the product of these functions instead of the sum, so that each function may be given equal weight. Should it be desired to give one type of error function thus combined higher priority, it may be multiplied to a higher exponent than the other.

It may be desirable to use different error functions in the GA search and a local minimization disconnected from the GA itself. As the GA is not suitable for local searches, it may be that some candidate giving low error rates overall may show bad results for a few systems. As some local work might give significant improvements to these worst fits, there is a higher incentive for introducing some variant of a maximum error function in local minimization rather than in the GA.

Convergence and ending the search

A survey of convergence and performance tracking in genetic algorithms was given in section 2.2. For determining when to end the genetic algorithm, one of the methods presented there may be used. However, it may be the case that a satisfactory level of performance/convergence may be judged by simple inspection, or it may be included in the implementation so the GA ends when a sufficiently low error function value has been achieved. When doing such a judgement, it should be taken into consideration that there may be significant improvement from local minimization of the best individual of the GA. Hence, although a good candidate from a non-Lamarckian GA may not fall below the "satisfactory error threshold", it may ultimately prove to be a very good result after some local minimization work.

Summary of the experimental procedure

Here, we give a stepwise summary of the experimental procedure to bring together all the material presented thus far. First, molecules are selected for the reference sets. Their geometries are optimized and the molecular dipole moments and molecular polarizabilities are calculated. These results are put in a form recognizable by an optimization routine.

The genetic algorithm is set up, choosing methods for selection, recombination and mutation, and introducing elitism if that is desired. The choice is made between real or binary representation. A suitable fitness function is chosen. For one of the quantities, say, the molecular dipole moment, the elements to be covered in the first parametrization run are selected. Molecules spanning a wide range of roles these elements can take are put in a reference set, and the genetic algorithm is run with the comparison between model predicted values and reference values being with respect to the molecules in the reference set. When a satisfactory individual has been found, the GA is halted and this individual is subjected to local optimization by some procedure. At convergence, the parameters are extracted. Next, another element is added and the parameters for the elements already subjected to optimization are held invariant while the GA/local minimization

procedure is repeated, so that the only parameters being optimized are the ones for the new element. More elements are added in this way, possibly the entire collection of parameters is given some more local minimization intermittently, and the GA/local minimization scheme is yet again repeated until no more elements under consideration remain. The entire procedure is then repeated for the molecular polarizability.

3 Implementation and experimental details

In this section, we proceed to describe the particulars of the parameter obtainment carried out in the present work. The first part of this section describes how reference values were gathered. Following this, we give an overview of how a genetic algorithm was implemented and set to work.

3.1 The obtainment of reference values

The calculation of reference values was done in accordance with the considerations put forth in section 2.3. The geometry optimizations and the calculations of molecular dipole moments at optimized geometries were carried out on the 2007 and 2008 versions of the DFT program ADF [56] using a TZP basis set and the BLYP functional. The calculation of molecular polarizabilities were done at optimized geometries by the 2007 version of ADF using current-DFT with an augmented TZP basis set.

The decision on which molecules to include in the reference sets were based on intuition because, as mentioned, there is no “analytic” or clear-cut method of making such decisions. The elements under consideration were limited to carbon, hydrogen, oxygen, fluorine and chlorine. In the following a summary of the systems comprising the reference sets is given. The naming of each molecule is not provided here as they number in the hundreds.

It was decided that a wide selection of hydrocarbon “skeletons” should be represented, both for systems encompassing only these two elements and systems with heteroatom substituents and functional groups. Selected for calculation were single-chain alkanes from methane up to 30 carbon atoms and conjugated dienes from ethene up to 26 and 33 carbon atoms, both with and without substituents. For molecules with substituents, single alkene chains with single double bonds were also used. These ranged from chain lengths of 3 to 20 carbon atoms. Aromatic compounds ranging from benzene to tetraphenyl were selected. In addition to this, reference values were also obtained from a series of asymmetric hydrocarbon molecules. A few of these were also given substituents. The method by which these asymmetric molecules were selected did not follow any regular procedure apart from the desire to cover a wide range and combination of molecule types.

For oxygen, a large selection of functional groups on different carbon skeletons were covered. For fluorine and chlorine, also, a wide selection of substitution patterns on carbon skeletons mentioned above were chosen. Most of the molecules considered were monosubstituted hydrocarbons, but several di-, tri- and larger extents of substitution were covered.

In addition to this, results were obtained for a broad range of compounds containing both functional groups of oxygen and halogen substituents. Care was taken to include molecules of every combination of substituent elements, i.e. containing only fluorine substituents and oxygen functional groups on a selection of hydrocarbon skeletons, only chlorine and oxygen, only fluorine and chlorine, and some containing all three elements.

In summary, it is believed that a broad coverage of functional groups of oxygen and their interplay with halogen substituents has been achieved. This is also believed to hold for the halogen substituents and pure hydrocarbon systems.

3.2 The genetic algorithm and local minimization. Experimental procedure

A program for the calculation of molecular dipole moments and polarizabilities by the model presented in section 2.1 has been implemented using PYTHON [38] and its extension SCIPY [39]. This environment has also been used for the genetic algorithm. For

local minimization, a Simplex method from the optimization library of SCIPY has been used. This method does not make use of derivatives, being convenient in this case as no analytical derivatives are available. Furthermore, the local optimization method used was constrained to taking only positive parameter values.

A non-Lamarckian version of the genetic algorithm was chosen, as the calculation of the aforementioned properties by the model was decided to be too demanding for a Lamarckian implementation as the size of the reference sets could be large (regularly approaching 100 systems). The genetic algorithm was chosen to use real value representation of parameters for ease of implementation. The breeding pool size was set to a fraction of 1/2 of the population size. For the selection procedures, 1/5 of the breeding pool was filled by cutoff selection, while the remaining 4/5 was filled by rank-based roulette selection. For recombination, 1/2 of the population size was created by single-point crossover, while the remaining 1/2 was created by 2-point crossover. The mutation rate was set to 0.03, and Gaussian mutation was used. Finally, elitism was introduced by replacing an arbitrary individual of the new generation with the best individual of the previous generation.

For the η_I^* parameters, the parameter range was set to [0.1, 20] and mutation standard deviation to 1.6. For the Φ_I^* parameters, the parameter range was set to [0.1, 12] and mutation standard deviation to 0.8. For the α_I^* parameters, the parameter range was set to [0.000001, 3.0], taking the inverse of the polarizability as a parameter, and the mutation standard deviation was set to 0.16. The permitting of this parameter to thus reach a very high values could be called an experimental mistake, but as will be revealed in the coming section, no α_I^* parameters ended near this limit and it was therefore without serious consequences. For the χ_I^* parameters, the parameter range was set to [0.1, 3.0] and mutation standard deviation to 0.4, but for the elements other than carbon and hydrogen, the electronegativity of carbon was added as well to force a higher electronegativity for these other elements. For the R_I^* parameters, the parameter ranges were set close to the notion of “half of an ‘equilibrium’ bond length”. The range for this parameter for carbon was set to [1.25, 1.45], for hydrogen the range was [0.55, 0.85]. For oxygen, the range was set to [1.0, 2.5], but this range may have been too wide. For fluorine, the range was set to [1.2, 1.5]. For chlorine, the range for the polarizability parametrization was mistakenly set to [1.2, 1.5]. For the dipole moment parametrization, this range was set to [1.0, 2.5], which may have been too broad. The mutation standard deviations were set to 0.4. This is likely a too high mutation rate when compared to the ranges for this parameter. For the C_I^* parameters, the range was set to [0.1, 6] for carbon and hydrogen, in order to roughly approach the value of the global C in the original model article[2] where these systems were studied [2] was set to 5.0, and [0.1, 40.0] for other elements, as there was little knowledge about their possible values. Mutation standard deviations for this parameter were set to 0.8. The initial GA population was created by assigning random values inside these parameter ranges for the specification of individuals. The GA was ended by simple inspection, looking at the similarity of fitness values in the population and the fitness of the best individual.

The parametrization proceeded by the procedure summarized at the end of section 2.3. Selections of the molecules for which reference values had been calculated were put to use as test sets for the genetic algorithm. The composition of these sets were made in accordance with the considerations discussed earlier in section 2.3, selecting a number of systems deemed adequate for good “element role” representation. For the molecular polarizability, the first parametrization run included molecules consisting of carbon, hydrogen and oxygen. The genetic algorithm, with a population size of 360, was run followed by local optimization of the best individual. Following this, the elements fluorine and chlorine were added in turn for the obtainment of their respective atom-type parameters, running the genetic algorithm with a population size of 210 for each run and

performing local optimization on the best individual. For these runs, parameters already obtained were held constant. In the genetic algorithm runs, the taxicab absolute error function given by eqn. (71) was used as a fitness function, as this was judged to be a good choice as it would to a high degree take individual elements into consideration. For the local optimization, this error function was used both by itself and also multiplied with a maximum relative error function (77), taking the error function given by eqn. (76) as its argument $E(\cdot)_i^{rel}$. The reason for including the relative error function in this way was that it was desired to reduce the worst error while still having some way of reducing the total error if it was to prove possible to satisfy both of these simultaneously. These two error functions were used in turn until a satisfactory error level had been reached. The final minimization did not include the maximum error, however, as it was regarded as more important to reduce the total error than the maximum error. Finally, a final local optimization was performed wherein the parameters for all the elements were allowed to vary. For this final run, the error function (71) was used exclusively.

For the dipole moment, it was elected to use molecules of the elements carbon, hydrogen and chlorine for the first parametrization, using for this and the following runs only molecules where the dipole moment was nonzero. The test sets were composed and varied in a way similar to that of the polarizability parametrization. The parametrization proceeded by the same method as for the polarizability, employing first the genetic algorithm with a population size of 360 followed by local minimization. Holding the parameters then obtained constant, parameters for fluorine were obtained in the same way with a GA population size of 210. In the GA runs, the error function used was eqn. (69). For the local minimization, the error function from eqn. (68) was used, at times multiplied with the maximum relative error function taking the function eqn. (73) as its argument. For the final run, eqn. (68) was used exclusively. The switch from the taxicab error function to its Euclidian counterpart was an experimental mistake, but it is believed that this did not have a large effect on the results.

The addition of oxygen could not be completed for the dipole moment parametrization. The reason for this is two-fold. Firstly, earlier test runs not using the modification of the global C parameter into individual, atom-type parameters, as detailed in section 2.1, showed that no satisfactory results could be obtained for a wide range of oxygen substitution patterns. For instance, using reference sets comprised of ethers, carboxylic acids, alcohols, esters, aldehydes and ketones on a broad selection of carbon skeletons would give satisfactory results only for some of these systems. No combination could be found wherein all types of systems in the set gave satisfactory results. Secondly, when the modification of the C parameter was introduced, time constraints regrettably demanded that this search be abandoned. Therefore, dipole moment parameters were obtained only for the elements carbon, hydrogen, fluorine and chlorine.

4 Results and discussion

4.1 Polarizability

In this section, we present the results of the parametrization of the molecular polarizability. We will start by presenting the parameters obtained, and move on to present results for a selection of systems. This presentation is done graphically for molecular chains and in tables for a selection of other molecules.

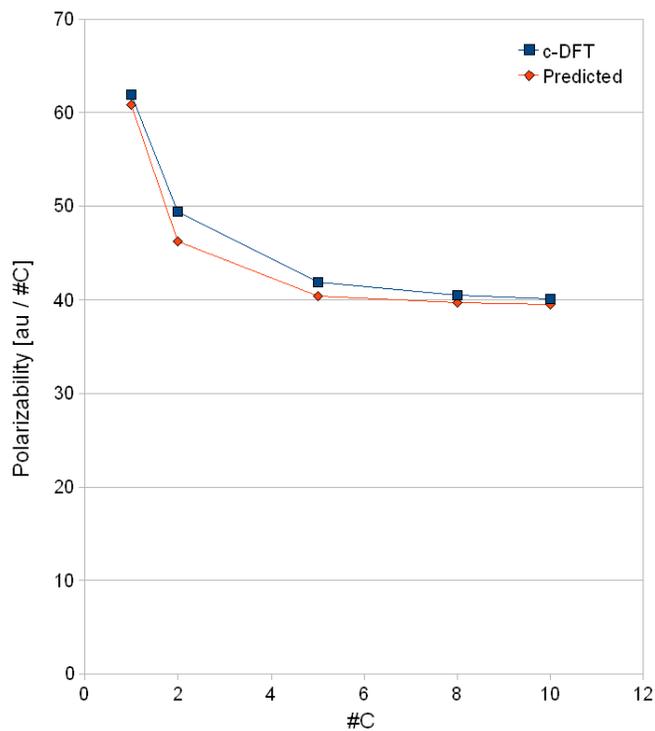


Figure 4.1: *Reference and predicted polarizability per carbon atom of nonsubstituted alkane chains from the trace of the polarizability tensor.*

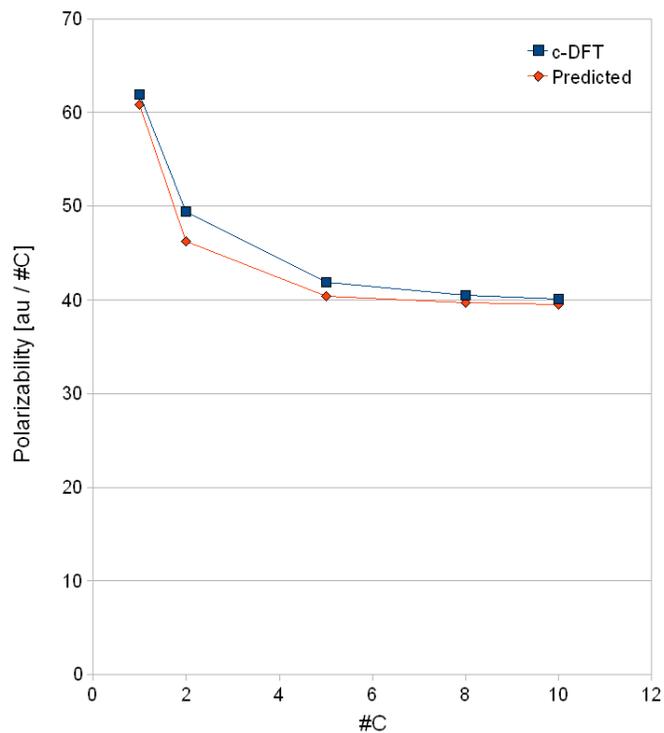


Figure 4.2: Reference and predicted polarizability per carbon atom of alkane chain alcohols (substituent at end of chain) from the trace of the polarizability tensor.

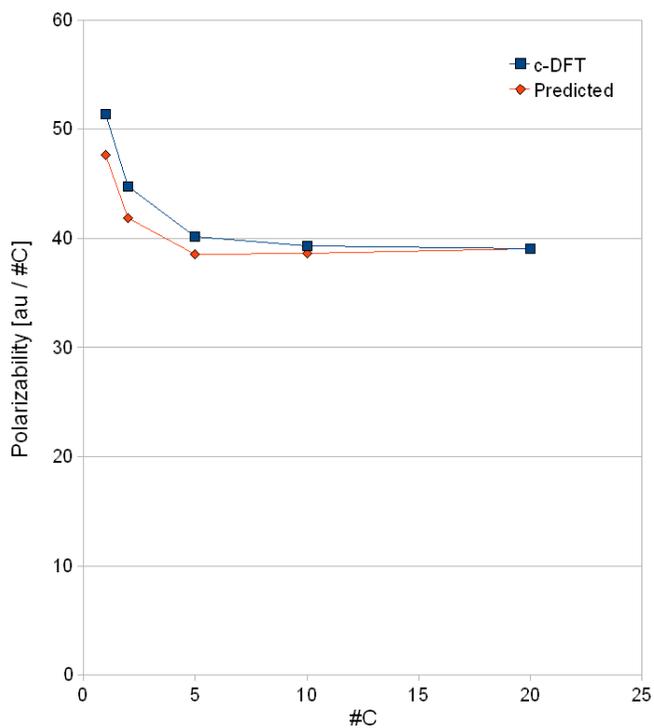


Figure 4.3: Reference and predicted polarizability per carbon atom of alkane chain aldehydes from the trace of the polarizability tensor.

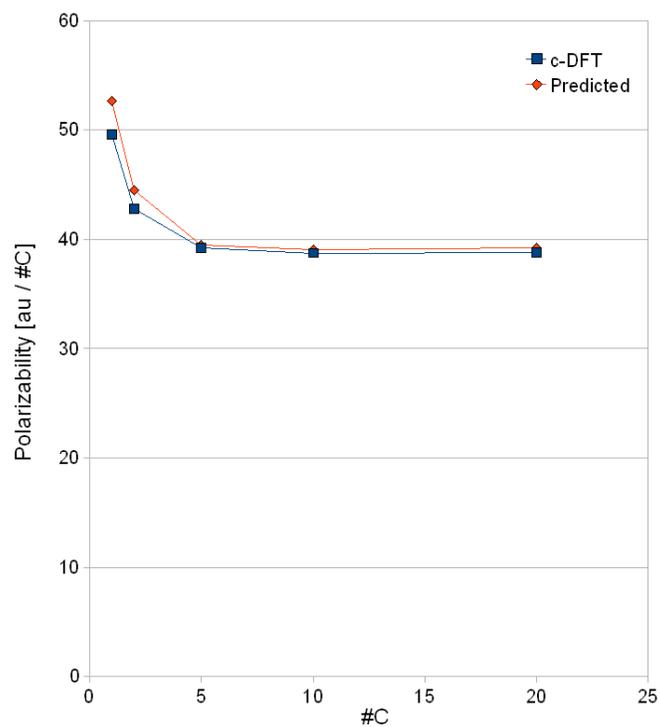


Figure 4.4: Reference and predicted polarizability per carbon atom of monofluorinated alkane chains (substituent at end of chain) from the trace of the polarizability tensor.

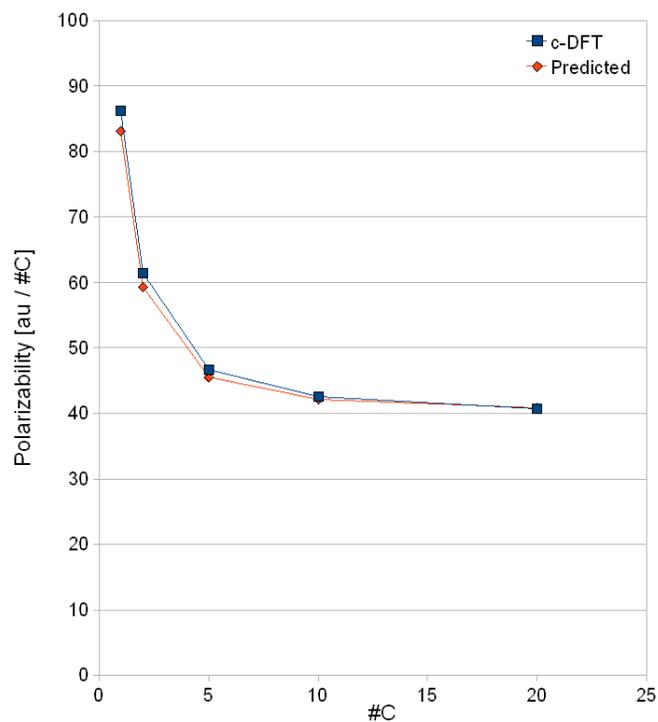


Figure 4.5: Reference and predicted polarizability per carbon atom of monochlorinated alkane chains (substituent at end of chain) from the trace of the polarizability tensor.

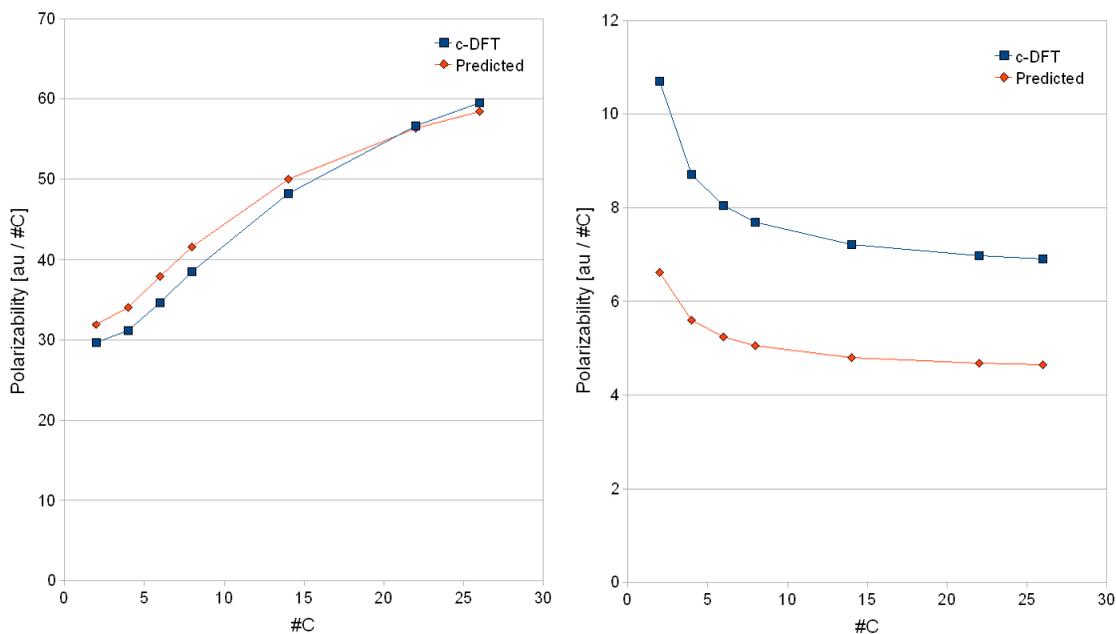


Figure 4.6: Reference and predicted polarizability per carbon atom of alkene chains with alternating double bonds. In-plane polarizability (left) from sum of two diagonal tensor elements. Out-of-plane polarizability (right) from the final diagonal tensor element.

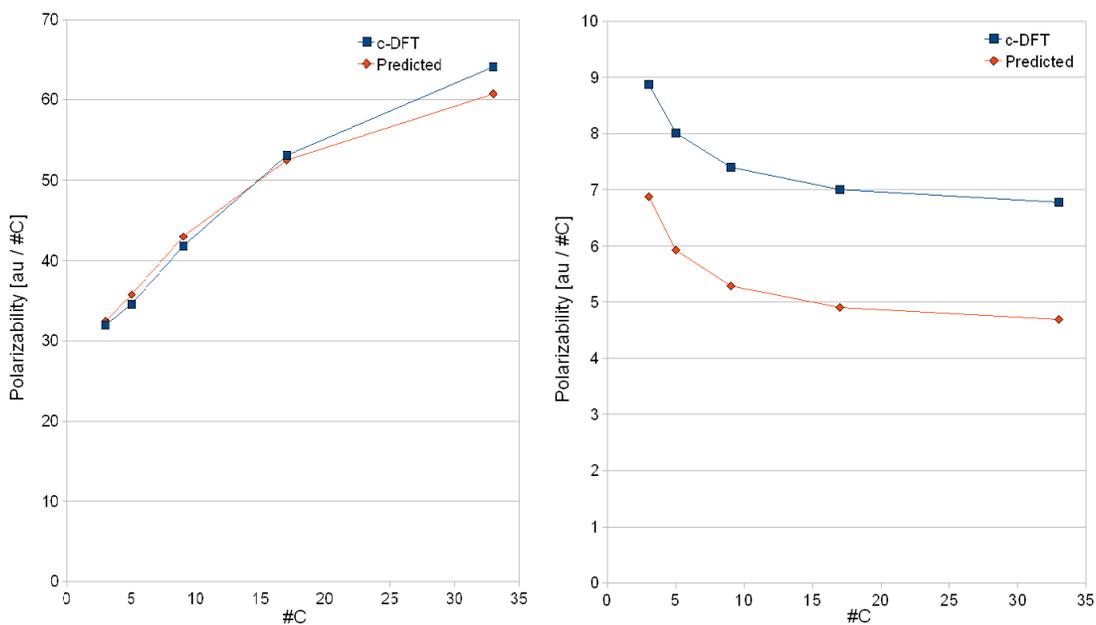


Figure 4.7: Reference and predicted polarizability per carbon atom of alternating double bond alkene chain aldehydes in trans position. In-plane polarizability (left) from sum of two diagonal tensor elements. Out-of-plane polarizability (right) from the final diagonal tensor element.

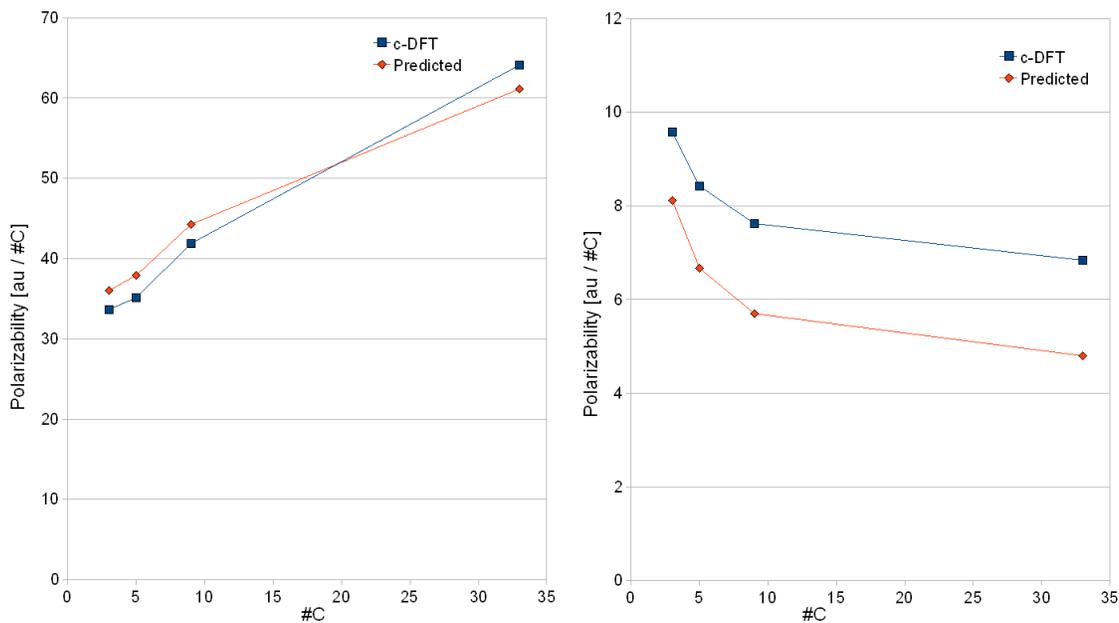


Figure 4.8: Reference and predicted polarizability per carbon atom of alternating double bond alkene chain carboxylic acids in cis position. In-plane polarizability (left) from sum of two diagonal tensor elements. Out-of-plane polarizability (right) from the final diagonal tensor element.

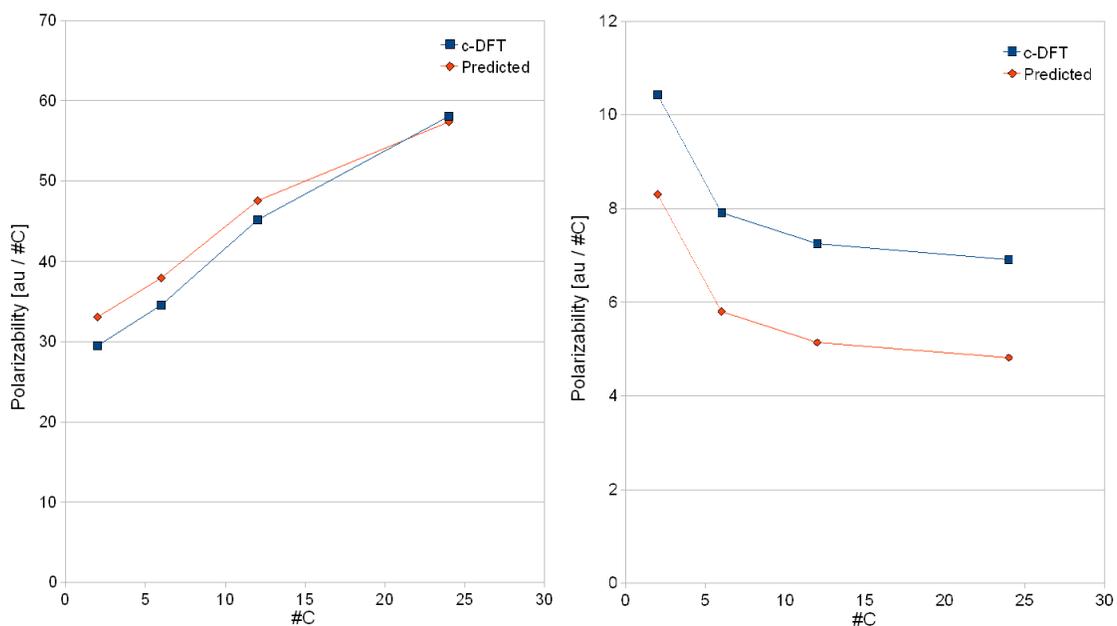


Figure 4.9: Reference and predicted polarizability per carbon atom of monofluorinated alternating double bond alkene chains (substituent at end of chain in trans position). In-plane polarizability (left) from sum of two diagonal tensor elements. Out-of-plane polarizability (right) from the final diagonal tensor element.

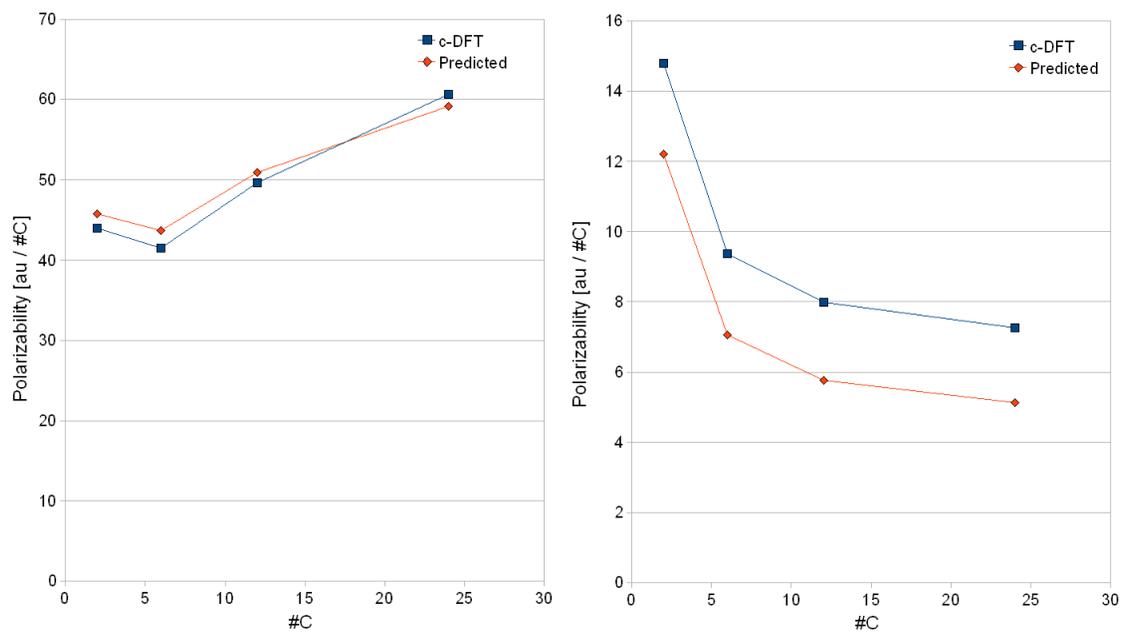


Figure 4.10: Reference and predicted polarizability per carbon atom of mono-chlorinated alternating double bond alkene chains (substituent at end of chain in trans position). In-plane polarizability (left) from sum of two diagonal tensor elements. Out-of-plane polarizability (right) from the final diagonal tensor element.

Table 4.1: Reference and predicted polarizability values for planar molecules. Reference and predicted in-plane polarizabilities $\alpha_{\text{in}}^{\text{ref}}$ and $\alpha_{\text{in}}^{\text{pred}}$ from the sum of two in-plane diagonal terms of polarizability tensor. Out-of-plane reference and predicted polarizabilities $\alpha_{\text{out}}^{\text{ref}}$ and $\alpha_{\text{out}}^{\text{pred}}$ from the final diagonal term. All values in atomic units.

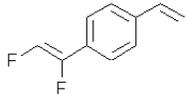
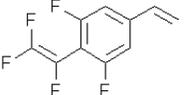
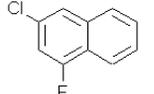
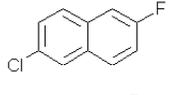
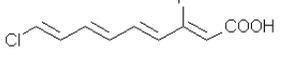
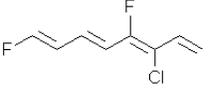
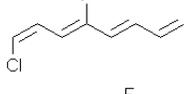
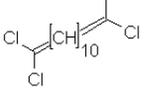
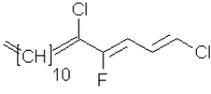
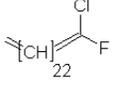
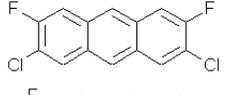
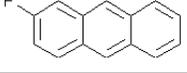
Molecule	$\alpha_{\text{in}}^{\text{ref}}$	$\alpha_{\text{in}}^{\text{pred}}$	Rel. err.	$\alpha_{\text{out}}^{\text{ref}}$	$\alpha_{\text{out}}^{\text{pred}}$	Rel. err.
	307.7	327.3	6.39%	66.6	52.3	-21.48%
	313.9	338.8	7.94%	66.3	62.9	-5.15%
	310.1	318.3	2.64%	69.4	56.1	-19.09%
	310.3	314.3	1.28%	69.4	56.1	-19.17%
	422.1	444.6	5.33%	76.3	65.7	-13.85%
	329.9	365.1	10.67%	67.4	57.5	-14.62%
	87.7	93.3	6.33%	29.2	27.9	-4.42%
	340.3	376.3	10.57%	68.8	54.7	-20.56%
	671.3	675.4	0.61%	110.3	93.5	-15.28%
	884.2	908.0	2.69%	128.4	100.6	-21.64%
	1439.6	1409.5	-2.09%	173.6	126.2	-27.33%
	499.2	492.7	-1.30%	97.3	84.2	-13.43%
	455.0	459.0	0.89%	90.2	66.7	-26.07%

Table 4.2: Reference and predicted polarizability values for disubstituted benzene. Reference and predicted in-plane polarizabilities $\alpha_{\text{in}}^{\text{ref}}$ and $\alpha_{\text{in}}^{\text{pred}}$ from the sum of two in-plane diagonal terms of polarizability tensor. Out-of-plane reference and predicted polarizabilities $\alpha_{\text{out}}^{\text{ref}}$ and $\alpha_{\text{out}}^{\text{pred}}$ from the final diagonal term. All values in atomic units.

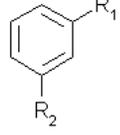
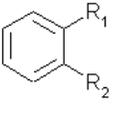
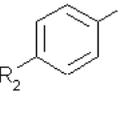
Molecule	R ₁	R ₂	$\alpha_{\text{in}}^{\text{ref}}$	$\alpha_{\text{in}}^{\text{pred}}$	Rel. err.	$\alpha_{\text{out}}^{\text{ref}}$	$\alpha_{\text{out}}^{\text{pred}}$	Rel. err.
	F	OH	166.9	172.6	3.45%	43.4	35.0	-19.21%
	F	CHO	196.2	199.1	1.46%	46.6	38.4	-17.68%
	F	COOH	201.1	209.5	4.20%	48.7	42.0	-13.89%
	F	Cl	187.1	191.8	2.52%	49.2	42.0	-14.63%
	Cl	OH	200.2	201.9	0.86%	51.7	42.3	-18.25%
	Cl	CHO	230.4	229.5	-0.40%	54.9	45.6	-17.05%
	Cl	COOH	234.4	239.9	2.35%	56.9	49.1	-13.75%
	F	OH	166.5	172.8	3.76%	43.2	35.1	-18.90%
	F	CHO	195.0	201.6	3.40%	46.6	38.4	-17.51%
	F	COOH	200.5	211.8	5.64%	48.6	42.1	-13.47%
	F	Cl	187.2	192.2	2.68%	49.2	42.1	-14.50%
	Cl	OH	198.1	201.8	1.88%	51.2	42.0	-17.87%
	Cl	CHO	226.8	228.3	0.66%	54.5	45.3	-16.78%
	Cl	COOH	231.5	238.7	3.11%	56.2	48.7	-13.40%
	F	OH	166.3	172.3	3.62%	43.1	35.0	-18.79%
	F	CHO	197.9	199.1	0.61%	46.5	38.4	-17.48%
	F	COOH	202.2	209.2	3.47%	48.6	42.0	-13.75%
	F	Cl	186.4	191.5	2.70%	49.0	42.0	-14.23%
	Cl	OH	199.8	202.1	1.15%	51.6	42.3	-18.03%
	Cl	CHO	234.8	230.6	-1.80%	54.9	45.6	-16.91%
	Cl	COOH	238.5	241.0	1.07%	57.1	49.2	-13.85%

Table 4.3: Reference and predicted polarizability values for non-planar molecules. Reference and predicted polarizabilities α^{ref} and α^{pred} from the trace of the polarizability tensor. All values in atomic units.

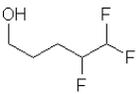
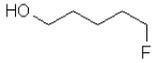
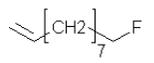
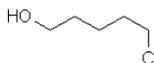
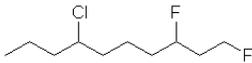
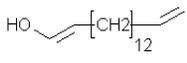
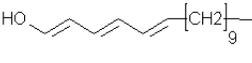
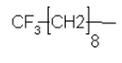
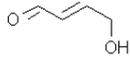
Molecule	α^{ref}	α^{pred}	Rel. err.
	206.3	238.5	15.61%
	62.5	63.9	2.18%
	207.7	208.3	0.30%
	383.7	392.3	2.23%
	244.5	238.4	-2.52%
	418.2	435.3	4.08%
	628.5	636.7	1.31%
	680.7	693.3	1.85%
	395.4	387.0	-2.13%
	99.0	97.7	-1.36%
	176.5	176.7	0.12%

Table 4.4: Parameters obtained from the parametrization of the molecular polarizability. All values in atomic units.

Element	η_I^*	Φ_I^*	α_I^*	R_I^*	C_I^*
H	5.8823	0.3685	1.8349	0.6666	0.5704
C	19.561	0.7429	5.8591	1.4145	1.8825
O	2.4700	0.1130	4.7813	1.3159	21.099
F	12.333	1.57×10^{-5}	4.3312	1.2548	11.659
Cl	37.475	0.1128	14.376	0.4023	58.552

From Table 4.4 above, it is seen that some parameter values are strongly different from the rest. For hydrogen and carbon, the parameters seem to be within a largely plausible range. For oxygen, fluorine and chlorine, there is a tendency for the R_I^* to be lower than what should be expected from its interpretation as “half of the bond length” between the element in question and another atom. The values of C_I^* for these elements are correspondingly larger. The risk of this happening was mentioned in section 2.1 where these new parameters were introduced. Another consequence of this is that the parameter Φ_i , which is coupled to R_I^* by eqn. (60), may deviate in value according to

the way in which R_J^* was altered by the coupling to C_J^* . Ideally, the values would all fall inside a plausible range and not vary much from those of the best individual found by the genetic algorithm, but as some known or unknown shortcoming in the model could lead to one of these parameters deviating in value, the others will be “forced along” into deviation by the coupling as the local minimization procedure tries to satisfy some possibly small effect which the model in reality may not be able to encompass by varying one of the parameters. Some room for the deviation of parameters should therefore be granted. If stricter control of parameter ranges is desired, a local minimization procedure constrained to a higher degree than just positive parameter range could be adopted.

Turning to the results, an average relative error of 7.90% over all reference set molecules was found using the error function (71) and dividing by the sum of all the absolute values of all polarizability tensor elements in the entire set. This number, although being a somewhat condensed measure of accuracy, is representative of the average performance of the model for a large selection of the systems for whose description the model was designed. This low error rate shows that the model is largely capable of giving accurate results approaching QM values for the systems it was designed to describe, and so must be said to be a good result. When presenting results for individual molecules or molecular chains, as in the following, the diagonal polarizability tensor elements are the only ones taken into consideration. The reason for this is that the trace of the polarizability tensor is rotationally invariant, and as the molecules are arbitrarily rotated, the diagonal elements are used for reasons of reproducibility. Furthermore, for planar systems, the out-of plane polarizability is independent of the in-plane elements, and therefore, if a molecule is rotated so that the out-of-plane vector is parallel to one Cartesian axis, the diagonal element of the polarizability tensor corresponding to this is invariant with respect to the rotation of the molecule around this axis. In this way, the invariant part of the in-plane polarizability can be identified as the sum of the two diagonal elements of the polarizability tensor corresponding to the molecular plane and the out-of-plane polarizability will correspond to the final diagonal element.

Please note that that average relative error given above includes anisotropic polarizability elements, for which the model generally may not be expected to give a good description due to the use of an isotropic atomic polarizability. It should therefore be kept in mind that the error rates encountered in the following are generally better than would be the case if anisotropic tensor elements were also taken into consideration.

Turning now to the description of individual groups of systems, the first systems under consideration are alkanes with or without end-of chain monosubstituents, shown in figures (4.1-4.5). It is seen from these graphs that the agreement is generally very good. The largest errors tend to occur for methane, ethane and their substituted equivalents. The trends with increasing chain length is reproduced correctly.

The results for alkene chains with alternating double bonds (henceforth “conjugated alkene chains”) with and without end-of chain monosubstituents are shown in figures (4.6-4.10). For these chains, the polarizability is split into an in-plane contribution summing two diagonal elements and an out-of-plane contribution from the last tensor element. The other tensor elements are zero for all these molecules. The in-plane accuracy is good while the out-of-plane accuracy is somewhat low. In the latter, the chain length variational trends from the reference sets are reproduced correctly, but there is a systematically lower value for all systems.

Looking at both alkane and conjugated alkene chains, it is seen that the polarizability per carbon atom seems to converge. This is not observed conclusively for the in-plane polarizability of the conjugated alkene chains as the data for longer chains was not obtained, but there can still be seen a tendency for this to occur for those systems as the values for successive long chain points increase at a lower rate than for successive

short chain points.

For the substituted planar systems in table 4.1 and 4.2, the results echo those of the conjugated alkene chains. For the disubstituted benzene systems, the in-plane polarizability is generally in good agreement with reference values. For the other systems in this category, the in-plane polarizability is not well reproduced for the ethene-phenyl-ethene and short conjugated alkene systems, but shows good agreement for the rest. The out-of-plane polarizability is consistently predicted as lower than the reference values. With two exceptions, this lowering is quite pronounced. This behavior is also observed for the conjugated alkene chains. The out-of-plane discrepancy is believed to be an effect of the employment of isotropic atomic polarizabilities. Turning again to the chain length variations for the out-of-plane polarizability, as the trend itself is reproduced correctly, there is considerable reason to believe that had some anisotropic atomic polarizability scheme been employed wherein the atomic polarizability was separated into in-plane and out-of-plane terms, the out-of-plane predicted values could have been scaled with the out-of-plane term to not only reproduce the chain variation trend, but also approach the reference values.

For the substituted nonplanar systems in table 4.3, the results are generally very good. A notable exception is found for 1,2,2-trifluoropentane-5-ol. It is not known why this molecule does not show good results, and bringing about an understanding of this could be an interesting task. There may be that some effect is present that the model is not capable of handling, but the good results for the other substituted systems makes the inaccuracy all the more puzzling. The similar systems fluoromethanol and 5-fluoropentanol both show low errors, and from this it may be hypothesized that the “culprit” is the high level of fluorine substitution. However, as the results for another somewhat similar system, 1,1,1-trifluorodecane, are good, it may be that the effects are brought about by some interaction between the hydroxy group and the fluorines. This, however, is unlikely, as the fluorine atoms and the hydroxy group are too far apart for any considerable interaction to occur. Another explanation may simply be that this error is a fluke brought about by some unfortunate consequence of the parameter values for all elements or for fluorine alone. In this case, a repetition of the entire parametrization could provide a better result.

As the molecules presented discussed above all must be said to have a significant degree of substitution, there is support for concluding that the model is largely capable of handling substitution in carbon chains in the calculation of the molecular polarizability. In other words, introducing other elements is not an obstacle in itself for the model, but particular types of configurations may still introduce inaccuracies. Looking at the planar systems with the largest in-plane discrepancies (assuming that the out-of-plane discrepancies largely stem from the lack of anisotropy discussed above), it is seen that a common feature of these systems is the presence of only a few ethene groups.

The error seems to decline when the conjugated alkene chain is lengthened beyond 8 carbon atoms. As short and intermediate conjugated alkene chains to some extent can be regarded as ethene oligomers, it may be argued that increasing the number of ethene units gradually makes the chain lose its “etheneic” character and assume a character of its own. The way in which this behavior arises is not well understood and is a prime candidate for further study. If this behavior can be described and taken into consideration in the model, it is likely that more accurate results can be obtained. The following section contains a discussion that may be applicable to this behavior. In conclusion, the molecular polarizability is well described by the model, and while some discrepancies taint the results they must all in all be called satisfactory, verging on very good.

4.2 Dipole moment

Table 4.5: Bond lengths for end-of-chain fluorinated conjugated alkenes. Bond lengths given for C-F bond, C-C bond of the fluorinated carbon and next carbon in chain, and the C-H bond of the fluorinated carbon. Units in Ångström. Angle θ_{ref} and θ_{pred} for reference/predicted molecular dipole moment vector and C-F bond in degrees. Negative angle means that dipole moment points away from the chain compared to bond vector.

Number of C atoms	R_{CF}	R_{CC}	R_{CH}	θ_{ref}	θ_{pred}
2	1.3769	1.3277	1.0881	-1.52°	-10.69°
4	1.3758	1.3398	1.0875	-3.46°	-7.37°
6	1.3758	1.3427	1.0872	-3.74°	-4.61°
8	1.3759	1.3441	1.0871	-3.87°	-2.73°
12	1.3759	1.3452	1.0870	-2.91°	-0.35°
24	1.3758	1.3460	1.0869	0.57°	1.98°

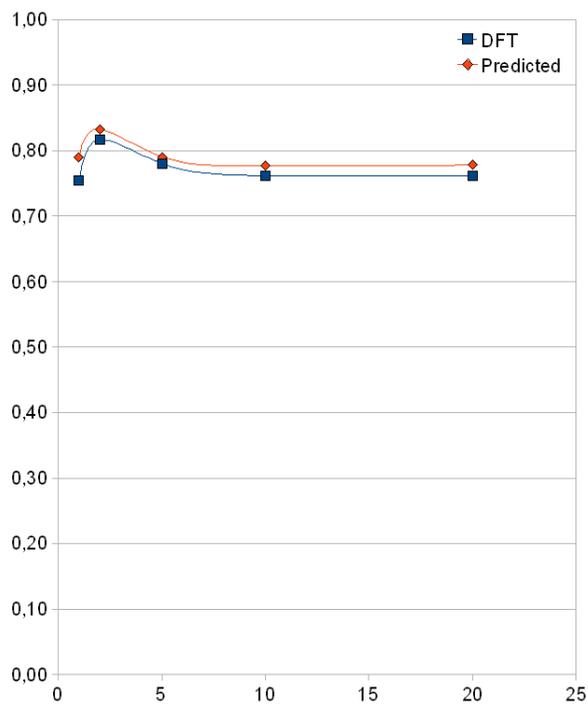


Figure 4.11: Reference and predicted total dipole moment of monofluorinated alkane chains (substituent at end of chain).

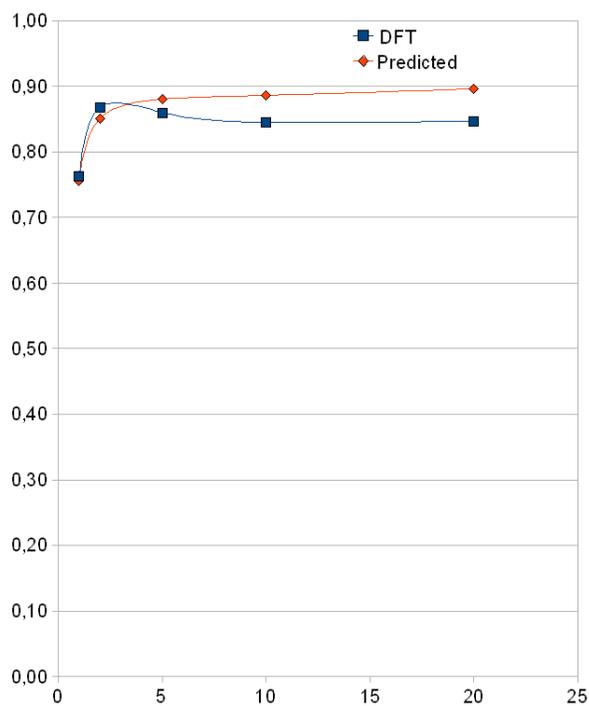


Figure 4.12: *Reference and predicted total dipole moment of monochlorinated alkane chains (substituent at end of chain).*

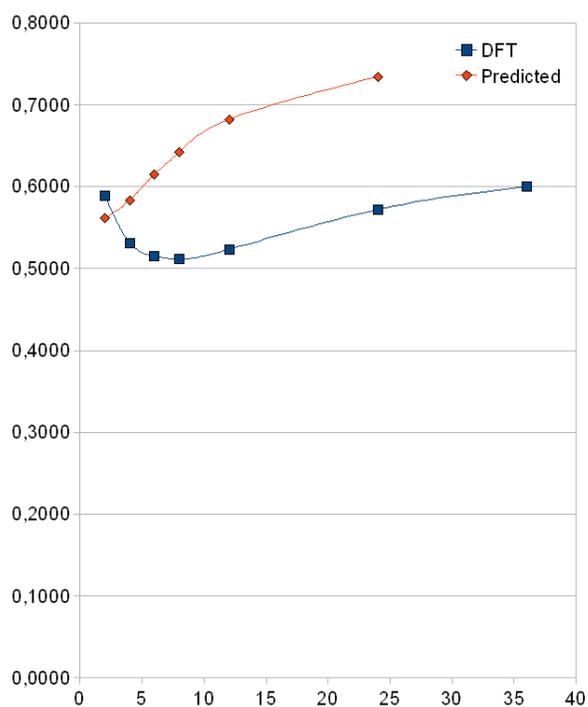


Figure 4.13: *Reference and predicted total dipole moment of monofluorinated conjugated alkene chains (substituent at end of chain in trans position).*

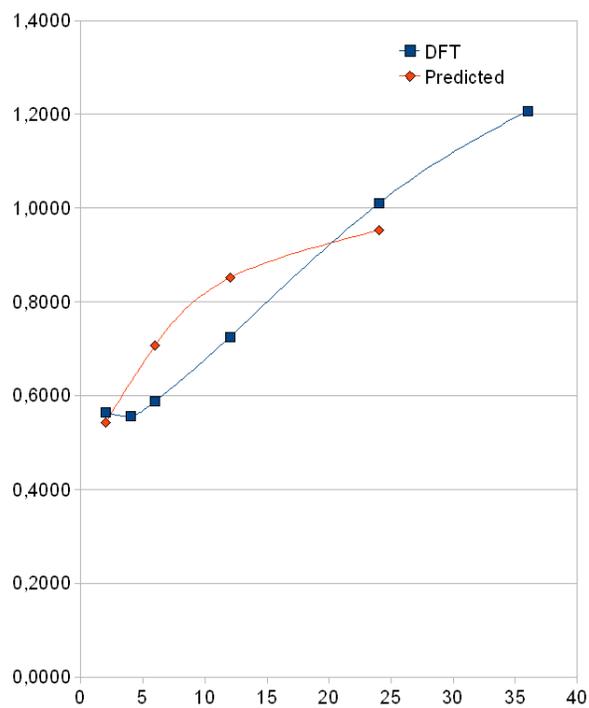


Figure 4.14: *Reference and predicted total dipole moment of monochlorinated conjugated alkene chains (substituent at end of chain in trans position).*

Table 4.6: Reference and predicted total molecular dipole moment values μ^{ref} and μ^{pred} for a selection of molecules. All values in atomic units.

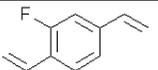
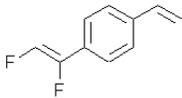
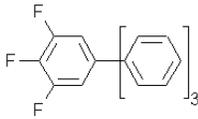
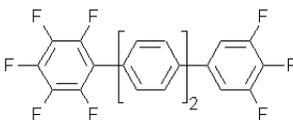
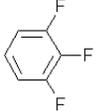
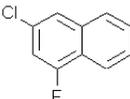
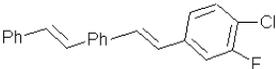
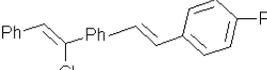
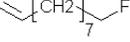
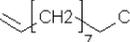
Molecule	μ^{ref}	μ^{pred}	Rel. err.
	0.583	0.494	-15.26%
	1.111	1.112	0.10%
	1.778	1.622	-8.78%
	0.327	0.507	55.14%
	1.216	1.032	-15.15%
	0.659	0.594	-9.84%
	1.036	1.009	-2.53%
	0.034	0.015	-57.68%
	0.829	0.813	-2.01%
	1.522	1.472	-3.26%
	0.600	0.887	47.87%
	0.765	0.761	-0.49%
	0.877	0.883	0.65%
	0.733	0.672	-8.35%
	0.827	0.804	-2.78%

Table 4.6 contd.

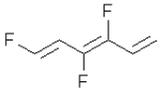
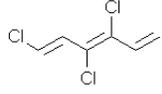
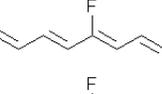
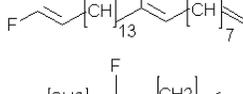
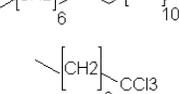
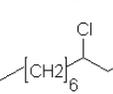
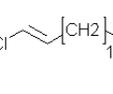
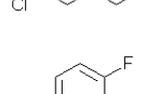
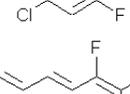
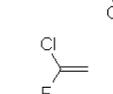
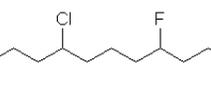
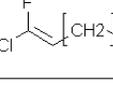
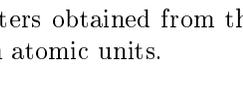
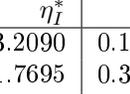
Molecule	μ^{ref}	μ^{pred}	Rel. err.
	0.509	0.565	10.98%
	0.516	0.663	28.58%
	0.444	0.531	19.66%
	0.270	0.526	94.75%
	0.701	0.700	-0.20%
	1.020	1.022	0.21%
	0.838	0.855	1.95%
	0.938	0.962	2.58%
	0.166	0.184	10.37%
	0.607	0.506	-16.69%
	0.257	0.413	60.53%
	0.550	0.478	-13.08%
	1.457	1.434	-1.63%
	0.872	0.871	-0.13%

Table 4.7: Parameters obtained from the parametrization of the molecular dipole moment. All values in atomic units.

Element	η_I^*	Φ_I^*	α_I^*	χ_I^*	R_I^*	C_I^*
H	3.2090	0.1253	3.7793	0.0	0.6043	2.5034
C	1.7695	0.3210	8.0321	0.9205	1.4318	5.6616
F	8.2666	70.572	3.5088	4.5610	1.6232	0.0171
Cl	4.5428	96.541	6.9638	1.6404	2.1490	0.2423

From table 4.7, it is seen that the parameters Φ_I^* and C_I^* for fluorine and chlorine are much higher and lower, respectively, than the corresponding parameters for carbon and hydrogen. Like the parameters for polarizability, it is believed that this can be explained by the “parameter wandering” in the local minimization procedure exacerbated by some problems to be described in the following. The rest of the parameters are largely within or near a plausible range. The parameter χ_H^* is set to zero as the model deal only with electronegativity differences, leaving a degree of freedom here used in this way.

Turning to the results, we remark that the total molecular dipole moment is the one presented as it is invariant of rotation and therefore reproducible for arbitrarily rotated molecules. The error rate for all molecules in the test set was found to be 12.19% from the Euclidian error function (68) divided by the sum of the total dipole moments of the entire test set. This error rate, like the one given in the polarizability section, is also a condensed measure of error, but it is deemed to be a little too high to be called satisfactory.

Turning now to the results for molecular chains, it is seen from figures 4.11 and 4.12 that the total dipole moment for end-of-chain monohalogenated alkane chains is reproduced well. The results are particularly good for fluorine. The accuracy for the chlorinated alkane chains drops off a bit as the chain grows larger, but is still within a satisfactory range. The results for the alkanic individual molecules in table 4.6 are comparably good, with some systems showing excellent agreement to the reference values. The results are also good for systems with one double bond and for some aromates.

For the conjugated alkene chains, however, the results are not so good. From figures 4.13 and 4.14, it is seen that the model fails to reproduce the trend for increasing chain lengths, and for the individual molecules in table 4.6, it is seen that some results for these systems are highly inaccurate. Now, in the model article [2], it is stated that the total molecular dipole moment for these end-of-chain substituted conjugated alkene chains, showing near-metallic behavior, is expected to increase monotonically with the chain length but finally converge for long chains. The increase is due to the fact that charge transfer to some extent occurs between the substituent and the carbon to which it is bonded. As the chain length increases, charge transfer occurs easily along the length of the chain as its behavior is nearly metallic, and so the chain is able to “draw” harder on the substituent atom the longer it is. This gives an increase in the molecular dipole moment. As can be seen from these figures, the model manages to provide this behavior, but the reference data do not follow this trend. The principal difference between the predicted trend and the reference trend is that the total dipole moment for the reference data decreases when going from ethene to larger chain lengths, reaching a minimum for substituted conjugated octatetraene for fluorine substitution and for conjugated butadiene for chlorine substitution. After this minimum, the reference values increase again before tending to convergence for larger chain lengths. The reference values for chlorine increase more rapidly than those for fluorine. Thus, the model for this parametrization run falls short of providing a good agreement for these systems, although it may seem that the “increase until convergence” effect mentioned above is present to some extent in the reference data, but accompanied by some effect not yet accounted for.

At this point it may be convenient to make an attempt to explain the reference data behavior, as it is an unexpected result. As the behavior in question is qualitatively the same for both fluorine and chlorine, it is sufficient to consider one type of system, and here, fluorine is chosen.

As the molecular dipole moment is dictated largely by local effects in the molecule, the area in the vicinity of the substituent should be the one investigated. The explanation, if any can be found, should come from effects near the substituent. Table 4.5 provides the

pertinent information. In this table, it is seen that the dipole moment vector and the C-F bond vector are nearly parallel for the reference values. This supports the conclusion that the dipole moment is primarily dictated by interaction over the C-F bond, as was supposed when suggesting to look at local effects. It is also seen that for all except the longest chain, the dipole moment vector points away from the chain. The reference angles to some extent follow the trend in figure 4.13, in the sense that broadness of these dipole moment/bond angles is ranked in the same way as the total dipole moment, with the exception of the last two points in the table.

Table 4.5 also reveals that the C-F bond length does not vary significantly over the chain length, with the exception being fluoroethene (the first point). Still, even this difference is not really very large. The same holds for the C-H bond length. However, the bond length between the carbon bonded to the substituent and the next carbon in the chain does vary to a larger degree than is observed for the other bond lengths under consideration. This would suggest that there is some mechanism dependent on this bond length that is at play in bringing about the initial decline in the reference curve, recalling that the subsequent increase for longer chains is another effect that is certainly within the grasp of the model. This C-C bond length, however, shows a monotonic increase with bond length, which means that if it is in some way responsible for bringing about the initial decline, its effect must be after a while be overwhelmed by the chain length effect. This is potentially a plausible explanation.

To phrase it more clearly: The initial decline may be brought about by the fact that R_{CC} increases with the chain length. This should mean that the character of this bond should assume a less metallic character. By this, charge transfer is hampered in its path down the carbon chain, thereby denying the total dipole moment of assuming an angle oriented more towards the chain. This “less metallic effect” grows with the length of the chain, but after a while, the chain has grown so large that the charge transfer that actually happens between the first two carbon atoms gives a large effect, at some point overpowering the tendency to decline brought about by the decrease in R_{CC} . Again, the fact that the angles θ_{ref} do not strictly adhere to this may either mean that this explanation is incorrect, or possibly that some secondary effect is at play, or again that they don’t need to be strictly related to the variation in trends. Another counterargument is that this effect could perhaps be said to be too weak to bring on the observed trend. This must stand unanswered at present. A final remark is that this effect could possibly explain some of the errors encountered for the polarizability, as the errors showed some tendency to appear for systems containing ethene in various places.

The question then remains: Supposing that the above explanation is sound, why can this not be accounted for in the model? A suggestion for this proposed here is that $g_{I,KM}$ or S_{IJ} may simply not be sensitive enough to incorporate this variation. The parameters R_I^* and C_I may by circumstance be “tied up” in the description of other and larger effects, or the functional form of $g_{I,KM}$ may be insufficient. This concludes the presentation of the results.

4.3 Performance of the parametrization method

The genetic algorithm can, for the obtainment of parameters for the molecular polarizability, be said to have performed well. As is apparent from the results given by the polarizability parameters, the broad search of the genetic algorithm combined with the subsequent local minimization of the best candidate is able to give good results for this property. However, the parameters for the calculation of the molecular dipole moment do not show results as good as those for the polarizability, and as mentioned earlier, this is believed to be due to effects for which the model is not capable of providing a description, thereby leading the local optimization to unsuccessfully approach a descrip-

tion of these effects by sending parameters into unlikely regions, taking other parameters along due to coupling. As this “parameter wandering” was observed for the polarizability parameters as well, it is clear that unconstrained local optimization may indeed be too unconstrained for this use.

Using a constrained local optimization method, as has not been done in this work, may keep the parameters from “wandering” due to their coupling. The obtainment of satisfactory results under such constraints, however, may be more dependent on the model’s ability to cover all important effects contributing to the quantity under study, as the introduction of constraints would in broad terms mean that parameters are not allowed to go outside a preset range in order to approach some effect for which the model gives an insufficient description. In further development of the model, there may come a point where the number of effects one can incorporate is exhausted, either due to having reached the limit of understanding or due to practical and computational concerns. If there still remains some source of error at that point, one may be content to allow some parameter deviation to run its course.

A viewpoint in criticism of the parametrization method itself is the argument that the obtainment of good parameters for the dipole moment could be more difficult than for the polarizability, and although the optimization scheme was able to complete the “easy task” of the polarizability, it fell short of completing the harder task to a satisfactory level, the implication being that the model could in fact be able to consistently reproduce accurate dipole moment values had it only been given good parameters from the optimization scheme. However, apart from the “parameter wandering” discussed earlier, it is believed that the parameter optimization scheme itself is sound and well suited for further use. The considerations put forth in section 4.1 could serve to vindicate the model’s role in the dipole moment errors, so that its success in the obtainment of polarizability values could be argued to be representative of its actual effectivity, so that the polarizability results are the ones by which the parametrization method is rightfully judged. The conclusion is therefore that the genetic algorithm has functioned well as a framework for parametrization when combined with a local optimization procedure. The latter may benefit from the introduction of constraints.

Improvements to the method may, for instance, include introducing other selection, recombination or mutation methods into the genetic algorithm or using constrained local optimization. There may also be benefits from carrying out a systematic exploration of the choice of fitness function, or simply tuning the parameters governing the working of the genetic algorithm itself such as mutation rate or population size. The setup of the genetic algorithm used in this work has been developed by a nonsystematic procedure, and these subjects could be pursued for improved performance.

4.4 Suggested directions for further work

From the preceding results and discussion there emerge some directions from further work, to be surveyed in this section.

An obvious extension of this work is simply the parametrization of more elements. Limited by time constraints in the present work, it is believed that parameters for several other elements are readily obtainable, providing that the model is capable of describing the effects produced by these elements. The confidence in the model’s ability of description is higher for the molecular polarizability than for the dipole moment, as should be apparent by the preceding results and discussion. A starting point for the further parametrization of elements could for example be nitrogen, sulphur and phosphorus, providing that dipole moment parameters for oxygen can be obtained. By including these elements and obtaining good results, the model can be said to be applicable to a large proportion of organic compounds.

From the results for the reference total dipole moment trends shown for end-of chain monohalogenated conjugated alkene systems and the discussion about them, it is apparent that a prime candidate for further study is bringing about an understanding of this phenomenon. If this effect could be understood and be made to appear in the model by some method, there will potentially be a large increase in accuracy to be gained.

Concerning the polarizability, the principal move for the improvement of accuracy would be the introduction of anisotropic atomic polarizability. From the results presented in this work, in particular for the planar systems, it can confidently be stated that if this is incorporated into the model, the result will readily improve, potentially approaching an excellent level of agreement with reference values.

An extension of the model could for instance be one made to cover frequency-dependent molecular polarizability. Another extension could be the inclusion of other properties such as the dipole moment gradient.

A possible application of the model in its present form would be using the atomic charges q_I , which in this work are used only for the calculation of the molecular dipole moment, for the prediction of reactive sites on substituted molecules. For example, if the molecule in question is a substituted benzene, the substitute will perturb the atomic charges of the carbon atoms in the benzene ring so that some carbon atoms become more positively charged than they would be in the unsubstituted case, and correspondingly, some carbon atoms become more negatively charged. A nucleophilic or electrophilic species being a candidate for substitution onto the benzene ring may show preference for substitution on carbon sites that are more positively or negatively charged, respectively, than the other carbon atoms. This is textbook material in organic chemistry. An accurate representation of atomic charges could then provide an easy and quick method for the prediction of such reactive sites. This interesting topic could regrettably not be pursued due to time constraints, but it is considered an exciting candidate for further study.

For the parametrization method itself, most points have been covered in the preceding section. To summarize, further work could involve using a constrained local minimization method or exploring further the various features of the genetic algorithm.

5 Conclusion

The parameter obtainment of atom-type parameters for the extended EEM/PDI combination model has been carried out using a genetic algorithm. The genetic algorithm is believed to function well for this purpose. The local minimization procedure may benefit from the introduction of constraints to prevent parameter wandering.

The combined EEM/PDI model is able to provide good agreement with reference values in the calculation of molecular polarizability. The model may in this regard benefit from the introduction of anisotropic atomic polarizabilities. The model provides a good description the dipole moment of halogenated alkane systems and for some such aromatic systems, but this agreement is not present for other aromatic systems and conjugated alkene chains.

References

- [1] Leach, A. R., "Molecular Modelling", 2nd ed., Pearson Education Limited (Harlow, Essex, UK), 2001.
- [2] Smalø, Hans S., Åstrand, P.-O., Jensen, Lasse: Nonmetallic electronegativity equalization nad point dipole interaction model including exchange interactions for molecular dipole moments and polarizabilities. *J. Chem. Phys.*, *accepted for publication* (2009).
- [3] Sanderson, R. T., *Science*, **114**, 670-672 (1951).
- [4] Sanderson, R. T., "Chemical bonds and bond energy", 2nd ed., Academic Press, New York, 1976.
- [5] Mortier, W. J., van Genechten, K., Gasteiger, J., *J. Am. Chem. Soc.*, **107**, 829-835 (1985).
- [6] Rappe, A. K., Goddard III, W. A., *J. Phys. Chem*, **95**, 3358-3363 (1991).
- [7] Cioslowski, J., Stefanov, B. B., *J. Chem. Phys.*, **99**, 5151-5162 (1993).
- [8] York, D. M., Yang, W., *J. Chem. Phys.*, **104**, 159-172 (1996).
- [9] Silberstein, L., *Phil. Mag.*, **33**, 92-128 (1917).
- [10] Silberstein, L., *Phil. Mag.*, **33**, 215-222 (1917).
- [11] Silberstein, L., *Phil. Mag.*, **33**, 521-533 (1917).
- [12] Applequist, J., Carl, J. R., Fung. K.-F., *J. Am. Chem. Soc*, **94**, 2952-2960 (1972).
- [13] Applequist, J., *Acc. Chem. Res.*, **10**, 79-85 (1977).
- [14] Judson, R., "Reviews in Computational Chemistry", **10**, Lipkowitz, K. B., Boyd, D. B. (Editors), VCH Publishers, Inc. (New York), 1997.
- [15] Thole, B. T., *Chem. Phys.*, **59**, 341-350 (1981).
- [16] Birge, R. R., *J. Chem. Phys.*, **72**, 5312-5319 (1980).
- [17] Birge, R. R., Schick, G. A., Bocian, D. F., *J. Chem. Phys.*, **79**, 2256-2264 (1983).
- [18] Jensen, L., Åstrand, P.-O., Osted, A., Kongsted, J., Mikkelsen, K. V., *J. Chem. Phys.*, **116**, 4001-4010 (2002).
- [19] Swart, M., Snijders, J. G., van Duijnen, P. Th., *J. Comp. Meth. Sci. Engin.*, **4**, 419-425 (2004).
- [20] Geerlings, P., De Proft, F., Langenaeker, W., *Chem. Rev.*, **103**, 1793-1873 (2003)
- [21] Chelli, R., Procacci, P., Righini, R., Califano, S., *J. Chem. Phys.*, **111**, 8569-8575 (1999)
- [22] Nistor, R. A., Polihronov, J. G., Muser, M. H., Mosey, N. J., *J. Chem. Phys.*, **125**, 094108 (2006)
- [23] Mathieu, D., *J. Chem. Phys.*, **127**, 224103 (2007)
- [24] Chen, J., Hundertmark, D., Martinez, T. J., *J. Chem. Phys.*, **129**, 214113 (2008)

- [25] Chen, J., Martinez, T. J., *Chem. Phys. Lett.*, **438**, 315-320 (2007)
- [26] Stern, H. A., Kaminski, A., Banks, J. L., Zhou, R., Berne, B. J., Friesner, R. A., *J. Phys. Chem. B*, **103**, 4730-4737 (1999)
- [27] Helgaker, T., Jørgensen, P., Olsen, J., "Molecular electronic-structure theory", Wiley, Chichester (2000)
- [28] Reimers, J. R., Hush, N. S., *J. Phys. Chem. B*, **105**, 8979-8988 (2001)
- [29] Mayer, A., *Phys. Rev. B*, **71**, 235333 (2005)
- [30] Margenau, M., Kestner, N. R., "Theory of Intermolecular Forces", Pergamon, Oxford (1969)
- [31] Wheatley, R. J., *Chem. Phys. Lett.*, **294**, 487-492 (1998)
- [32] Wheatley, R. J., Meath, W. J., *Mol. Phys.*, **79**, 253-275 (1993)
- [33] Jensen, L., Åstrand, P.-O., Mikkelsen, K. V., *Int. J. Quant. Chem.*, **84**, 513-522 (2001)
- [34] Mulliken, R. S., *J. Chem. Phys.*, **23**, 1833-1840 (1955)
- [35] Nocedal, J., Wright, S. J., "Numerical Optimization", 2nd ed., Springer (New York), 2006
- [36] Goldberg, D., "Genetic algorithms in search, optimization and learning", Addison-Wesley, Reading, Massachusetts, 1989
- [37] DeJong, K. A., "An analysis of the behavior of a class of genetic adaptive systems", Doctoral Thesis, University of Michigan, 1976.
- [38] van Rossum, G. et al., PYTHON Language Website, www.python.org
- [39] Jones, E., Oliphant, T., Peterson, P., SCIPY: Open source scientific tools for PYTHON, www.scipy.org
- [40] Kongsted, J., Osted, A., Jensen, L., Åstrand, P.-O., Mikkelsen, K. V., *J. Phys. Chem. B*, **105**, 10243-10248 (2001)
- [41] Hansen, T., Jensen, L., Åstrand, P.-O., Mikkelsen, K. V., *J. Chem. Theory Comput.*, **1**, 626-633 (2005)
- [42] Jensen, L., Sylvester-Hvid, K. O., Mikkelsen, K. V., Åstrand, P.-O., *J. Phys. Chem. A*, **107**, 2270-2276 (2003)
- [43] Jensen, L., Åstrand, P.-O., Sylvester-Hvid, K. O., Mikkelsen, K. V., *J. Phys. Chem. A*, **104**, 1563-1569 (2000)
- [44] Champagne, B., Perpete, E. A., Jacquemin, D., van Gisbergen, S. J. A., Baerends, E.-J., Soubra-Ghaoui, C., Robins, K. A., Kirtman, B., *J. Chem. Phys.*, **104**, 4755-4763 (2000)
- [45] Champagne, B., Perpete, E. A., Jacquemin, D., van Gisbergen, S. J. A., Baerends, E.-J., Snijders, J. G., Soubra-Ghaoui, S., Robins, K. A., Kirtman, B., *J. Chem. Phys.*, **109**, 10489-10498 (1998)
- [46] van Faassen, M., de Boeij, P. L., van Leeuwen, R., Berger, J. A., Snijders, J. G., *Phys. Rev. Lett.*, **88**, 186401 (2002)

- [47] van Faassen, M., de Boeij, P. L., van Leeuwen, R., Berger, J. A., Snijders, J. G., *J. Chem. Phys.*, **118**, 1044-1053 (2003)
- [48] Salek, P., Helgaker, T., Vahtras, O., Ågren, H., Jonsson, D., Gauss, J., *Mol. Phys.* **103**, 439-450 (2005)
- [49] Peach, M. J. G., Helgaker, T., Salek, P., Keal, T. W., Lutns, O. B., Tozer, D. J., Handy, N. C., *Phys. Chem. Chem. Phys.*, **8**, 558-562 (2006)
- [50] Champagne, B., Bulat, F. A., Yang, W., Bonness, S., Kirtman, B., *J. Chem. Phys.*, **125**, 194114 (2006)
- [51] Vydrov, O. A., Scuseria, G. E., *J. Chem. Phys.*, **125**, 234109 (2006)
- [52] Sekino, H., Maeda, Y., Kamiya, M., Hirao, K., *J. Chem. Phys.*, **126**, 014107 (2007)
- [53] van Faassen, M., Jensen, L., Berger, J. A., de Boeij, P. L., *Chem. Phys. Lett.*, **395**, 274-278 (2004)
- [54] te Velde, G., Bickelhaupt, F. M., van Gisbergen, S. J. A., Guerra, C. F., Baerends, E. J., Snijders, J. G., Ziegler, T., *J. Comput. Chem.*, **22**, 931-967 (2001)
- [55] Guerra, C. F., Snijders, J. G., te Velde, G., Baerends, E. J., *Theor. Chem. Acc.*, **99**, 391-401 (1998)
- [56] SCM, Theoretical Chemistry, Vrije Universiteit, Amsterdam, The Netherlands, ADF2007.01, ADF2008.01