

PREDICTING
SOLID-STATE QUBIT
MATERIAL HOST

by

Oliver Lerstøl Hebnes

THESIS
for the degree of
MASTER OF SCIENCE



Faculty of Mathematics and Natural Sciences
University of Oslo

March 18, 2021

Contents

I	Introduction	1
1	Introduction	3
1.1	Motivation	4
1.2	Holy grail	4
1.3	Structure of thesis	4
II	Theory	5
2	Semiconductor quantum platforms	7
2.1	Quantum technologies	7
2.1.1	Quantum computation	8
2.1.2	Quantum communication	10
2.1.3	Quantum sensing	11
2.1.4	Quantum computing requirements	11
2.1.5	Available quantum platforms	12
2.2	A brief overview of materials science	14
2.3	Introduction to semiconductor physics	16
2.3.1	Point defects in semiconductors	18
2.3.2	Optical defect transitions	18
2.4	Semiconductor candidates for quantum technology	20
3	Introduction to density functional theory	25
3.1	The Schrödinger equation	25
3.2	The many-particle Schrödinger equation	26
3.3	The Born-Oppenheimer approximation	28
3.4	The Hartree and Hartree-Fock approximation	29
3.5	The variational principle	29
3.6	The density functional theory	30
3.6.1	The Hohenberg-Kohn theorems	31
3.6.2	The Kohn-Sham equation	33
3.6.3	The exchange-correlation energy	35
3.6.4	Self-consistent field methods	37

3.6.5	Limitations of the DFT	37
4	Machine learning	39
4.1	Supervised learning	40
4.2	Evaluating accuracy of a model	40
4.2.1	Bias-variance tradeoff	41
4.2.2	Accuracy, precision and recall	43
4.2.3	Cross-validation	44
4.3	Decision trees	45
4.3.1	Growing a classification tree	46
4.3.2	Classification algorithm	48
4.3.3	Pruning a tree	48
4.3.4	Pros and cons for decision trees	48
4.4	Ensemble methods	49
4.4.1	Bagging	49
4.4.2	Boosting	50
4.5	Dimensionality reduction	53
4.5.1	Principal component analysis	54
III	Methodology and implementation	57
5	Material Science Databases	59
5.1	Fundamentals of a database	59
5.1.1	API and HTTP requests	59
5.1.2	Open-source Python libraries for material analysis . . .	61
5.2	Databases and cloud services	62
5.2.1	Novel Materials Discovery	63
5.2.2	Materials project	63
5.2.3	AFLOW	64
5.2.4	Open Quantum Materials Database	64
5.2.5	JARVIS	65
6	Extraction and featurization of data	67
6.1	Practical data extraction with Python-examples	68
6.1.1	Materials Project	69
6.1.2	Citrine Informatics	70
6.1.3	AFLOW	71
6.1.4	AFLOW-ML	72
6.1.5	JARVIS-DFT	73
6.2	Matminer featurization	74
6.3	Preprocessing	75
6.4	Screen procedure	75

IV	Results	77
7	Validation	79
7.1	The perovskite dataset	79
7.1.1	Features	80
V	Appendices	83
A	Featurization	85
A.1	Table of featurizers	85
A.2	Erroneous entries	90

Part I

Introduction

Chapter 1

Introduction

This is the introduction. Another coffee.

Trengs det egentlig undertitler her? Tenketenketenk

1.1 Motivation

1.2 Holy grail

1.3 Structure of thesis

Part II

Theory

Chapter 2

Semiconductor quantum platforms

This chapter will provide a brief overview of the current state-of-the-art in quantum technological advances. This will not only give us insights in how the technology is being used today, but also grant us the opportunity to discuss key concepts that are fundamental to understand for this thesis. Thereafter we will look into how materials are built up, and what the characteristics of a semiconductor is.

2.1 Quantum technologies

Quantum technology (QT) refers to practical applications and devices that utilize the principles of quantum physics as a foundation. Technologies in this spectrum are based on concepts such as *superposition*, *entanglement* and *coherence*, which are all closely related to one another.

A quantum superposition refers to that any two or more quantum eigenstates can be added together into another valid quantum state, such that every quantum state can be represented as a sum, or a superposition, of two or more distinct states. This is according to the wave-particle duality which states that every particle or another quantum entity may be described as either a particle or a wave. When measuring the state of a system residing in a superposition of eigenstates, however, the system falls back to one of the basis states that formed the superposition, destroying the original configuration.

Quantum entanglement refers to when a two- or many-particle state cannot be expressed independently of the state of the other particles, even when the particles are separated by a significant distance. As a result, the many-particle state is termed an entangled state [1].

Quantum coherence arises if two waves coherently interfere with each other and generate a superposition of the two states with a phase relation. Likewise, loss of coherence is known as *decoherence*.

Another concept that the reader should be familiar with is the famous Heisenberg uncertainty principle. It states that

$$\sigma_x \sigma_p \leq \frac{\hbar}{2}, \quad (2.1)$$

where σ_x is the standard deviation for the position and σ_p is the standard deviation in momentum. This means that we cannot accurately predict both the position and momentum of a particle at the same time. Thus, we often calculate the probability for a particle to be in a state which results in concepts such as an electron sky surrounding an atom core. However, remember that equation (2.1) is an inequality, which means that it is possible to create a state where neither the position nor the momentum is well defined.

2.1.1 Quantum computation

The start of the digital world's computational powers can be credited to Alan Turing. In 1937, Turing [2] published a paper where he described the *Turing machine*, which is regarded as the foundation of computation and computer science. It states that only the simplest form of calculus, such as boolean Algebra (1 for true and 0 for false), is actually computable. This required developing hardware that could handle classical logic operations, and was the basis of transistors that are either in the state ON or OFF depending on the electrical signal. Equipped with a circuit consisting of wires and transistors, commonly known as a computer, we could develop software to solve all kinds of possible applications.

Driven by the development of software, conventional computers have in accordance to Moore's law [3], doubled the amount of transistors on integrated circuit chips every two years as a result of smaller transistors. Furthermore, the clock frequency has enhanced with time, resulting in a doubling of computer performance every 18 months [4]. Alas, miniaturization cannot go on forever as transistors are mass-produced at 5 nm today and are expected to reach a critical limit of 3 nm in the following years [5].

To sustain the digital world's increasing computational demand, other alternatives than the conventional classical computer must be explored. This is where quantum computation comes into the picture. The term quantum computer is a device that exploits quantum properties to solve certain computational problems more efficiently than allowed by Boolean logic [6].

The idea is to pass information in the form of a quantum bit, or *qubit* for short. They are the building blocks of quantum computers, and as opposed to the conventional 0 or 1-bits that classical computers are based on, they can inhabit any superposition of the states 0 or 1. This is illustrated in figure 2.1.

The architecture of a gate-based quantum computer is dependent on a set of quantum logic gates that perform unitary transformations on sets of

qubits [7, 8]. Other implementations of quantum computers exist, such as the adiabatic quantum computer. This approach is not based on gates, but on defining the answer of a problem as the ground state of a complex network of interactions between qubits, and then controlling the interactions to adiabatically evolve the system to the ground state [9].

It has been demonstrated that exponentially complex problems can be reduced to polynomially complex problems for quantum computers [4]. For example, a quantum search algorithm found by Grover [10] offers a quadratic speed-up compared to classical algorithms, while Shor's quantum integer factorization algorithm [11] presents an exponential speed-up. Intriguingly, Google reported in 2019 that they ran a random number generator algorithm on a superconducting processor containing 53 qubits in 200 seconds, which would most likely take several times longer for a classical supercomputer to solve [12]. It is anticipated that quantum computers will excel in exceedingly complex problems, while many simpler tasks may not see any speed-up at all compared to the classical regime. Hence, quantum- and classical computers are envisioned to coexist for each their purpose.

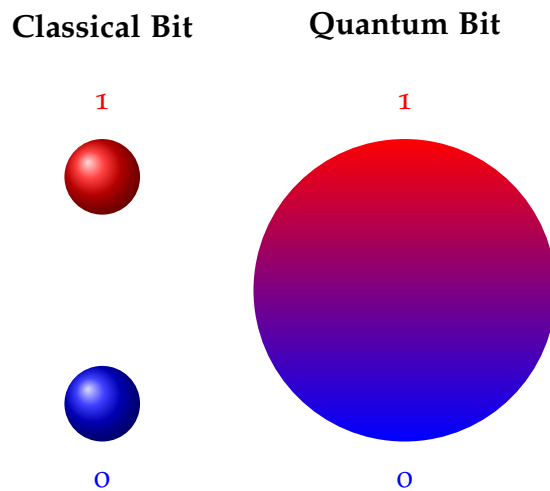


Figure 2.1: Conceptual illustration of the two-level classical bit, which are restricted to the boolean states 1 (true) or 0 (false), and the quantum bit that can be in any superposition of the states 0 or 1.

Quantum computing is a highly sought-after goal, but there are extensive challenges that need to be addressed. Controlling a complex many-qubit system is difficult, since it is not always possible to establish interactions between qubits [7] and maintain entanglement over both time and distance. Additionally, decoherence and other quantum noise occurs as a result of the high volatility of quantum states, making quantum state manipulation prone to errors. The *quantum error correction* protocols and the theory of *threshold theorem* deals with this vulnerability, stating that noise most likely does not pose any fundamental barrier to the performance of large-scale computations [4].

2.1.2 Quantum communication

Quantum communication refers to the transfer of a state of one quantum system to another. Since information can be stored in qubits, we picture *flying qubits* that transfer information from one location to another [13]. The benefits of using flying qubits are in particular valued in quantum cryptography, since the quantum nature of qubits can be exploited to add extra layers of security [4].

Consider the example of encrypting a digitally transmitted conversation. It is difficult to avoid someone eavesdropping on a conversation, however, the problem is diminished if the eavesdropper does not speak the language, keeping the information in the conversation safe. This is the original idea of encryption, such that the information has been encrypted into something incomprehensible for any eavesdropper. A common practice is to encrypt information and share a public key, which everyone can read, and a private key, only known for the sender and receiver of information. This should be sufficient to keep the information secure, given that the complexity of the private key is impenetrable.

Importantly, we live in a digital world where most of our actions are increasingly being stored as information, and we could imagine that the eavesdropper in the latter example stored the conversation. Even if the content of the conversation was encrypted, it still presents a challenge, since encrypted information stored today could be deciphered in ten or twenty years' time. Consequently, finding an encryption method that could make information either impossible to eavesdrop on or make the security unbreakable forever is very desirable. This is the ultimate goal of quantum cryptography [4].

Consider the example of information encoded into a qubit as a superposition of two quantum states. Now, if a wild eavesdropper would try to measure the information, the nature of quantum physics tells us that the original configuration would be destroyed and the receiver would be alerted of the eavesdropper. Furthermore, if the eavesdropper would try to make a copy of the message, the copying itself would be limited of the no-cloning theorem [14] which declare that quantum states cannot be copied.

A clever approach to ensure confidentiality is to send the encryption key before sending the actual encrypted information. If the key is received unperturbed, the key remains secret and can be safely employed. If it turns out perturbed, confidentiality is still intact since the key does not contain any information and can be discarded. This approach is termed the *quantum key distribution* (QKD) [14, 15]. It should be noted that this requires both the sender and receiver to have access to methods for sending, receiving and storing qubit states, such as a quantum computer. Additionally, the sender and receiver will need to initially exchange a common secret which is later expanded, making quantum key *expansion* a more exact term for QKD [4, 15].

Most applications and experiments use optical fibers for sending information via photons, with the distance regarded as the main limitation. This is because classical repeaters are unable to enhance quantum information because of the no-cloning theorem, making photon loss in optical fiber cables inevitable. Thus, quantum communication must reinvent the repeater concept, using hardware that preserves the quantum nature [16] and are compatible with wavelengths used in telecommunication. Nonetheless, secure QKD up to 400 km has recently been demonstrated using optical fibres in academic prototypes [17].

2.1.3 Quantum sensing

Measurements are part of our digital world today to a great extent. There would be no way to exchange goods, services or information without reliable and precise measurements [16]. Thus, improving the accuracy of sensors for every measurement done is desirable. One method to improve measurement accuracy, resolution and sensitivity can be by utilizing quantum sensors. Quantum sensors exploit quantum properties to measure a physical quantity [18]. This is possible because quantum systems are highly susceptible to perturbations to its surroundings, and can be used to detect physical properties such as either temperature or an electrical or magnetic field [18].

For a quantum system to be able to function as a quantum sensors, a few criterias needs to be met. Firstly, the quantum system needs to have discrete and resolvable energy levels. The quantum system also needs to be controllably initialised into a state that can be identified and coherently manipulated by time-dependent fields. Lastly, the quantum system needs to be able to interact with the physical property one wants to measure through a coupling parameter [18].

It is also possible to also exploit quantum entanglement to improve the precision of a measurement. This gain of precision is used to reach what is called the Heisenberg-limit, which states that the precision scales as the number of particles N in an idealized quantum system [16, 18], while the best classical sensors scale with \sqrt{N} .

2.1.4 Quantum computing requirements

As ever-promising the concepts of quantum technology are, the physical realizations are in the preliminary stage of development. Here we will concretize critical principles for a physical realisation of a quantum platform.

“I always said that in some sense, these criteria are exactly the ones that you would teach to kindergarten children about computers, quantum or otherwise” DiVincenzo [19]

DiVincenzo formulated in the year of 2000 seven basic criteria for a physical qubit system with a logic-based architecture [7].

1. A scalable physical system with well characterized qubits
2. The ability to initialize the state of the qubits to a simple initial system
3. Have coherence times that are much longer than the gate operation time
4. Have a universal set of quantum gates
5. Have the ability to perform qubit-specific measurements
6. The ability to convert stationary qubits to flying qubits
7. The ability to faithfully transmit flying qubits between specified locations

The first five criteria (1-5) must be met for a quantum platform to be considered a quantum computer, while the two last criteria (6-7) were added for quantum communication, since its applications provide a unique advantage compared to its classical counterpart.

2.1.5 Available quantum platforms

Many different quantum platforms have been physically implemented, and this section will serve as a brief overview of the current status. For a more thorough review of qubit implementations, the reader is directed to Refs. [8, 16].

Superconducting circuits can be used in quantum computing, since electrons in superconducting materials can form Cooper pairs via an effective electron-electron attraction when the temperature is lower than a critical limit. Below the limit, electrons can move without resistance in the material [20]. Exploiting this intrinsic coherence, qubits can be made by forming microwave circuits based on loops of two superconducting elements separated by an insulator, also known as Josephson tunnel junctions [16]. Today, superconducting Josephson junctions are the most widely used quantum platform, but they require very low temperature (mK) to function, making them costly to use. Additionally, the current devices experience a relatively short coherence time, causing challenges in scaling up.

Single photons is an eligible quantum platform that can be implemented as qubits with one-qubit gates being formed by rotations of the photon polarization. Its use in fiber optics are less prone to decoherence, but faces

challenges since the more complex photon-photon entanglement and control of multi-qubits is strenuous [8].

By fixing the nuclear spin of solid-state systems, it is possible to implement a quantum platform that experience long spin coherence. This enables the manipulation of qubits that utilize electromagnetic fields, making one-qubit gates realizable.

The isolated atom platform is characterized by its well-defined atom isolation. Here, every qubit is based on energy levels of a trapped ion or atom. Quantum entanglement can be achieved through laser-induced spin coupling, however scaling up to large atom numbers induce problems in controlling large systems and cooling of the trapped atoms or ions.

A quantum dot (QD) can be imagined as an artificial atom which is confined in a solid-state host. As an example, a quantum dot can occur when a hole or an electron is trapped in the localized potential of a semiconductor's nanostructure. QDs exhibit similar coherence potential as the isolated atom platform, but without the drawback of confining and cooling of the given atom or ion [16]. Moreover, it is possible to limit decoherence due to nuclear spins by dynamic decoupling of nuclear spin noise and isotope purification [8].

A QD can normally be defined lithographically using metallic gates, or as self-assembled QDs where a growth process creates the potential that traps electrons or holes. The difference between them is a question of controllability and temperature, since the metallic gates is primarily controlled electrically and operate at < 1 K, while self-assembly QDs are primarily controlled optically at ~ 4 K [8]. Despite requiring very low temperatures, QDs have the potential for fast voltage control and optical initialization. As with trapped ions, electrostatically defined quantum dots experience a short-range exchange interaction, imposing a limitation for quantum computing and quantum error correction protocols. A potential solution could include photonic connections between quantum dots. On the contrary, self-assembled quantum dots couple strongly to photons due to their large size in comparison to single atoms. However, the size and shapes of self-assembled quantum dots are decided randomly during the growth process, causing an unfavourable large range of optical absorption and emission energies [8].

Lastly, we will turn towards point defects in bulk semiconductors as a physical implementation of a quantum platform. Point defects shares many of the attributes of quantum dots, such as discrete optical transitions and controllable coherent spin states, but are vulnerable to small changes in the lattice

of the semiconductor. Thus, it can be difficult to isolate a point defect from the surrounding environment. However, one can utilize the strength of the solid-state semiconductor host to isolate to some extent the point defect, yielding extended coherence times and greater optical homogeneity than other quantum dot systems. Before we dwell into the intricacies of point defect qubits as a building block for QT, we will provide the necessary background for the crystal- and electronic structure of semiconductors.

2.2 A brief overview of materials science

The interactions between atoms and characteristics of matter form the foundation of materials science. The applications of materials science are extensive, with examples such as a bottle of water or to a chair to sit in.

Solid materials, like plastic bottles, are formed by densely packed atoms. These atoms can randomly occur through the material without any long-range order, which would categorize the material as an *amorphous solid*. Amorphous solids are frequently used in gels, glass and polymers [21]. However, the atoms can also be periodically ordered in small regions of the material, classifying the material as a *polycrystalline solid*. All ceramics are polycrystalline with a broad specter of applications ranging from kitchen-porcelain to orthopedical bio-implants [22]. A third option is to have these atoms arranged with infinite periodicity, making the material a *crystalline solid* or more commonly named a *crystal*. The three options are visualised in figure 2.2. Hereon, we will focus on crystalline solids.

The periodicity in a crystal is defined in terms of a symmetric array of points in space called the *lattice*, which can be simplified as either a one-dimensional array, a two-dimensional matrix or a three dimensional vector space, depending on the material. At each lattice point we can add an atom to make an arrangement called a *basis*. The basis can be one atom or a cluster of atoms having the same spatial arrangement. Every crystal has periodically repeated building blocks called *cells* representing the entire crystal. The smallest cell possible is called a *primitive cell*, but such a cell only allows lattice points at its corners and it is often quite rigid to work with when the structure becomes complex. As a solution, we will consider the *unit cell*, which allows lattice points on face centers and body centers.

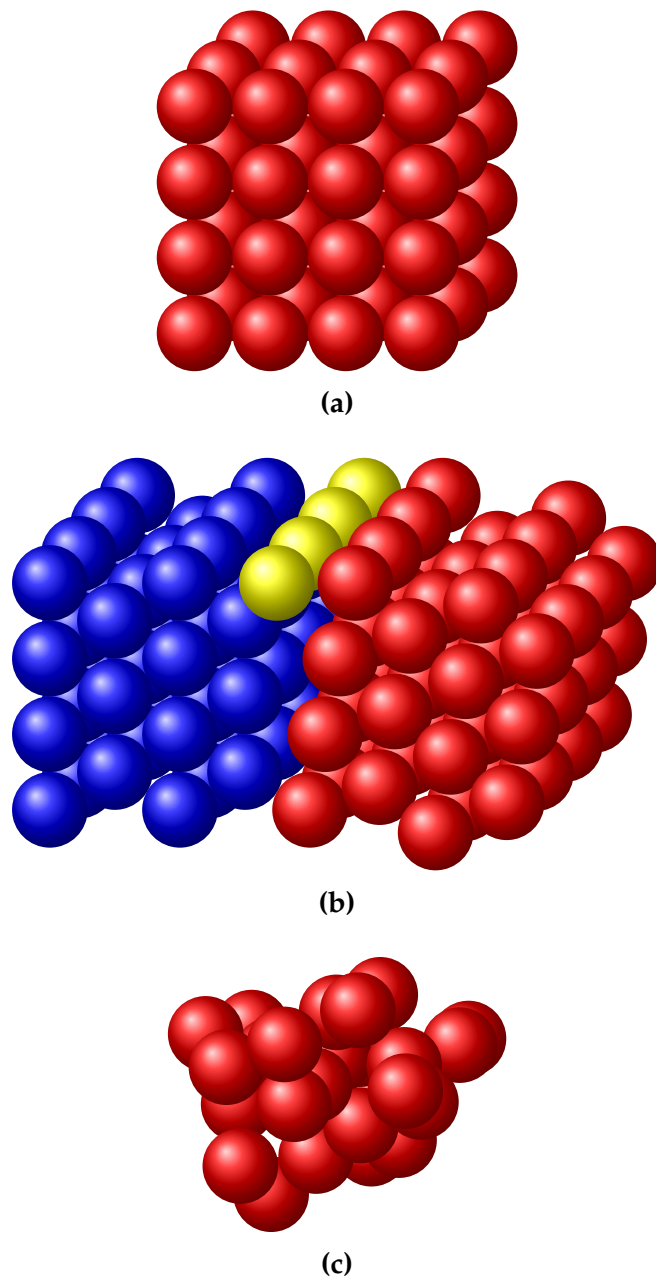


Figure 2.2: Schematic representation of different degrees of ordered structures, where (a) is a crystalline of a simple cubic lattice, (b) is a polycrystalline hexagonal lattice, and (c) is an amorphous lattice.

One example of a crystal structure is the perovskite structure. Compounds with this structure are characterized by having an ABX_3 stoichiometry whose symmetries belong to one of 15 space groups identified by Lufaso & Woodward [23], such as the cubic, orthorhombic and tetragonal. For our purpose, we will be looking into when the X atom is oxygen, and refer to the oxygen-perovskite ABO_3 . The A atom is nine- to 12-fold coordinated by oxygen, while the B atom is sixfold coordinated by oxygen, and the BO_6 octahedra are connected to the corners in all three directions as visualized in figure 2.3.

The motivation behind the research on perovskites is related to the large amount of available ABO_3 chemistries, where a significant portion of these take the perovskite structure. Perovskites have a broad specter of applications, ranging from high-temperature superconductors [24] and ionic conductors [25] to multiferroic materials [26]. Additionally, adding a perovskite-type compound to solar cells has reportedly resulted in higher performance efficiencies while being cheap to produce and simple to manufacture [27, 28]. However, this includes the use of hybrid organic-inorganic compounds and excludes the use of oxygen.

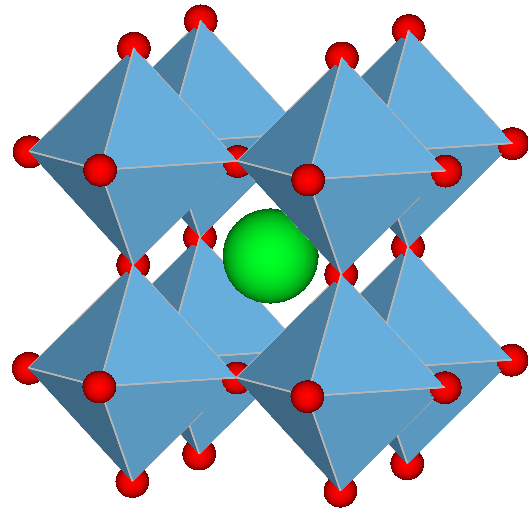


Figure 2.3: A crystal structure of $SrTiO_3$, which is a cubic perovskite. The red atoms are oxygen, whereas the green atom is strontium, and inside every corner-sharing BO_6 octahedral unit is a titanium atom.

2.3 Introduction to semiconductor physics

Isolated atoms have distinct energy levels, where the Pauli exclusion principle [29] states for fermions that each energy level can at most accommodate two electrons of opposite spin. In a solid, the discrete energy levels of the isolated atom spread into continuous energy bands since the wavefunctions of the electrons in the neighboring atoms overlap. Hence, an electron is not necessarily localized at a particular atom anymore. Every material has a unique band structure, similar to every human having their unique fingerprint.

Knowing which energy bands are occupied by electrons is the key in understanding the electrical properties of solids. The highest occupied electron band at 0 K is called the valence band (VB), while the lowest unoccupied electron band is called the conduction band (CB). The energy gap of forbidden energy levels between the maximum VB and the minimum CB is known as the band gap, and its energy is denoted as E_g . If a material can be classified as a semiconductor depends on the band gap and the electrical conductivity. As an example, Silicon is commonly thought of as a semiconductor, and has a band gap of about 1.12 eV at 275 K [30].

To be able to accelerate electrons in a solid using an electrical field, they must be able to move into new energy states. At 0 K, the entire valence band of a semiconductor is full with electrons and there are no available states nearby, making it impossible for current to flow through the material. This can be solved by using either thermal or optical energy to excite electrons from the valence band to the conduction band, in order to *conduct* electricity. At room temperature, some semiconductors will have electrons excited to the conduction band solely from thermal energy matching the energy band gap [21].

In some scenarios, thermal or optical energy is not sufficient for an excitation since the energy bands are also dependent on the crystal momentum. A difference in the momentum of the minimal-energy state in the conduction band and the maximum-energy state in the valence band results in an *indirect bandgap* as seen in figure 2.4a. If there is no difference at all, the material has a *direct bandgap*, which is visualized in figure 2.4b.

Electrons in semiconductor materials can be described according to the Fermi-Dirac distribution

$$f(E) = \frac{1}{1 + e^{(E-E_F)/kT}},$$

where k is Boltzmann's constant, T is temperature, E is the energy and E_F is the Fermi level. The Fermi-Dirac distribution gives the probability that a state will be occupied by an electron, and at $T = 0$ K, every energy state lower than E_F is occupied by electrons while the opposite is true for energy states above E_F [21].

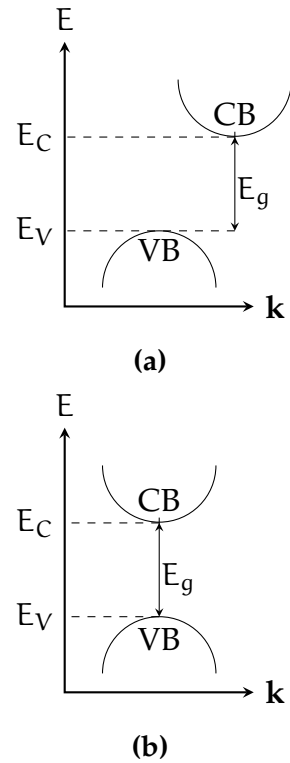


Figure 2.4: A schematic drawing of an (a) indirect- and a (b) direct bandgap.

2.3.1 Point defects in semiconductors

In real life, a perfect crystal without any symmetry-breaking flaw does not exist. These flaws are known as defects and can occur up to three dimensions. An example one-dimensional defect is known as a *line defect*, while two dimensional defects can be *planar defects*, and in three dimensions we have *volume defects*. Lastly, defects can also occur in zero dimensions and are then termed *point defects*. Point defects normally occur as either vacancies, interstitial placement inbetween lattice sites or as substitution of another existing atom in the lattice.

Defects can greatly influence both the electronic and optical properties of a material. A substitutional defect can at first be regarded as an impurity or an antisite, but they can also be intentionally inserted, an approach known as *doping*. Doping can result in an excess of electrons or holes, making the semiconductor either an n- or p-type, respectively. Consequently, the semiconductor will have energy levels in the (forbidden) band gap that originates from the defects. If the energy levels introduced are closer than ~ 0.2 eV to the band edges, they are termed *shallow* defects.

Shallow defects can contribute with either excess electrons to the conduction band, or excess holes to the valence band. However, the induced charge carriers (electrons or holes) interact strongly with the band edges, resulting in a delocalized wavefunction regarding the position in the lattice.

For the opposite case, if the energy levels rests closer to the middle of the semiconductor's gap, the introduced defects are known as *deep level* defects. Deep levels normally occur due to either dangling bonds or impurities, and have highly localized electron wavefunctions. This might assure the isolation required for long coherence times, which is an appealing promise in quantum technological advances.

Deep levels can be unfortunate in semiconductors since they can interact with the charge carriers, potentially destroying the desired electronic or optical property of the material. Deep level defects can function as electron-hole recombination centers, or to trap charge carriers, yielding the commonly used name deep level *traps*. Both of the given situations results in a lower concentration of charge carriers, which showcase why deep levels can be unwanted in semiconductor devices. However, deep level defects show extraordinary properties in Q

2.3.2 Optical defect transitions

Optical transitions refers to excitation of charge carriers due to either emission or absorption of electromagnetic radiation, and can be done with a laser light or electron beam. Figure 2.5 represents a configuration coordinate (CC) diagram of a defect transition. The y-axis is a function of the energy E , while

the x-axis is a function of the configuration coordination Q . The lowest point in the lower parabola is known as the ground state (GS) configuration Q_{GS} , which is the most stable atomic position, while for the upper parabola it is known as the excited state configuration Q_{ES} . The dotted lines represent vibronic excitations to the energy of the ground state Q_{GS} for the lower parabola, while it represents Q_{ES} for the higher parabola.

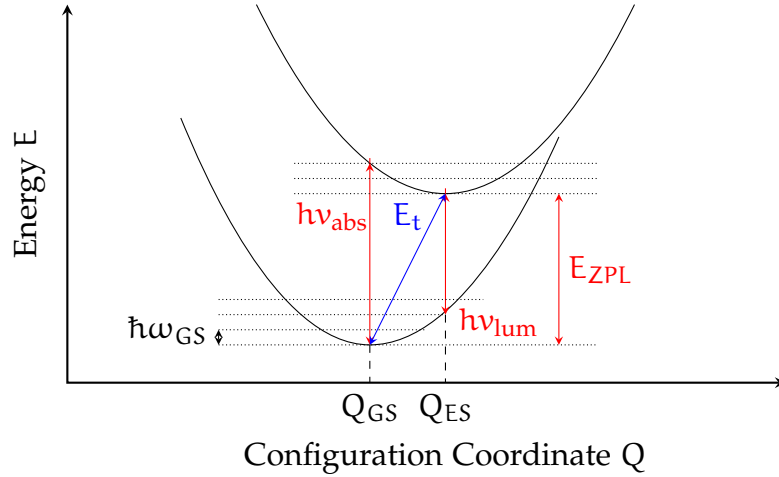


Figure 2.5: A schematic representation of a configuration coordination diagram based on Ref. [31].

The optical transitions in figure 2.5 are marked with red arrows. During slow transitions, such as during thermodynamic defect transitions, the original configuration have time to rearrange due to phonon vibrations. This is schematically drawn as the blue arrow, where the energy E_t equals the ionization energy or the position of the defect level. Optical transitions, on the other hand, are marked in red and occur in a short time range such that the original configuration does not change. They can appear in the exchange of charge carriers with the band edges, and in a defect's internal excited state, with the latter scenario being most relevant for this thesis.

Consider a defect that rests in the ground state configuration Q_{GS} . Suddenly, it absorbs a photon with energy $h\nu_{abs}$ and occupies an excited vibronic state of the upper parabola after a vertical transition. Through lattice reconfigurations, the defect will move towards the bottom of the upper parabola, also known as Q_{ES} . Eventually, it will relax to the lower parabola by emitting a photon with energy $h\nu_{lum}$, also known as a zero-phonon line (ZPL) of energy E_{ZPL} . On the other hand, any transitions between vibronic excitation levels are phonon-related. How strong the electron-phonon interaction is can be quantified by the Huang-Rhys factor S [32]. If the two parabolas in figure 2.5 have the same configuration of Q , emission into the ZPL is enabled and $S \sim 0$. The stronger the coupling, the smaller amount of emission in the ZPL.

The optical properties of a host material can be greatly influenced by defects, in particular the ES to GS transition that can occur in a defect, as discussed for figure 2.5. If the defect were to facilitate the emission of single photons with a detectable time inbetween together with a distinguishable ZPL, the defect would be referred to as a single photon source (SPS). The criteria for SPS are not met in many materials, since charge-state transitions often comprise interactions with either the VB or the CB. Thus, most SPSs' GS and ES levels are situated within the band gap of a host material. Consequently, mostly wide-band gap semiconductors are used as host materials for SPSs.

2.4 Semiconductor candidates for quantum technology

The properties of point defects are promising in a quantum technological perspective. We have seen that point defects facilitate deep energy levels within the band gap of the semiconductor, and provide isolation in the solid-state matrix as a result from a high degree of localization of the defect orbitals. If the host material have a small spin-orbit coupling, it could provide long coherence times for a deep level trap in localized and high-spin states. Additionally, point defects have the potential to be single-photon sources, giving rise to sharp and distinguishable optical transitions, where a significant amount of the emission can be of the energy E_{ZPL} . This is in particular seen in wide bandgap semiconductors, and combined with a weak electron-phonon interaction, can have the capacity to be fabricated as a high-fidelity SPS with a significant ZPL part.

The most studied point defect system is the nitrogen-vacancy (NV^{-1}) in diamond. Figure 2.6 schematically shows the different stages of constructing the negative charge state. Panel 2.6a shows the electronic states that correspond to the difference for an isolated atom and a lattice of atoms, as a superposition of sp^3 orbitals that generates valence and conduction bands. In panel 2.6b, a vacancy has been created by removing a carbon atom, and the four orbitals interact with each other resulting in two new states with a_1 and t_2 symmetry due to dangling bonds. Substituting a carbon atom with a nitrogen atom further splits the t_2 -states into two new states. The states $a(1)$ and e_x/e_y are of importance, as they are the GS and the ES of the qubit defects, respectively. Here, an optical spin-conserving transition can occur due to a laser light of correct wavelength [33], as exemplified from the discussion from the last section.

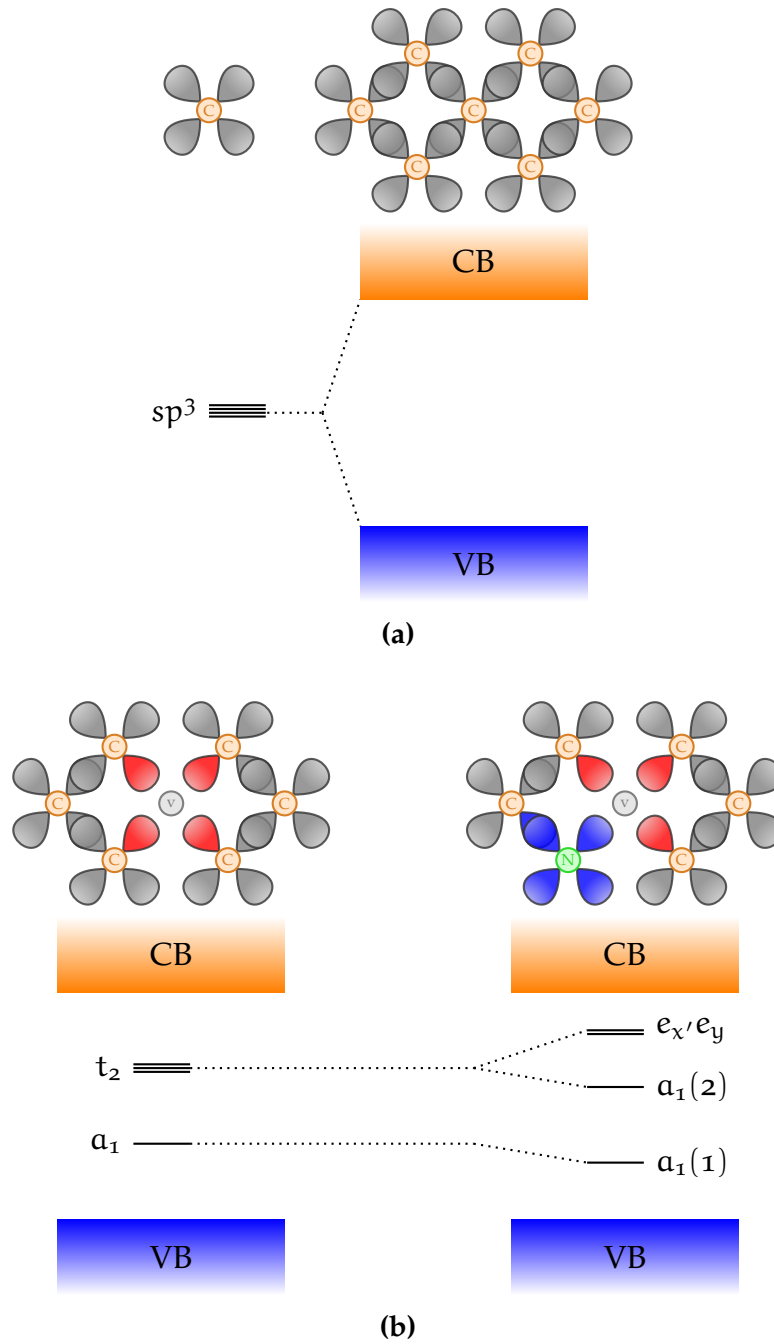


Figure 2.6: A schematic representation of the electronic structure of the NV^{-1} defect in a tetrahedrally coordinated semiconductor, exemplified by diamond. Figure used from Ref. [33].

The nitrogen-vacancy in diamond is a prominent single-photon source up to room temperatures. This involves initializing, manipulating and reading out of the qubit state using optical and electric excitations, and electric and magnetic fields [33]. The potential qubit system have promising applications in quantum- communication and computation, with a demonstrated entanglement between two NV center spins that are separated by 3 m [34]. Nevertheless, perhaps the most propitious application can be seen in quantum sensing as high-sensitivity magnetometer with nanoscale resolution [35].

Unfortunately, the NV-center display restricted capabilities for quantum communication and computation. The amount of emission into the zero-phonon line is 4% at 6 K [36], which is low. The emission of the qubit center is not completely compatible with current optical fiber technologies, since the emission is in the red wave-length specter. Additionally, fabricating materials of diamond is far from unchallenging and serves as a significant incentive to find other promising qubit candidates.

Therefore, we turn to the search of other qubit systems that offers similar capabilities, but that are more user-friendly. In particular, we need to search for new promising materials that can host a potential qubit. Weber *et al.* [6] proposed in 2010 four criteria that should be met for a solid-state semiconductor material hosting a qubit defect, whereas some of the criteria has already been discussed. An ideal crystalline host should have [6]

- (H1) A wide band gap to accomodate a deep center.
- (H2) Small spin-orbit coupling in order to avoid unwanted spin flips in the defect bound states.
- (H3) Availability as high-quality, bulk, or thin-film single crystals.
- (H4) Constituent elements with naturally occuring isotopes of zero nuclear spin.

Table (2.1) lists several material host candidates that exhibit promising band gap capable of accommodating a deep level defect. The spin-orbit splitting is an indication of the strength of the spin-orbit interaction, and is taken at the Γ point from the valence-band splitting. A smaller value may indicate less susceptibility to decoherence.

Criterion (H3) is important for scalability and further potential for a large-scale fabrication. The given candidate hosts provided in table (2.1) can all be grown as single crystals, but with varying quality and size.

Normally, nuclear spin is a major source of decoherence for all semiconductor-based quantum technologies. This would exclude the use of all elements in odd groups in the periodic table, since these elements exhibit nonzero nuclear spin. As a result, the spin-coherence time of a paramagnetic deep center [6]

might increase. However, nuclear spin can also induce additional quantum degrees of freedom for applications in the right configuration [37]. Therefore, criterion (H₄) is not a strict requirement but is a general recommendation for reducing decoherence time.

Weber *et al* [6] use criteria (H₁) – (H₄) to specifically find analogies to the NV⁻¹ center in other material systems, thus leaving the discussion of other criteria out, such as the choice of crystal system. The atomic configuration and crystal structure of a material strongly influences the properties of a defect, since a defect's orbital and spin structure is dependent on its spatial symmetry [37]. In particular, it is the point group that decides which multiplicity a given energy level should have [38]. A higher defect symmetry group generally facilitates degenerate states, which may give rise to high spin states according to Hund's rules [37, 39]. Inversion symmetry in the host crystal can also be beneficial, resulting in reduced inhomogenous broadening and spectral diffusion of optical transitions as a consequence of being generally insensitive to external electric fields [37].

Material	Band gap E_g (eV)	Spin-orbit splitting Δ_{so} (meV)	Stable spinless nuclear isotopes?
3C-SiC	2.39	10	Yes
4C-SiC	3.26 [40]	6.8	Yes
6C-SiC	3.02	7.1	Yes
AlN	6.13	19 [41]	No
GaN	3.44	17.0	No
AlP	2.45	50 [42]	No
GaP	2.27	80	No
AlAs	2.15	275	No
ZnO	3.44 [43]	-3.5	Yes
ZnS	3.72 [44]	64	Yes
ZnSe	2.82	420	Yes
ZnTe	2.25	970	Yes
CdS	2.48	67	Yes
C (Diamond)	5.5	6	Yes
Si	1.12	44	Yes

Table 2.1: Table taken from Gordon *et al.* [33] that lists a number of tetrahedrally coordinated hosts whose band gaps are larger than 2.0 (eV), and compares it to diamond and Si. All experimental values are from Ref. [30], except for where explicitly cited otherwise.

Chapter 3

Introduction to density functional theory

To fully understanding the underlying physics behind computational material science, we will need to investigate how we can calculate the forces acting inside a crystal. Since these forces are happening on a microscopic scale, we will need to utilise the theory of quantum mechanics.

In this chapter, we will only summarize the necessary theory behind density functional theory, leaving most of the quantum-mechanical world untouched. However, the fundamental theory remains the same and we will start our venture with the Schrödinger equation.

3.1 The Schrödinger equation

In principle, we can describe all physical phenomena of a system with the wavefunction $\Psi(\mathbf{r}, t)$ and the Hamiltonian $\hat{H}(\mathbf{r}, t)$, where \mathbf{r} is the spatial position and t is the time. Unfortunately, analytical solutions for the time-dependent Schrödinger equation,

$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{r}, t) = \hat{H}(\mathbf{r}, t) \Psi(\mathbf{r}, t), \quad (3.1)$$

are extremely rare. More conveniently, we can generate a general wavefunction by a summation of eigenfunctions,

$$\Psi(\mathbf{r}, t) = \sum_{\kappa} c_{\kappa} \psi_{\kappa}(\mathbf{r}, t), \quad (3.2)$$

where c_{κ} is a constant and ψ_{κ} is the κ -th eigenfunction. A general wavefunction does not necessarily describe stationary states, and consequently

does not have distinct energies but is rather represented statistically from the expectation value

$$E = \sum_{\kappa} |c_{\kappa}|^2 E_{\kappa}. \quad (3.3)$$

Solving the Schrödinger equation for a general wavefunction is rather troublesome, but luckily we can use the eigenfunctions instead, transforming equation 3.1 into the time-independent Schrödinger equation for eigenfunctions

$$\hat{H}\psi_{\kappa}(\mathbf{r}) = E_{\kappa}\psi_{\kappa}(\mathbf{r}), \quad (3.4)$$

where E_{κ} is the eigenvalue of the κ -th eigenstate $\psi_{\kappa}(\mathbf{r})$. The eigenfunctions have distinct energies, and the state with the lowest energy is called the ground state. They have the attribute that they are orthogonal and normalized with respect to

$$\langle \psi_{\kappa}(\mathbf{r}) | \psi_{\kappa'}(\mathbf{r}) \rangle = \delta_{\kappa\kappa'}. \quad (3.5)$$

The symmetry of an eigenfunction depends on the symmetry of the potential $V_{\text{ext}}(\mathbf{r})$ and the boundary conditions [45].

3.2 The many-particle Schrödinger equation

As we extend the theory to include many-particle systems, we will gradually explain and add the different contributions that make up the many-body Hamiltonian. During this process, we will neglect any external potential applied to the system.

If we place a simple electron with mass m_e in its own system, it will be in possession of kinetic energy. Instead of just one electron, we can place N_e electrons, and they will together have the total kinetic energy

$$T_e = - \sum_{j=1}^{N_e} \frac{\hbar^2 \nabla_j^2}{2m_e}. \quad (3.6)$$

All the electrons are negatively charged, causing repulsive Coulomb interactions between each and every electron, totalling to

$$U_{ee} = \sum_{j=1}^{N_e} \sum_{j' < j} \frac{q^2}{|r_j - r_{j'}|}. \quad (3.7)$$

The summation voids counting each interaction more than once. Simultaneously, we can place N_n nuclei with mass m_n in the same system, accumulating the kinetic energy

$$T_n = - \sum_{a=1}^{N_n} \frac{\hbar^2 \nabla_a}{2m_n}. \quad (3.8)$$

As in the example with electrons, the nuclei are also experiencing repulsive interactions between every single nucleus, adding up the total interactions as

$$U_{nn} = \sum_{a=1}^{N_n} \sum_{a' < a} \frac{q^2 Z_a Z_{a'}}{|R_a - R_{a'}|}. \quad (3.9)$$

where Z_a is the atom number of nuclei number a .

The system now contains N_e electrons and N_n nuclei, thus we need to include the attractive interactions between the them,

$$U_{en} = - \sum_{j=1}^{N_e} \sum_{a=1}^{N_n} \frac{q^2 Z_a}{|r_j - R_a|}. \quad (3.10)$$

Together, these equations comprise the time-independent many-particle Hamiltonian

$$\begin{aligned} \hat{H} = & - \sum_{j=1}^{N_e} \frac{\hbar^2 \nabla_j}{2m_e} - \sum_{a=1}^{N_n} \frac{\hbar^2 \nabla_a}{2m_n} + \sum_{j=1}^{N_e} \sum_{j' < j} \frac{q^2}{|r_j - r_{j'}|} \\ & + \sum_{a=1}^{N_n} \sum_{a' < a} \frac{q^2 Z_a Z_{a'}}{|R_a - R_{a'}|} - \sum_{j=1}^{N_e} \sum_{a=1}^{N_n} \frac{q^2 Z_a}{|r_j - R_a|}. \end{aligned} \quad (3.11)$$

A few problems arise when trying to solve the many-particle Schrödinger equation. Firstly, the amount of atoms in a crystal is very, very massive. As an example, we can numerically try to calculate the equation 3.7 for a 1mm^3 silicon-crystal that contains $7 \cdot 10^{20}$ electrons. For this particular problem, we will pretend to use the current fastest supercomputer Fugaku [46] that can calculate 514 TFlops, and we will assume that we need 2000 Flops to calculate each term inside the sum [45], and we need to calculate it $N_e \cdot N_e/2$ times for the (tiny) crystal. The entire electron-electron interaction calculation would take $2.46 \cdot 10^{19}$ years to finish for a tiny crystal. Thus, the large amount of particles translates into a challenging numerical problem.

Secondly, the many-particle Hamiltonian contains operators that has to be applied to single-particle wavefunctions, and we have no prior knowledge of how Ψ depends on the single-particle wavefunctions ψ_k .

3.3 The Born-Oppenheimer approximation

The many-particle eigenfunction describes the wavefunction of all the electrons and nuclei and we denote it as Ψ_{κ}^{en} for electrons (e) and nuclei (n), respectively. The Born-Oppenheimer approximation assumes that nuclei, of substantially larger mass than electrons, can be treated as fixed point charges. According to this assumption, we can separate the eigenfunction into an electronic part and a nuclear part,

$$\Psi_{\kappa}^{en}(\mathbf{r}, \mathbf{R}) \approx \Psi_{\kappa}(\mathbf{r}, \mathbf{R})\Theta_{\kappa}(\mathbf{R}), \quad (3.12)$$

where the electronic part is dependent on the nuclei. This is in accordance with the assumption above, since electrons can respond instantaneously to a new position of the much slower nucleus, but this is not true for the opposite scenario. To our advantage, we already have knowledge of the terms in the many-particle Hamiltonian, and we can begin by separating the Hamiltonian into electronic and nuclear parts:

$$\hat{H}^{en} = \overbrace{T_e + U_{ee} + U_{en}}^{\hat{H}^e} + \overbrace{T_n + U_{nn}}^{\hat{H}^n}. \quad (3.13)$$

Starting from the Schrödinger equation, we can formulate separate expressions for the electronic and the nuclear Schrödinger equations.

$$\hat{H}^{en}\Psi_{\kappa}^{en}(\mathbf{r}, \mathbf{R}) = E_{\kappa}^{en}\Psi_{\kappa}^{en}(\mathbf{r}, \mathbf{R}) \quad |\times \int \Psi^*(\mathbf{r}, \mathbf{R})d\mathbf{r} \quad (3.14)$$

$$\int \Psi_{\kappa}^*(\mathbf{r}, \mathbf{R})(\hat{H}^e + \hat{H}^n)\Psi_{\kappa}(\mathbf{r}, \mathbf{R})\Theta_{\kappa}(\mathbf{R})d\mathbf{r} = E_{\kappa}^{en} \underbrace{\int \Psi_{\kappa}^*(\mathbf{r}, \mathbf{R})\Psi_{\kappa}(\mathbf{r}, \mathbf{R})d\mathbf{r}}_1 \Theta_{\kappa}(\mathbf{R}). \quad (3.15)$$

Since $\Theta_{\kappa}(\mathbf{R})$ is independent of the spatial coordinates to electrons, we get E_{κ} as the total energy of the electrons in the state κ .

$$E_{\kappa}(\mathbf{R})\Theta_{\kappa}(\mathbf{R}) + \int \Psi_{\kappa}^*(\mathbf{r}, \mathbf{R})\hat{H}^n\Psi_{\kappa}(\mathbf{r}, \mathbf{R})\Theta_{\kappa}(\mathbf{R})d\mathbf{r} = E_{\kappa}^{en}\Theta_{\kappa}(\mathbf{R}). \quad (3.16)$$

Now, the final integration term can be simplified by using the product rule, which results in

$$\left(T_n + T_n' + T_n'' + U_{nn} + E_{\kappa}(\mathbf{R})\right)\Theta_{\kappa}(\mathbf{R}) = E_{\kappa}^{en}\Theta_{\kappa}(\mathbf{R}). \quad (3.17)$$

If we neglect T'_n and T''_n to lower the computational efforts, we obtain the Born-Oppenheimer approximation with the electronic eigenfunction as

$$(T_e + U_{ee} + U_{en}) \Psi_\kappa(\mathbf{r}, \mathbf{R}) = E_\kappa(\mathbf{R}) \Psi_\kappa(\mathbf{r}, \mathbf{R}) \quad (3.18)$$

and the nuclear eigenfunction as

$$(T_n + U_{nn} + E_\kappa(\mathbf{R})) \Theta_\kappa(\mathbf{R}) = E_\kappa^{en}(\mathbf{R}) \Theta_\kappa(\mathbf{R}). \quad (3.19)$$

How are they coupled, you might ask? The total energy in the electronic equation is a potential in the nuclear equation.

3.4 The Hartree and Hartree-Fock approximation

The next question in line is to find a wavefunction $\Psi(\mathbf{r}, \mathbf{R})$ that depends on all of the electrons in the system. The Hartree [45] approximation to this is to assume that electrons can be described independently, suggesting the *ansatz* for a two-electron wavefunction

$$\Psi_\kappa(\mathbf{r}_1, \mathbf{r}_2) = A \cdot \psi_1(\mathbf{r}_1) \psi_2(\mathbf{r}_2), \quad (3.20)$$

where A is a normalization constant. This approximation simplifies the many-particle Schrödinger equation a lot, but comes with the downside that the particles are distinguishable and do not obey the Pauli exclusion principle for fermions.

The Hartree-Fock approach, however, overcame this challenge and presented an anti-symmetric wavefunction that made the electrons indistinguishable [1]:

$$\Psi_\kappa(\mathbf{r}_1, \mathbf{r}_2) = \frac{1}{\sqrt{2}} (\psi_1(\mathbf{r}_1) \psi_2(\mathbf{r}_2) - \psi_1(\mathbf{r}_2) \psi_2(\mathbf{r}_1)). \quad (3.21)$$

For systems containing more than one particles, the factor $1/\sqrt{2}$ becomes the Slater determinant and is used to normalize the wave function.

3.5 The variational principle

So far, we have tried to make the time-independent Schrödinger equation easier with the use of an *ansatz*, but we do not necessarily have an adequate guess for the eigenfunctions and the *ansatz* can only give a rough estimate in most scenarios. Another approach, namely the *variational principle*, states that the energy of any trial wavefunction is always an upper bound to the exact ground state energy by definition E_0 .

$$E_0 = \langle \psi_0 | H | \psi_0 \rangle \leq \langle \psi | H | \psi \rangle = E \quad (3.22)$$

The eigenfunctions of H form a complete set, which means any normalized Ψ can be expressed in terms of the eigenstates

$$\Psi = \sum_n c_n \psi_n, \quad \text{where} \quad H\psi_n = E_n \psi_n \quad (3.23)$$

for all $n = 1, 2, \dots$. The expectation value for the energy can be calculated as

$$\begin{aligned} \langle \Psi | H | \Psi \rangle &= \left\langle \sum_n c_n \psi_n \left| H \right| \sum_{n'} c_{n'} \psi_{n'} \right\rangle \\ &= \sum_n \sum_{n'} c_n^* c_{n'} \langle \psi_n | H | \psi_{n'} \rangle \\ &= \sum_n \sum_{n'} c_n^* E_n c_{n'} \langle \psi_n | \psi_{n'} \rangle \end{aligned}$$

Here we assume that the eigenfunctions have been orthonormalized and we can utilize $\langle \psi_m | \psi_n \rangle = \delta_{mn}$, resulting in

$$\sum_n c_n^* c_n E_n = \sum_n |c_n|^2 E_n.$$

We have already stated that Ψ is normalized, thus $\sum_n |c_n|^2 = 1$, and the expectation value conveniently is bound to follow equation 3.22. The quest to understand the variational principle can be summarized in a sentence - it is possible to tweak the wavefunction parameters to minimize the energy, or summed up in a mathematical phrase,

$$E_0 = \min_{\Psi \rightarrow \Psi_0} \langle \Psi | H | \Psi \rangle. \quad (3.24)$$

3.6 The density functional theory

Hitherto we have tried to solve the Schrödinger equation to get a ground state wave function, and from there we can obtain ground state properties, such as the ground state total energy. One fundamental problem that exists when trying to solve the many-electron Schrödinger equation is that the wavefunction is a complicated function that depends on $3N_e$ variables¹.

Hohenberg and Kohn [47] showed in 1964 that the ground-state density $n_0(r) = |\Psi_0(r)|$ determines a general external potential, which includes U_{en} , up to an additive constant, and thus also the Hamiltonian [48]. From another

¹not including spin

point of view, the theory states that all physical ground-state properties of the many-electron system are unique functionals of the density [45]. A consequence of this is that the number of variables is reduced from $3N_e$ to 3, significantly reducing the computational efforts.

However, the scheme is not without limitations, as the density functional theory (DFT) can only be used to find all the ground-state physical properties if the exact functional of the electron density is known. And 56 years after Hohenberg and Kohn published their paper, the exact functional still remains unknown.

We will start this chapter with a discussion of the Hohenberg-Kohn theorems, before we delve further into the Kohn-Sham equation.

3.6.1 The Hohenberg-Kohn theorems

THEOREM 1. *For any system of interacting particles in an external potential V_{ext} , the density is uniquely determined.*

PROOF. Assume that two external potentials $V_{\text{ext}}^{(1)}$ and $V_{\text{ext}}^{(2)}$, that differ by more than a constant, have the same ground state density $n_0(\mathbf{r})$. The two different potentials correspond to distinct Hamiltonians $\hat{H}_{\text{ext}}^{(1)}$ and $\hat{H}_{\text{ext}}^{(2)}$, which again give rise to distinct wavefunctions $\Psi_{\text{ext}}^{(1)}$ and $\Psi_{\text{ext}}^{(2)}$. Utilizing the variational principle, we find that no wavefunction can give an energy that is less than the energy of $\Psi_{\text{ext}}^{(1)}$ for $\hat{H}_{\text{ext}}^{(1)}$, that is

$$E^{(1)} = \langle \Psi^{(1)} | \hat{H}^{(1)} | \Psi^{(1)} \rangle < \langle \Psi^{(2)} | \hat{H}^{(1)} | \Psi^{(2)} \rangle \quad (3.25)$$

and

$$E^{(2)} = \langle \Psi^{(2)} | \hat{H}^{(2)} | \Psi^{(2)} \rangle < \langle \Psi^{(1)} | \hat{H}^{(2)} | \Psi^{(1)} \rangle. \quad (3.26)$$

Assuming that the ground state is not degenerate, the inequality strictly holds. Since we have identical ground state densities for the two Hamilto-

nian's, we can rewrite the expectation value for equation 3.25 as

$$\begin{aligned}
E^{(1)} &= \langle \Psi^{(1)} | \hat{H}^{(1)} | \Psi^{(1)} \rangle \\
&= \langle \Psi^{(1)} | T + U_{ee} + U_{\text{ext}}^{(1)} | \Psi^{(1)} \rangle \\
&= \langle \Psi^{(1)} | T + U_{ee} | \Psi^{(1)} \rangle + \int \Psi^{*(1)}(\mathbf{r}) V_{\text{ext}}^{(1)}(\mathbf{r}) \Psi^{(1)}(\mathbf{r}) d\mathbf{r} \\
&= \langle \Psi^{(1)} | T + U_{ee} | \Psi^{(1)} \rangle + \int V_{\text{ext}}^{(1)}(\mathbf{r}) n(\mathbf{r}) d\mathbf{r} \\
&< \langle \Psi^{(2)} | \hat{H}^{(1)} | \Psi^{(2)} \rangle \\
&= \langle \Psi^{(2)} | T + U_{ee} + U_{\text{ext}}^{(1)} + \overbrace{U_{\text{ext}}^{(2)} - U_{\text{ext}}^{(2)}}^0 | \Psi^{(2)} \rangle \\
&= \langle \Psi^{(2)} | T + U_{ee} + U_{\text{ext}}^{(2)} | \Psi^{(1)} \rangle + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r} \\
&= E^{(2)} + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r}.
\end{aligned}$$

Thus,

$$E^{(1)} = E^{(2)} + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r} \quad (3.27)$$

A similar procedure can be performed for $E^{(2)}$ in equation 3.26, resulting in

$$E^{(2)} = E^{(1)} + \int (V_{\text{ext}}^{(2)} - V_{\text{ext}}^{(1)}) n(\mathbf{r}) d\mathbf{r}. \quad (3.28)$$

If we add these two equations together, we get

$$\begin{aligned}
E^{(1)} + E^{(2)} &< E^{(2)} + E^{(1)} + \int (V_{\text{ext}}^{(1)} - V_{\text{ext}}^{(2)}) n(\mathbf{r}) d\mathbf{r} \\
&\quad + \int (V_{\text{ext}}^{(2)} - V_{\text{ext}}^{(1)}) n(\mathbf{r}) d\mathbf{r} \\
E^{(1)} + E^{(2)} &< E^{(2)} + E^{(1)},
\end{aligned} \quad (3.29)$$

which is a contradiction. Thus, the two external potentials cannot have the same ground-state density, and $V_{\text{ext}}(\mathbf{r})$ is determined uniquely (except for a constant) by $n(\mathbf{r})$. \square

THEOREM 2. *There exists a variational principle for the energy density functional such that, if n is not the electron density of the ground state, then $E[n_0] < E[n]$.*

PROOF. Since the external potential is uniquely determined by the density and since the potential in turn uniquely determines the ground state wavefunction (except in degenerate situations), all the other observables of the system are uniquely determined. Then the energy can be expressed as a functional of the density.

$$E[n] = \overbrace{T[n] + U_{ee}[n]}^{F[n]} + \overbrace{\int V_{en} n(r) dr}_{U_{en}[n]} \quad (3.30)$$

where $F[n]$ is a universal functional because the treatment of the kinetic and internal potential energies are the same for all systems, however, it is most commonly known as the Hohenberg-Kohn functional.

In the ground state, the energy is defined by the unique ground-state density $n_o(r)$,

$$E_o = E[n_o] = \langle \Psi_o | H | \Psi_o \rangle. \quad (3.31)$$

From the variational principle, a different density $n(r)$ will give a higher energy

$$E_o = E[n_o] = \langle \Psi_o | H | \Psi_o \rangle < \langle \Psi | H | \Psi \rangle = E[n] \quad (3.32)$$

Thus, the total energy is minimized for n_o , and so has to be the ground-state energy. \square

3.6.2 The Kohn-Sham equation

So far, we have tried to make the challenging Schrödinger equation less challenging by simplifying it, with the last attempt containing the Hohenberg-Kohn's theorems where the theory states that the total ground-state energy can, in principle, be determined exactly once we have found the ground-state density.

In 1965, Kohn and Sham [49] reformulated the Hohenberg-Kohn theorems by generating the exact ground-state density $n_o(r)$ using a Hartree-like total wavefunction

$$\Psi(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{N_e}) = \psi_1^{KS}(\mathbf{r}_1) \psi_2^{KS}(\mathbf{r}_2) \dots \psi_{N_e}^{KS}(\mathbf{r}_{N_e}), \quad (3.33)$$

where $\psi_j^{KS}(\mathbf{r}_j)$ are some auxiliary independent single-particle wavefunctions. However, the Kohn-Sham wavefunctions cannot be the correct single-particle wavefunctions since our ansatz implies an exact density

$$n(\mathbf{r}) = \sum_{j=1}^{N_e} |\psi_j^{KS}(\mathbf{r})|^2. \quad (3.34)$$

Recalling that equation 3.30 describes the total energy as a functional of the density,

$$E[n] = T[n] + U_{ee}[n] + U_{en}[n], \quad (3.35)$$

we try to modify it to include the kinetic energy $T_s[n]$ and the interaction energy $U_s[n]$ of the auxiliary wavefunction, and the denotation s for single-particle wavefunctions.

$$\begin{aligned} E[n] &= T[n] + U_{ee}[n] + U_{en}[n] + (T_s[n] - T_s[n]) + (U_s[n] - U_s[n]) \\ &= T_s[n] + U_s[n] + U_{en}[n] + \underbrace{(T[n] - T_s[n]) + (U_{ee}[n] - U_s[n])}_{E_{xc}[n]} \end{aligned}$$

Here we have our first encounter with the *exchange-correlation energy*

$$E_{xc}[n] = \Delta T + \Delta U = (T[n] - T_s[n]) + (U_{ee}[n] - U_s[n]), \quad (3.36)$$

which contains the complex many-electron interaction. For non-interacting system, $E_{xc}[n]$ is conveniently zero, but in interacting systems it most likely is a complex expression. However, one can consider it as our mission to find good approximations to this term, as the better approximations, the closer we get to the exact expression.

The exact total energy functional can now be expressed as

$$\begin{aligned} E[n] &= \underbrace{\sum_j \int \psi_j^{KS*} \frac{-\hbar^2 \nabla^2}{2m} \psi_j^{KS} d\mathbf{r}}_{T_s[n]} + \underbrace{\frac{1}{2} \iint q^2 \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d\mathbf{r}d\mathbf{r}'}_{U_s[n]} \\ &\quad + \underbrace{\int V_{en}(\mathbf{r})n(\mathbf{r})d\mathbf{r}}_{U_{en}[n]} + \underbrace{(T[n] - T_s[n]) + (U_{ee}[n] - U_s[n])}_{E_{xc}[n]}. \end{aligned} \quad (3.37)$$

given that the exchange-correlation functional is described correctly. By utilizing the variational principle, we can now formulate a set of Kohn-Sham single-electron equations,

$$\left\{ -\frac{\hbar^2}{2m_e} \nabla_s^2 + V_H(\mathbf{r}) + V_{j\alpha}(\mathbf{r}) + V_{xc}(\mathbf{r}) \right\} \psi_s^{KS}(\mathbf{r}) = \epsilon_s^{KS} \psi_s^{KS}(\mathbf{r}) \quad (3.38)$$

where $V_{xc}(\mathbf{r}) = \partial E_{xc}[n] / \partial n(\mathbf{r})$ and $V_H(\mathbf{r}) = \int q^2 \frac{n(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d\mathbf{r}'$ is the Hartree potential describing the electron-electron interaction. It is worth to notice that $V_H(\mathbf{r})$ allows an electron to interacts with itself, resulting in a self-interaction contribution, however this will be taken care of in V_{xc} .

Finally, we can define the total energy of the system according to Kohn-Sham theory as

$$E[n] = \sum_j \epsilon_j^{\text{KS}} - \frac{1}{2} \iint q^2 \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}d\mathbf{r}' + E_{\text{xc}}[n] - \int V_{\text{xc}}(\mathbf{r})n(\mathbf{r})d\mathbf{r}. \quad (3.39)$$

If V_{xc} is exact, and $E[n]$ gives the true total energy, we still do not know if the energy eigenvalues ϵ_s^{KS} are the true single-electron eigenvalues. However, there exists one exception, which is that the highest occupied eigenvalue of a finite system has to be exact if the density is exact.

The only task that is left for us now is to find the exact expression for $E_{\text{xc}}[n]$ as a functional of the density $n(\mathbf{r})$. With that expression, we would be able to calculate the total energies of any material, and most likely solve a few of the biggest puzzles in the history of humankind. Unfortunately, the exchange-correlation potential is unknown for most systems.

3.6.3 The exchange-correlation energy

There is one scenario for which we can derive the exact expression of the exchange-correlation functional, namely the *homogeneous electron gas* (HEG). However, this has a natural cause, since by definition $n(\mathbf{r})$ is constant for this situation. Given that it is the variations of electron density that are the foundation of material properties, the usefulness of HEG is limited. The *local density approximation* (LDA) is an approximation based on this approach, where the local density is the only variable used to define the exchange-correlation functional. Specifically, we can set the exchange-correlation potential at each position to be the known exchange-correlation potential from homogeneous electron gas at the electron density observed at that position [49]:

$$V_{\text{xc}}(\mathbf{r}) = V_{\text{xc}}^{\text{electron gas}} [n(\mathbf{r})]. \quad (3.40)$$

This is the simplest and most known approximation to the exchange-correlation functional, and accordingly it has a few drawbacks. One of them is the incomplete cancellation of the self-interaction term, which leads to a repulsion that may cause artificial repulsion between electrons, and hence increased electron delocalization [50]. In addition, LDA has proven challenging to use when studying atoms and molecules because of their rapidly varying electron densities, however, the LDA is seen as succesful for bulk materials because of the slowly varying electron density [51]. Considering the relatively low computational cost and high accuracy, the LDA overall makes a good model for estimation of the exchange-correlation functional for bulk-materials.

In the light of the merits of the LDA, an extensive search for new approximations was launched. The *generalized gradient approximation* (GGA) is an

extension of the LDA, which includes the gradient of the density

$$V_{xc}^{GGA}(\mathbf{r}) = V_{xc} [n(\mathbf{r}), \nabla n(\mathbf{r})]. \quad (3.41)$$

The GGA is a good approximation for the cases where the electron density varies slowly, but faces difficulties in many materials with rapidly varying gradients in the density, causing the GGA to fail. Thus, the annotation *generalized* in GGA is set to include the different approaches to deal with this challenge. Two of the most commonly implemented GGA functionals are the non-empirical approaches Perdew-Wang 91 (PW91) [52] and Perdew-Burke-Ernzerhof (PBE) [53].

Both LDA and GGA are commonly known to severely underestimate the band gaps of semiconductor materials, in addition to incorrectly predicting charge localizations originating from narrow bands or associated with local lattice distortions around defects [54]. The latter limitation is thought to be due to self-interaction in the Hartree potential in equation 3.38.

Hybrid functionals intermix exact Hartree-Fock exchange with exchange and correlation from functionals based on the LDA or GGA. Hartree-Fock theory completely ignore correlation effects, but account for self-interaction and treats exchange as exact. Since LDA/GGA and Hartree-Fock supplement each other, they can be used as a combination for hybrid-functionals resulting in some cancellation of the self-interaction error. Becke [55] introduced a 50% Hartree-Fock exact exchange and 50% LDA energy functional, while Perdew *et al.* [56] altered it to 25% – 75% and favoring PBE-GGA instead of LDA.

The inclusion of Hartree Fock exchange improves the description of localized states, but requires significantly more computational power for large systems. Another method called the GW approximation includes screening of the exchange interaction [57], but has a computational price that does not necessarily defend its use. Thus, the real challenge is to reduce the computational effort while still producing satisfactory results. Heyd *et al.* [58] suggested to separate the non-local Hartree-Fock exchange into a short- and long-range portion, incorporating the exact exchange in the short-range contribution. The separation is controlled by an adjustable parameter ω , which was empirically optimised for molecules to $\omega = 0.15$ and solids to $\omega = 0.11$ and are known as the HSE03 and HSE06 (Heyd-Scuseria-Ernzerhof), respectively [59]. The functionals are expressed as

$$E_{xc}^{HSE} = \alpha E_x^{HFSR}(\omega) + (1 - \alpha) E_x^{PBE,SR}(\omega) + E_x^{PBE,LR}(\omega) + E_c^{PBE} \quad (3.42)$$

where $\alpha = 1/4$ is the Hartree-Fock mixing constant and the abbreviations SR and LR stands for short range and long range, respectively.

Hence, hybrid-functionals are *semi-empirical* functionals that rely on experimental data for accurate results. They give accurate results for several

properties, such as energetics, bandgaps and lattice parameters, and can fine-tune parameters fitted to experimental data for even higher accuracy.

Furthermore, the computational effort required for the hybrid-functionals are significantly larger than for non-empirical functionals such as LDA or GGA. Krukau *et al.* [59] reported a substantial increase in computational cost when reducing the parameter ω from 0.20 to 0.11 for 25 solids, and going lower than 0.11 demanded too much to actually defend its use.

Write about TBMBJ functional and OptB88vDW functional (used by JARVIS).

3.6.4 Self-consistent field methods

So, the remaining question is, how do we solve the Kohn-Sham equation? First, we would need to define the Hartree potential, which can be found if we know the electron density. The electron density can be found from the single-electron wave-functions, however, these can only be found from solving the Kohn-Sham equation. This *circle of life* has to start somewhere, but where?

The process can be defined as an iterative method, *a computational scheme*, as visualized in figure 3.1.

3.6.5 Limitations of the DFT

If we had known the exact exchange-correlation functional, the density functional theory would yield the exact total energy. Alas, that is not the case and we are bound to use approximations in forms of functionals. The accuracy of calculations is dependent on which functional being used, and normally a higher accuracy means the use of a more complex and computationally demanding computational functional.

Nonetheless, density functional theory is considered a very successful approach and Walter Kohn was awarded the Nobel Price in chemistry in 1998 for his development of the density-functional theory [60]. One can only hope that the future will be as bright as the past, and that this successful theory provides incentives for further growth in the next generation.

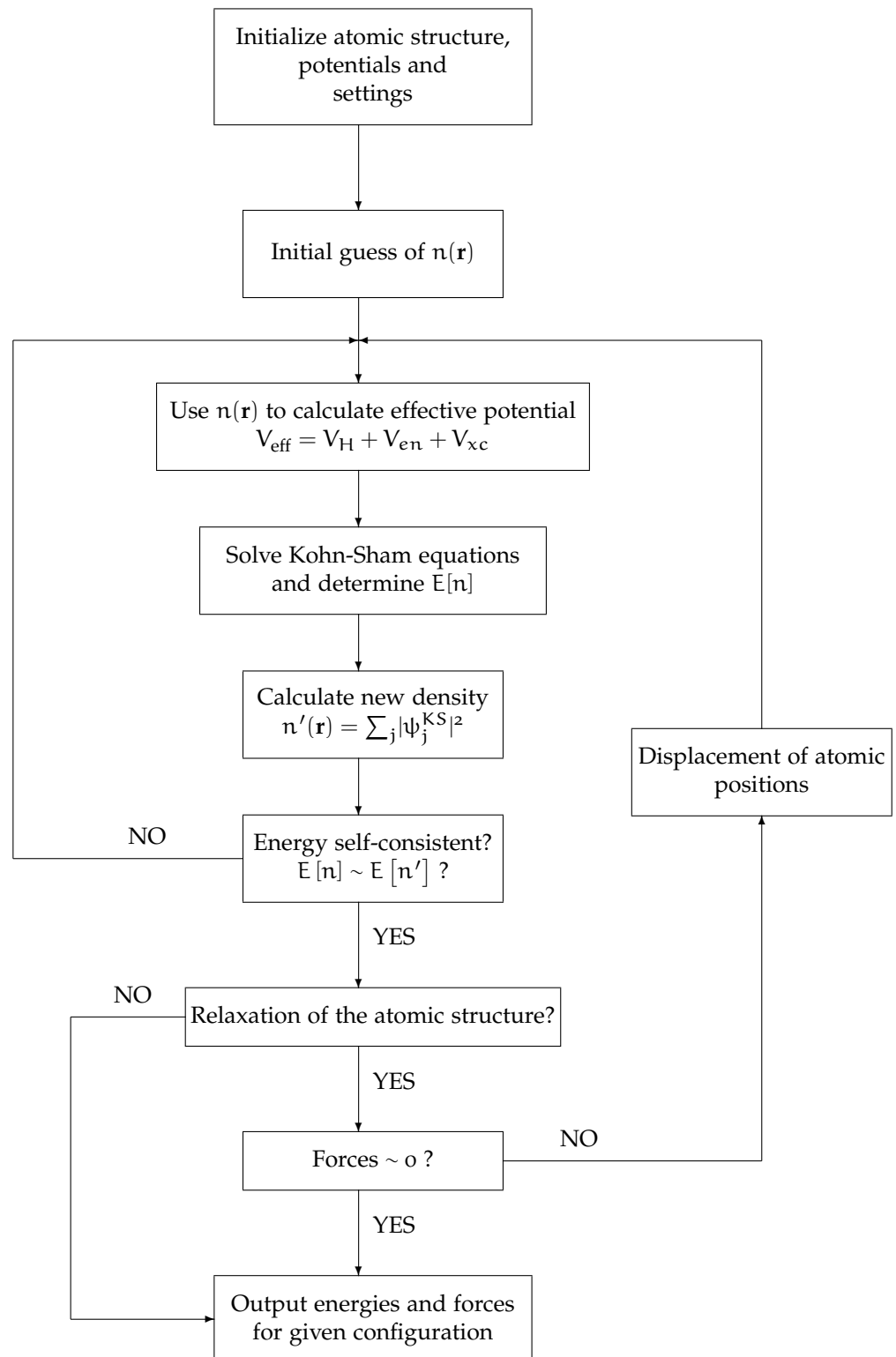


Figure 3.1: A flow chart of the self-consistent field method for DFT.

Chapter 4

Machine learning

The enormous amount of data generated in the digital world today is beyond comprehension. In 2019, more than 500 hours of video was uploaded to Youtube every second, totalling to over 82 years of content every day ¹. In addition, more than 1 trillion web pages exists and has information that needs to be stored somewhere.

However, an increasing amount of data comes hand in hand with an increasing demand of knowledge about the data. If we are unable to extract information from the data, the data serves no intention and exists as an excess. Therefore, we need methods to process and automate data analysis, which is what the promises of *machine learning* covers. Machine learning can reveal patterns in data with ease where a human would face difficulties, and use this information to predict or generate new data. Many tools in machine learning is based on probability theory, which can be applied to problems involving uncertainty. Thus, machine learning is also commonly (and boringly) named as *statistical learning*.

There are mainly two types of machine learning, either *supervised* or *unsupervised* learning. In unsupervised learning we are given inputs $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$, and it is necessary to use the tools of machine learning to find peculiar and interesting patterns. In the supervised approach, on the other hand, the model tries to learn a mapping from inputs \mathbf{x} to outputs y , given a labeled set of pairs $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$. The set \mathcal{D} is known as the training set, and N is the number of entries. Each training input \mathbf{x}_i has D -dimensions that describes each entry, where each dimension is known as a *feature*. The features could be exemplified as height or weight, or it could be something complex that has no practical meaning. The flexibility of the shape of a feature is also shared with the output. It can in principle be anything, but it is mostly assumed that the output is either *categorical* or *nominal* restricted by a finite set $y_i \in \{1, \dots, \mathcal{C}\}$. The problem is defined as *classification* if the output is categorical, or *regression*

¹Source: <https://www.youtube.com/intl/no/about/press/> extracted 15.02.2021

if the output is real-valued [61].

4.1 Supervised learning

Recall, supervised learning in classification has as goal to learn the target output $y \in \{1, \dots, \mathcal{C}\}$ from the inputs \mathbf{x} . The number of classes is \mathcal{C} , and depicts if the classification is *binary* ($\mathcal{C} = 2$), *multiclass* ($\mathcal{C} > 2$), or *multi-label* if the class labels are not mutually exclusive (exemplified with the weather can be both sunny and cold at the same time). Normally, classification is used when the problem is formulated as a multiclass classification, and hereon we will adapt to the formulation as well [61].

In order to be able to learn from data, we will need to formulate a function approximate. By assuming $y = f(\mathbf{x}) + \epsilon$ for some unknown function f and a random error term ϵ with mean zero, we can try to approximate f from a labeled training set, which we can use to make the predictions $\hat{y} = \hat{f}(\mathbf{x})$. With the estimated \hat{f} , we can make predictions on unlabeled data and achieve a *generalized model*. The estimated function \hat{f} is often considered as a black box, since we are not necessary interested in the exact shape of the function but rather the predictions.

As simple the idea behind supervised classification appears, a generalized model remains deeply dependent on the available data. Imagine a training set containing two entries. The first one is a young and tall person who is labeled as healthy. The other entry is an old and short person who is labeled sick. The pattern in this simple scenario is abundantly clear, but will face a challenge if it were to predict on a test set containing a person who is young and short.

Therefore, it is desirable to compute the probability of an entry belonging to one class. The probability distribution is given by $p(y|\mathbf{x}, \mathcal{D})$, where the probability is conditional on the input vector (test set) \mathbf{x} and the training set \mathcal{D} . If the output is probabilistic, we can compute the estimation to the true label as

$$\hat{y} = \hat{f}(\mathbf{x}) = \underset{c=1}{\operatorname{argmax}} f(\mathbf{x})p(y = c|\mathbf{x}, \mathcal{D}), \quad (4.1)$$

which represents the most probable class label and is known as the *maximum a posteriori* estimate.

4.2 Evaluating accuracy of a model

The idea of finding one algorithm that is far more superior than any other algorithm, for all types and sizes of datasets, is of an imaginary sort. There is

no algorithm that has this property, since one model might be recognized as best on one particular dataset, while others are far better on other datasets. The same goes with evaluating the model - there is no metrics that stand alone as the best metric to evaluate a model. Choosing how to actually evaluate a model can be the most challenging part of a statistical learning procedure.

4.2.1 Bias-variance tradeoff

To illustrate a challenge in choosing the correct parameters, we give an example using the mean squared error (MSE) as a *cost function*, that is a measurement of how wrong the model is in term of estimating the relationship between \mathbf{X} and \mathbf{y} . Additionally, assume that our data can be represented by an unknown function with the addition of some noise ϵ ,

$$\mathbf{y} = f(\mathbf{x}) + \epsilon,$$

where ϵ is normally distributed with a mean equal to zero and standard deviation equal to σ^2 . Furthermore, we also assume that the function $f(\mathbf{x})$ can be approximated to a model $\hat{\mathbf{y}}$, where the model is defined by a design matrix \mathbf{X} and parameters β .

$$\hat{\mathbf{y}} = \mathbf{X}\beta$$

The parameters β are in turn found by optimizing the mean squared error (MSE) via the cost function

$$C(\mathbf{X}, \beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2 = E[(\mathbf{y} - \hat{\mathbf{y}})^2].$$

One can show that the cost function can be rewritten as

$$E[(\mathbf{y} - \hat{\mathbf{y}})^2] = \frac{1}{n} \sum_i (f_i - E[\hat{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\hat{y}_i - E[\hat{\mathbf{y}}])^2 + \sigma^2.$$

and since the variance of \mathbf{y} and ϵ are both equal to σ^2 , we can use the relations $E[\mathbf{y}] = \mathbf{f}$, $E[\epsilon] = 0$ and $\text{Var}(\mathbf{y}) = \text{Var}(\epsilon) = \sigma^2$. The mean value of ϵ is by definition equal to zero. In addition, the function \mathbf{f} is not a stochastic variable, and the same argument can also be used for $\hat{\mathbf{y}}$. This allows us to insert the expression \mathbf{y} into the cost function

$$E[(\mathbf{y} - \hat{\mathbf{y}})^2] = E[(\mathbf{f} + \epsilon - \hat{\mathbf{y}})^2].$$

By using the infamous trick of both adding and subtracting simultaneously, we arrive at

$$\mathbb{E}[(\mathbf{y} - \hat{\mathbf{y}})^2] = \mathbb{E}[(\mathbf{f} + \boldsymbol{\epsilon} - \hat{\mathbf{y}} + \mathbb{E}[\mathbf{y}] - \mathbb{E}[\mathbf{y}])^2].$$

And simply by using the relations mentioned above concerning the expectation value for \mathbf{y} and the variances for both \mathbf{y} and $\boldsymbol{\epsilon}$, the cost function can be rewritten to:

$$\begin{aligned} \mathbb{E}[(\mathbf{y} - \hat{\mathbf{y}})^2] &= \mathbb{E}[(\mathbf{f} - \mathbb{E}[\hat{\mathbf{y}}])^2] + \text{Var}(\hat{\mathbf{y}}) + \sigma_{\epsilon}^2 \\ &= \frac{1}{n} \sum_i (f_i - \mathbb{E}[\hat{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\hat{\mathbf{y}}_i - \mathbb{E}[\hat{\mathbf{y}}])^2 + \sigma_{\epsilon}^2. \end{aligned}$$

The first term on the right hand side is the squared bias, the amount by which the average of our estimate differs from the true mean, while the second term represents the variance of the chosen model. The last term is the variance of the error $\boldsymbol{\epsilon}$, also known as the irreducible error. In general, an estimated function \hat{f} will never be a perfect estimate for f since we can not reduce the error associated introduced by $\boldsymbol{\epsilon}$. Therefore, any model will always be restricted to an upper bound of accuracy due to the irreducible error.

The more complex model one has, the lower the bias becomes, and the higher the variance becomes, as seen in figure 4.1.. If the model is not complex enough due to high bias and low variance, the algorithm can end up not learning the relevant relations between features and output, and is known as *underfitting*. On the other hand, for a very complex model with low bias and high variance, the algorithm might find trends in random noise from the training data instead of the relevant features, resulting in *overfitting*. An ideal model would be one that simultaneously achieves low variance and low bias. Therefore, we have to do a trade-off between how much bias and variance we would like in the model.

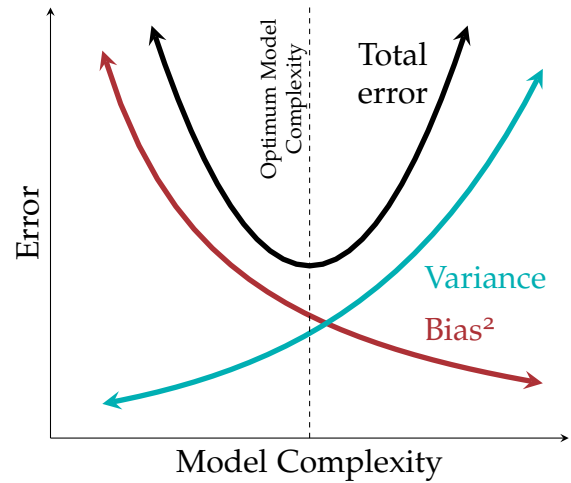


Figure 4.1: A schematic representation of the bias-variance tradeoff. To the left of the optimum model complexity increases the chance of having an underfitted model, while to the right increases the chance of having an overfitted model.

4.2.2 Accuracy, precision and recall

Given a model that has dealt with the intricacy of increasing complexity, we would like to evaluate the model's output quality. For a binary supervised classification problem we can measure the accuracy by finding how many correct predictions has been made. Prediction accuracy can provide a fine initial analysis, but it has some significant drawbacks seen in unbalanced datasets. This can be easily explained with a dataset consisting of 99 : 1 ratio of class, since just guessing the majority class will result in a very high 99% accuracy. Perhaps it is the 1% that is the most important class, thus the accuracy score severely lacks information for the model.

Therefore, we turn to other evaluation metrics such as a *confusion matrix*. A confusion matrix is a method for measuring the performance of classifiers. It is set up as a table with 4 different categories, where two of the categories are the predicted outcomes of the model/classifier and the two final categories are the true outcomes, which one wants to fit. An example of a confusion matrix for a binary classifier predicting the presence of a pregnancy is shown in table 4.1.

		Prediction outcome		
		p	n	total
Actual value	p'	True Positive	False Negative	P'
	n'	False Positive	True Negative	N'
total		P	N	

Table 4.1: A confusion matrix for a binary classifier.

For the binary confusion matrix there are two possible predicted outcomes, either positive or negative. This gives rise to some terminology.

- True Positive (TP): The classifier correctly predicts a positive event.
- True Negative (TN): The classifier correctly predicts a negative event.
- False Positive (FP): The classifier incorrectly predicts a positive event when the true event was negative.

- **False Negative (FN):** The classifier incorrectly predicts a negative event when the true event was positive.

From the confusion matrix one can then start estimating the performance of the model, by calculating different factors, such as

- **Recall** is a measure of how many truly relevant results are predicted. A high recall relates to a low false negative rate, and is defined as

$$\text{Recall} = \frac{TP}{TP + FN}. \quad (4.2)$$

- **Precision** is a measure of result relevancy. A high precision relates to a low false positive rate, and is defined as

$$\text{Precision} = \frac{TP}{TP + FP}. \quad (4.3)$$

Similar to the bias-variance tradeoff, it is common to compare the recall with the precision to identify the tradeoff for different thresholds. High scores for both reveals that a classifier returns accurate results combined with returning a majority of all positive results. Therefore, an ideal classifier will return many correctly predicted results.

Sometimes a classifier can have drastically different values for the precision and recall. This leads to another estimator for the performance of a classifier, which is known as the F1-score. The F1-score is defined as the harmonic mean of precision and recall,

$$F1\text{-score} = \frac{2 \cdot \text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}},$$

and can be used to find a good tradeoff between recall and precision. The highest value of F1-score is 1 and is considered an ideal classifier, while the lowest is 0.

4.2.3 Cross-validation

When evaluating different parameters for models, commonly done in a grid-search scheme, there is an abundant risk of performing an overfit to the test set since we can tweak the parameters to a model so it can perform optimally. To solve this problem, we can exclude a part of the dataset as a validation set (in addition to a test set). Therefore, we can train a model on the training set, and evaluate the parameters on the validation set. After a lot of trial and

error and the experiment seems successful, we can do one final evaluation on the test set.

Unfortunately, this reduces the number of samples that can be used for training drastically. A fix for this is to apply a cross-validation. Cross-validation is a technique used to evaluate predictive models by partitioning the original sample into a training set to train the model, and a test set to evaluate it.

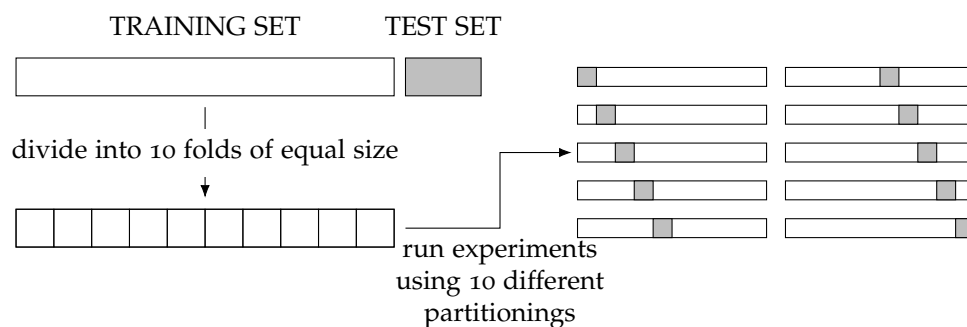


Figure 4.2: A schematic representation of a 10-fold cross-validation scheme.

It is common to apply cross validation into fold, yielding the name of k-fold cross-validation. In k-fold cross-validation, the sample is partitioned into k equal sized subsamples, as visualized in figure 4.2. Of the k samples, a single sample is used as validation set while the remaining k-1 samples are used as training data. The process is then repeated k-times, such that each of the k-th subsample is used as validation set exactly one time. Therefore, all observations are used for both training and validation, and each observation is used for validation exactly once. The k results from the folds can then be averaged to produce an estimate. These train and test samples are allowed to have an imbalanced dataset, such as each class is not necessarily represented equally in each fold. Since supervised algorithms tend to weight each instance equally, it may result in overrepresented classes getting too much weight. Even worse could be the result of an fold where one class is not represented at all, resulting in a model that can't predict a class.

To deal with vulnerability, one can employ a stratified k-fold cross-validation. Stratification is a process that seeks to ensure that each fold is representative of all strata in the data, making each fold having approximately equal class-representation.

4.3 Decision trees

Classification and regression trees (CART), also called decision trees, is one of the more basic supervised algorithms, and can be used for both regres-

sion and classification tasks, as the name suggests. The strength of decision trees lays within the simplicity that allows to build more complex networks. We will in this section provide special emphasis to classification trees, but with some remarks to the regression trees to provide a brief perspective of distinctions.

The idea behind decision trees is to find the features that contain the most information regarding the target feature, and then split up the dataset along the values of these features. This feature selection enables the target feature values for the resulting underlying dataset to be as *pure* as possible. The features that can reproduce the best target features are normally said to be the most informative features.

A decision tree can be divided a *root node*, *interior nodes*, and the final *leaf nodes*. Each entity are connected by *branches*. The decision tree is able to learn an underlying structure of the training data and can, given some assumptions, make predictions on unseen query instances. It is the leaf nodes that accomodates the predictions we will make for new entries that is presented to our trained model.

The process behind a decision tree can be seen as a top-down approach. First, we make a leaf provide the classification of a given instance. Then, a node specifies a test of some attribute of the instance, while a branch corresponds to a possible value of an attribute. This allows an entry to be classified by starting at the root node of the tree with corresponding testing of the attribute specified by this node. Subsequently, the instance move down the tree branch corresponding to the value of the attribute. Then the steps can be repeated for a new subtree rooted at the new node.

A classification tree mainly differs to a regression tree by the response of the prediction, since it produces a qualitative response rather than a quantitative one. For a regression tree, the response is given by the mean response of the training observations that belong to the same terminal node. For a classification tree, on the contrary, the response is given by the most commonly occuring class of training observations in the region which it belongs. Thus, the interpretation process includes both the class prediction corresponding to a particular terminal node region, but also in the class proportions among the training observations that fall into that region.

4.3.1 Growing a classification tree

In growing a classification tree, a process called recursive binary splitting is applied which involves mainly two steps.

1. Split the set of possible values (x_1, x_2, \dots, x_p) into J distinct an non-overlapping regions R_1, R_2, \dots, R_J .

2. If an observation falls within the region R_j , we make the prediction given by the most commonly occurring class of training observations in R_j .

The computational aspect of recursively doing this for every possible combination of features does not defend its use, and therefore the common strategy is to use a top-down approach. Thus, it begins at the top of the tree and consecutively splits the predictor space. This is indicated by two new branches further down the tree. It should be noted that the top-down approach is a greedy approach, since the best split is made at each step of the tree-growing process, instead of trying to pick a split that will lead to a better tree in a future step.

We can define a *probability density function* (PDF) p_{mk} that represents the number of observations k in a region R_m with N_m observations. This likelihood function can be represented in terms of the proportion $I(y_i = k)$ of observations of this class in region R_m as

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k) \quad (4.4)$$

Hitherto, the splitting of the nodes have been decided by the misclassification error

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i \neq k) = 1 - p_{mk}. \quad (4.5)$$

where the *indicator* function I equals one if we misclassify and equals zero if we classify correctly. However, other methods exists such as the Gini index

$$g = \sum_{k=1}^K p_{mk}(1 - p_{mk}) \quad (4.6)$$

and the information entropy

$$s = - \sum_{k=1}^K p_{mk} \log p_{mk}. \quad (4.7)$$

The two latter approaches are more sensitive to node purity than the misclassification error, i.e. only containing one class, and are in general preferred [61].

4.3.2 Classification algorithm

The CART algorithm splits the data set in two subsets using a single feature k and a threshold t_k . The pair of quantities (k, t_k) that constitute the purest subset using the gini factor G results in the cost function

$$C(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}, \quad (4.8)$$

where $G_{\text{left/right}}$ measures the impurity of left or right subset and $m_{\text{left/right}}$ is the number of instances on either the left or the right subset. The algorithm tries to minimize the cost function to find the pair by splitting the training set in two, and then following the same logic for the next subsets. It will continue to do this recursively until it reaches the maximum depth hyperparameter, or if the next split does not reduce impurity.

4.3.3 Pruning a tree

A decision tree has the ability to turn into a very complex model, making it significantly prone to overfitting. Therefore, techniques that deals with this vulnerability must be implemented.

Pre-puning is a method that stops the growing of a tree if the decrease in error is not sufficient to justify an increasing complex model by adding an extra subtree. However, this implementation has the liability for features that have small predictive power as this might cause a model without any splits at all.

Post-pruning, or just pruning, is the standard method which involves growing the tree to full size, and then prune it. To determine how far to prune it, we can use a cross-validated scheme to evaluate the amount of terminal nodes that has the lowest error.

4.3.4 Pros and cons for decision trees

Decision trees have several clear advantages compared to other algorithms. They are easy to understand, also known as a *white box*, and can be visualised effortlessly for small trees. The algorithm is completely invariant to scaling of the data since each feature is processed separately. Additionally, decision trees can handle both continuous and categorical data and can model interactions between different descriptive features.

As auspicious the advantages of decision trees seems, they are inevitable prone to overfitting and hence does not generalize the data well. Even with pre-puning, post-pruning and setting a maximum depth of terminal nodes, the algorithm is still prone for overfitting [62]. Another important issue concerns training on unbalanced datasets where one class occurs more frequently

than other classes, since this will lead to biased trees because the algorithm will favor the more occurring class. Furthermore, small changes in the data may lead to a completely different tree.

Many of these issues can be addressed by using ensemble methods such as either bagging, random forest, or boosting, and can result in a solid improvement of the predictive performance of trees.

4.4 Ensemble methods

By using a single decision tree, we often end up with an overfitted model that possess a high variance. Luckily, we can apply methods that aggregate different machine learning algorithms to reduce variance. If each of the algorithms get slightly different results, as they learn different part of the data, we can combine the results into something that is better than any one algorithm alone. These approaches falls under the category of ensemble methods, and will be elaborated further in this section.

4.4.1 Bagging

Bootstrap aggregation, or just *bagging*, is an ensemble method that involves averaging many estimates. If we have M trained trees on different subsamples of the data, chosen randomly, we can compute the ensemble

$$f(\mathbf{x}) = \sum_{b=1}^M \frac{1}{B} f_b(\mathbf{x}) \quad (4.9)$$

where f_b is the b 'th tree.

Simply re-running the same algorithm on different subsamples can result in a small variance reduction compared to a single tree due to highly correlated predictors, which showcase the need for better approaches.

Random forests provides an improvement of normal bagged trees by choosing a random sample of m predictors as split candidates from the full set of p predictors. The split is restricted in choosing only one of the m predictors. Normally, the value of m is chosen as

$$m \approx \sqrt{p}, \quad (4.10)$$

which means that at each split in a tree, the algorithm is restricted to a very small portion of the available predictors. The specific about the algorithm can

be found in Algorithm 1.

Algorithm 1: Random forest algorithm.

```

for For  $b = 1 : B$  do
    Draw a bootstrap sample from the training data;
    Select a tree  $T_b$  to grow based on the bootstrap data;
    while node size smaller than maximum node size do
        Select  $m \leq p$  variables at random from  $p$  predictors;
        Pick the best split point among the  $m$  features using CART
        algorithm and create a new node;
        Split the node into daughter nodes;
    end
end
Output the ensemble of trees  $\{T_b\}_1^B$  and make predictions

```

By inducing randomness into the model, we arrive at a suprisingly capable model that has a high predictive accuracy [63]. This can be exemplified by supposing that there is one strong predictor in a dataset, together with several other fairly strong predictors. Most of the trees will use this strong predictor at the top split, which means that the bagged trees will look quite similar to each other and will have highly correlated predictions.

However, even with higher prediction accuracy, it comes as a compromise since we loose the easy ability of model interpretation. A single tree can be easy to understand, but interpretation of a huge jungle of trees does not necessarily seem appealing for even an experienced data scientist. Furthermore, a random forest does not substantially reduce the variance as averaging many uncorrelated trees would do, as we will soon find out.

4.4.2 Boosting

Boosting is an ensemble method that fits an additive expansion in a set of elementary basis functions. The basic idea is to combine several weak classifiers, that are only just better than a random guess, in order to create a good classifier. This can be done in an iterative approach where we apply a weak classifier to modify the data. For each iteration, we make sure to weight the observations that are misclassified with a factor. The method is known as adaptive, since the algorithm is able to adapt during the learning process.

In *forward stagewise additive modeling* we want to find an adaptive model

$$f_M(\mathbf{x}) = \sum_{i=1}^M \beta_m b(\mathbf{x}; \gamma_m) \quad (4.11)$$

where β_m are expansion parameters that will be determined in a minimization process, and $b(\mathbf{x}; \gamma_m)$ are functions of the multivariable parameter \mathbf{x} that

are described by the parameters γ_m . We will in this example consider a binary classification problem with the outcomes $\gamma_i \in \{-1, 1\}$ where $i = 0, 1, 2, \dots, n-1$ as the set of observables. The predictions are produced by the classification function $G(\mathbf{x})$.

Then, the error rate of the training sample is given as

$$\overline{\text{err}} = \frac{1}{n} \sum_{i=0}^{n-1} I(\hat{y}_i \neq G(\mathbf{x}_i)). \quad (4.12)$$

After defining a weak classifier, we can apply it iteratively to repeatedly modified versions of the data which produce a sequence of different weak classifiers $G_m(\mathbf{x})$. The function $f_M(\mathbf{x})$ will be expressed in terms of

$$G_M(\mathbf{x}) = \text{sign} \sum_{i=1}^M \alpha_m G_m(\mathbf{x}). \quad (4.13)$$

The iterative procedure can be defined as

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m G_m(\mathbf{x}). \quad (4.14)$$

The cost function that leads to the *discrete AdaBoost* algorithm [64] is the exponential cost function

$$C(\mathbf{y}, \mathbf{f}) = \sum_{i=0}^{n-1} \exp(-\hat{y}_i (f_{m-1}(\mathbf{x}_i) + \beta G(\mathbf{x}_i))), \quad (4.15)$$

or with the weight $w_i^m = \exp(-\hat{y}_i f_{m-1}(\mathbf{x}_i))$ we can rewrite to

$$C(\mathbf{y}, \mathbf{f}) = \sum_{i=0}^{n-1} w_i^m \exp(-\hat{y}_i \beta G(\mathbf{x}_i)). \quad (4.16)$$

We can optimize G for any $\beta > 0$ with

$$G_m(\mathbf{x}) = \text{sign} \sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G(\mathbf{x}_i)). \quad (4.17)$$

This is the classifier that minimize the weighted error rate in predicting y . Furthermore, we can rewrite the cost function to

$$C = \exp(-\beta) \sum_{\hat{y}_i = G(\mathbf{x}_i)} w_i^m + \exp(\beta) \sum_{\hat{y}_i \neq G(\mathbf{x}_i)} w_i^m \quad (4.18)$$

$$= (\exp(\beta) - \exp(-\beta)) \sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G(\mathbf{x}_i)) + \exp(-\beta) \sum_{i=0}^{n-1} w_i^m. \quad (4.19)$$

Substituting G_m into C and solving for β , we obtain

$$\beta_m = \frac{1}{2} \log \frac{1 - \overline{\text{err}}}{\overline{\text{err}}} \quad (4.20)$$

with the error redefined as

$$\overline{\text{err}} = \frac{1}{n} \frac{\sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G_m(\mathbf{x}_i))}{\sum_{i=0}^{n-1} w_i^m}. \quad (4.21)$$

Finally, this leads to an update of

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m G_m(\mathbf{x}) \quad (4.22)$$

and the weights at the next iteration becomes

$$w_i^{m+1} = w_i^m \exp(-\hat{y}_i \beta_m G_m(\mathbf{x}_i)). \quad (4.23)$$

With the above definitions, we can define the discrete Adaboost algorithm in Algorithm 2.

Algorithm 2: Discrete Adaboost algorithm.

Initialize weights $w_i = 1/n$, $i = 0, \dots, n-1$, such that $\sum_{i=0}^{n-1} w_i = 1$;

for $m = 1 : M$ **do**

 Fit the classifier $f_m(\mathbf{x}) \in \{-1, 1\}$ using weights w_i on the training data;

 Compute the error $\overline{\text{err}} = \frac{1}{n} \frac{\sum_{i=0}^{n-1} w_i^m I(\hat{y}_i \neq G_m(\mathbf{x}_i))}{\sum_{i=0}^{n-1} w_i^m}$;

 Define a quantity $\alpha_m = \log [(1 - \overline{\text{err}}_m)/\overline{\text{err}}_m]$;

 Set new weights to $w_i \leftarrow w_i \exp(\alpha_m I(y_i \neq G(\mathbf{x}_i)))$;

end

Compute the new classifier $G(\mathbf{x}) = \sum_{i=0}^{n-1} \alpha_m I(y_i \neq G(\mathbf{x}_i))$;

It is possible to apply different cost functions resulting in a variety of boosting algorithms, which AdaBoost is an example with the cost function being the exponential cost function. But instead of deriving new versions of boosting based on different cost functions, we can find one generic method. The approach is known as *gradient boosting* [65].

Initially, we want to minimize

$$\hat{\mathbf{f}} = \operatorname{argmin}_{\mathbf{f}} L(\mathbf{f}), \quad (4.24)$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))$ are the parameters of the models, and L is a chosen loss function.

This can be solved stagewise, using an approach named *gradient descent*. At step m , let \mathbf{g}_m be the gradient evaluated at $f(\mathbf{x}_i) = f_{m-1}$:

$$\mathbf{g}_m(\mathbf{x}_i) = \left[\frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \right]_{f(\mathbf{x}_i)=f_{m-1}(\mathbf{x}_i)}. \quad (4.25)$$

Then we can update

$$\mathbf{f}_m = \mathbf{f}_{m-1} - \rho_m \mathbf{g}_m, \quad (4.26)$$

where ρ_m is the step length and can be find by approximating the real function

$$h_m(\mathbf{x}) = -\rho \mathbf{g}_m(\mathbf{x}). \quad (4.27)$$

So far, this only optimize f at a fixed set of points, but we can modify it by fitting a weak classifier to approximate the negative gradient. The resulting algorithm is shown in Algorithm 3 as the gradient boost algorithm.

Algorithm 3: Gradient boost algorithm.

```

Initialize the estimate  $f_0(\mathbf{x})$ ;
for  $m = 1 : M$  do
    Compute the negative gradient vector  $\mathbf{u}_m = -\partial C(\mathbf{y}, \mathbf{f}) / \partial \mathbf{f}(\mathbf{x})$  at
         $\mathbf{f}(\mathbf{x}) = \mathbf{f}_{m-1}$ ;
    Fit the base learner to the negative gradient  $h_m(\mathbf{u}_m, \mathbf{x})$ ;
    Update the estimate  $\mathbf{f}_m(\mathbf{x}) = \mathbf{f}_{m-1}(\mathbf{x}) + \nu h_m(\mathbf{u}_m, \mathbf{x})$ ;
end
Output the final estimation  $\mathbf{f}_M(\mathbf{x}) = \sum_{m=1}^M \nu h_m(\mathbf{u}_m, \mathbf{x})$ 

```

4.5 Dimensionality reduction

Supervised learning introduces models that are easy to understand, visualize and has well-defined tools and models. There are several different methods to evaluate a model with many different types of scores such as accuracy, precision and f1-scores. Unfortunately, this does not transfer to unsupervised learning. In unsupervised learning, there is no simple goal for the analysis and the evaluation tends to of a subjective matter. Therefore, unsupervised learning is often used as an *exploratory data analysis* [66]. For data consisting of hundreds or thousands of features, it is possible to apply unsupervised learning to find correlated features and reduce dimensionality of the data, potentially reducing computational effort and time usage drastically. This is the idea of *principal component analysis* (PCA).

4.5.1 Principal component analysis

Principal component analysis is an unsupervised algorithm that tries to find a low-dimension representation of a dataset that contains as much of the variance in the data as possible. Each of the dimensions found by PCA are a linear combination of the features in the dataset, and are known as *principal components*.

We can write the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, with p features and n entries, in terms of its column vectors as

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_0 & \mathbf{x}_1 & \mathbf{x}_2 & \dots & \dots & \mathbf{x}_{p-1} \end{bmatrix}, \quad (4.28)$$

with a given vector

$$\mathbf{x}_i^T = \begin{bmatrix} x_{0,i} & x_{1,i} & x_{2,i} & \dots & \dots & x_{n-1,i} \end{bmatrix}. \quad (4.29)$$

Then we can compute the *covariance matrix* of the design matrix \mathbf{X} , which is a measurement of the joint variability of the p features in \mathbf{X} . The covariance is defined as

$$\text{cov}[\mathbf{v}, \mathbf{u}] = \frac{1}{n} \sum_{i=0}^{n-1} (v_i - \bar{v})(u_i - \bar{u}), \quad (4.30)$$

where \mathbf{v} and \mathbf{u} are two vectors with n elements each. The covariance matrix is defined by applying the covariance for every pairwise feature, resulting in a $p \times p$ matrix. On the diagonal, the covariance of two equal features becomes the variance of one,

$$\text{cov}[\mathbf{u}, \mathbf{u}] = \frac{1}{n} \sum_{i=0}^{n-1} (u_i - \bar{u})(u_i - \bar{u}) \quad (4.31)$$

$$= \frac{1}{n} \sum_{i=0}^{n-1} (u_i - \bar{u})^2 \quad (4.32)$$

$$= \text{var}[\mathbf{u}]. \quad (4.33)$$

The covariance accepts values between 0 and infinity, which rises a computational issue due to loss of numerical precision. Therefore, we scale the covariance matrix using the correlation function

$$\text{corr}[\mathbf{u}, \mathbf{v}] = \frac{\text{cov}[\mathbf{u}, \mathbf{v}]}{\sqrt{\text{var}[\mathbf{u}]\text{var}[\mathbf{v}]}}. \quad (4.34)$$

Since all values are inbetween -1 and 1 , we avoid any loss of numerical precision. The resulting covariance matrix $\mathbf{C} \in \mathbb{R}^{p \times p}$ becomes

$$\mathbf{C}[\mathbf{x}] = \begin{bmatrix} \text{var}[\mathbf{x}_0] & \text{cov}[\mathbf{x}_0, \mathbf{x}_1] & \text{cov}[\mathbf{x}_0, \mathbf{x}_2] & \dots & \dots & \text{cov}[\mathbf{x}_0, \mathbf{x}_{p-1}] \\ \text{cov}[\mathbf{x}_1, \mathbf{x}_0] & \text{var}[\mathbf{x}_1] & \text{cov}[\mathbf{x}_1, \mathbf{x}_2] & \dots & \dots & \text{cov}[\mathbf{x}_1, \mathbf{x}_{p-1}] \\ \text{cov}[\mathbf{x}_2, \mathbf{x}_0] & \text{cov}[\mathbf{x}_2, \mathbf{x}_1] & \text{var}[\mathbf{x}_2] & \dots & \dots & \text{cov}[\mathbf{x}_2, \mathbf{x}_{p-1}] \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \text{cov}[\mathbf{x}_{p-1}, \mathbf{x}_0] & \text{cov}[\mathbf{x}_{p-1}, \mathbf{x}_1] & \text{cov}[\mathbf{x}_{p-1}, \mathbf{x}_2] & \dots & \dots & \text{var}[\mathbf{x}_{p-1}] \end{bmatrix}, \quad (4.35)$$

for all vectors \mathbf{x}_i where $i = 0, 1, \dots, p-1$. The correlation matrix becomes

$$\mathbf{K}[\mathbf{x}] = \begin{bmatrix} 1 & \text{corr}[\mathbf{x}_0, \mathbf{x}_1] & \text{corr}[\mathbf{x}_0, \mathbf{x}_2] & \dots & \dots & \text{corr}[\mathbf{x}_0, \mathbf{x}_{p-1}] \\ \text{corr}[\mathbf{x}_1, \mathbf{x}_0] & 1 & \text{corr}[\mathbf{x}_1, \mathbf{x}_2] & \dots & \dots & \text{corr}[\mathbf{x}_1, \mathbf{x}_{p-1}] \\ \text{corr}[\mathbf{x}_2, \mathbf{x}_0] & \text{corr}[\mathbf{x}_2, \mathbf{x}_1] & 1 & \dots & \dots & \text{corr}[\mathbf{x}_2, \mathbf{x}_{p-1}] \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \text{corr}[\mathbf{x}_{p-1}, \mathbf{x}_0] & \text{corr}[\mathbf{x}_{p-1}, \mathbf{x}_1] & \text{corr}[\mathbf{x}_{p-1}, \mathbf{x}_2] & \dots & \dots & 1 \end{bmatrix}. \quad (4.36)$$

The covariance matrix can be rewritten as a function of the design matrix,

$$\mathbf{C}[\mathbf{x}] = \frac{1}{n} \mathbf{X} \mathbf{X}^T = \mathbb{E}[\mathbf{X} \mathbf{X}^T], \quad (4.37)$$

where $\mathbb{E}[\mathbf{X} \mathbf{X}^T]$ is the expectation value.

Further on, we assume that we can do apply a number of orthogonal transformations by some orthogonal matrices $\mathbf{S} = [\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_{p-1}] \in \mathbb{R}^{p \times p}$ with the column vectors $\mathbf{s}_i \in \mathbb{R}^p$. Additionally, we assume that there is a transformation

$$\mathbf{C}[\mathbf{y}] = \mathbf{S} \mathbf{C}[\mathbf{x}] \mathbf{S}^T = \mathbb{E}[\mathbf{S} \mathbf{X} \mathbf{X}^T \mathbf{S}^T], \quad (4.38)$$

such that the new matrix $\mathbf{C}[\mathbf{y}]$ is diagonal with elements $[\lambda_0, \lambda_1, \lambda_2, \dots, \lambda_{p-1}]$. By multiplying with \mathbf{S}^T , we arrive at the given eigenvalue i of the covariance matrix that

$$\mathbf{S}_i^T \lambda_i = \mathbf{C}[\mathbf{x}] \mathbf{S}_i^T. \quad (4.39)$$

Dimensions with large eigenvalue have a large variation and can therefore be used to find features with useful information since we multiply the eigenvalue with the eigenvectors. When the eigenvalues are small, it means that the eigenvectors shrink accordingly and there is a small variation in these specific features.

So far, we have been leading up to the classical PCA theorem. Assume an orthogonal transformation $\mathbf{W} \in \mathbb{R}^{p \times p}$. We can then define the reconstruction error

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{n} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}}_i)^2, \quad (4.40)$$

with $\bar{\mathbf{x}}_i = \mathbf{W}\mathbf{z}_i$, where \mathbf{z}_i is a row vector with dimension \mathbb{R}^n of the matrix $\mathbf{Z} \in \mathbb{R}^{p \times n}$.

The PCA theorem states that minimizing the above reconstruction error corresponds to setting $\mathbf{W} = \mathbf{S}$, which is the orthogonal matrix which diagonalizes the covariance matrix [61]. The optimal number of features that corresponds to the encoding is given by the a set of vectors \mathbf{z}_i with at most l vectors. This is defined as the orthogonal projection of the data onto the columns by the eigenvectors of the covariance matrix. Instead of using the covariance matrix, it is preferable to use the correlation matrix to keep the numerical precision for any raw implementation. Additionally, it is important to mention that PCA is very sensitive to the data, which is why one should always remember to center the data around before applying PCA. We recommend the reader to read Ref. [61] p. 387 for proof of the classical PCA theorem, as we will not elaborate any further. The algorithm is shown in code-listing 4.

Algorithm 4: Principal component analysis algorithm.

Set up the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ with p features and n entries;
 Center the data by subtracting the mean value for each column;
 Compute the covariance matrix $\mathbb{E}[\mathbf{X}\mathbf{X}^T]$;
 Find the eigenpairs of \mathbf{C} with eigenvalues $[\lambda_0, \lambda_1, \dots, \lambda_{p-1}]$ and eigenvectors $[\mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_{p-1}]$;
 Order the eigenvalues, and therefore also the eigenvectors, in descending order. Keep only those l eigenvalues larger than a selected threshold value.

Instead of choosing an arbitrarily number of dimensions to reduce down to, it is common to choose the number of dimensions that accumulate a sufficiently amount of variance, such as e.g. 95%. However, it remains an subjective analysis in how many principal components one should include as it will depend on both the specific application and specific data set. If it is impossible to give a motivation for reducing a large dataset to just two or three principal components, there might still be a reason for why to apply PCA to a dataset. PCA can be applied as a preprocessing method to reduce the dimensionality of a dataset, and therefore might drastically improve the efficiency of further supervised learning approaches.

Part III

Methodology and implementation

Chapter 5

Material Science Databases

There are multiple different databases for material science discovery available for every day use, some of them completely open-source while others are commercial. This chapter will give a brief overview of databases available for computational material science, and will serve as a toolbox for how to request information and what kind of python packages exist to process that information.

5.1 Fundamentals of a database

A quick search online will reveal the tremendous escalation of effort for big-data driven material science the last few years, resulting in several databases that stores ab-initio calculation details and results. We will here distinguish between a *cloud service*, which is a place to store independent databases for research and commercial purposes, and a *database*, which is an organized collection of structured information. As an example, a cloud service can store several databases, but a database cannot host a cloud service.

To limit the quest of databases, we have restricted the search for databases and cloud services to include inorganic compounds obtained experimentally or by first-principles calculations, in particular DFT-calculations using *Vienna ab initio simulation package* (VASP) [67]. VASP is a software for atomic scale materials programming. Table 5.2 and 5.3 shows a selection of databases and cloud services that meets the given criteries, respectively.

5.1.1 API and HTTP requests

To extract information from a database it is convenient to interact through an *API* (Application Programming Interface), which defines important variables such as the kind of requests to be made, how to make them and the data format for transmission. Importantly, this permits communication between

different software medias. An API is entirely customizable, and can be made to extend existing functionality or tailor-made for specific user-demanding modules.

The APIs that will be encountered is handled by the use of *HTTP* (Hypertext Transfer Protocol), which in its simplest form is a protocol that allows the fetching of resources. The protocol is client-server based, such as the client is requesting information and the server is responding to the request.

The most common HTTP-methods are GET, POST and HEAD, which are used to either retrieve, send, or get information about data, respectively. The latter request is usually done before a GET-method for requests considering large amount of data, since this can be a significant variable for the client's bandwidth and load time. Following a request, the server normally responds with one of the status codes in table 5.1.

Status code	Description
2xx	OK - request was successful
3xx	Resource was redirected
4xx	Request failed due to either unsuccessful authentication or client error.
5xx	Request failed due to server error.

Table 5.1: Numeric status code for response. The leftmost digit decide the type of response, while the two follow-up digits depends on the implemented API.

A RESTful (Representational State Transfer) allows users to communicate with a server via a HTTP using a REST Architectural Style [68]. This enables the utilisation of Uniform Resource Identifiers (URI), where each object is represented as a unique resource and can be requested in a uniform manner. Importantly, this allows the use of both URIs and HTTP methods in an API, such that an object is represented by a unique URI whereas a HTTP-method can act on the object. This action will then return either the result of the action, or structured data that represents the object.

To provide a Python example, we can check the response by doing a GET request at the database Materials Project RESTful API in code listing 5.1. We use the preamble to version 2 of Materials Project, and add an API-check and an API-key. The response is shown in code listing 5.2. From the output, it is possible to tell that the supplied API-key is not valid, however, the request is valid.

```

1 import requests
2 preamble = "https://www.materialsproject.org/rest/v2/"
3 url = preamble + "api_check"
4 params = {"API_KEY": "unique_api_key"}
```

```

5 response = requests.get(url=url, params=params)
6 print(response.json())

```

Listing 5.1: Practical example of getting a response from Materials Project database.

```

1 {"valid_response": True,
2  "response":
3    {"api_key_valid": False,
4     "details": "API_KEY is not a valid key.",
5     "version":
6       {"db": "2020_09_08",
7        "pymatgen": "2020.8.13",
8        "rest": "2.0"}
9   }
10 }

```

Listing 5.2: Practical example of response from Materials Project request based on 5.1. The request was done 28. january 2020.

Database	API	Free educational access	Number of entries
AFLOW	REST	True	3.27 M
OQMD [69, 70]	RESTful API (qmpy, matminer)	True	0.82 M
MP [71]	MAPI [72]	True	0.71 M
ICSD [73]	RESTful API	False	0.21 M
Jarvis-DFT	API	True	0.04 M

Table 5.2: A selection of databases of computational material science sorted after number of compounds. Abbreviations used are Novel Materials Discovery (NOMAD), Automatic-FLOW for Materials Discovery (AFLOW), Materials Project (MP), Inorganic Crystal Structure Database (ICSD) and Open Quantum Materials Database (OQMD). The number of entries can give the wrong perception of size of each respective database, as it does not visualise how many calculations have been done for each entry, nor if there might be duplicates.

5.1.2 Open-source Python libraries for material analysis

Many of the databases share convenient modules that are used to adapt, visualize, calculate or predict properties, making it easier for scientists to utilise the databases. The Atomic Simulation Environment (ASE) is an environment

Cloud service	API/REST	Open educational access
NoMaD	API	True
CMR [74]	ASE	RESTful API
MatNavi	API	True
PRISMS	REST	True
Citrine	API	True
MPDS	API	False
MDF	API	False

Table 5.3: A selection of cloud services that offers database-storage. Abbreviations used are Computational Materials Repository (CMR), NIMS Materials Database (MatNavi), PRedictive Integrated Structural Materials Science (PRISMS), Materials Platform for Data Science (MPDS) and the Materials Data Facility (MDF).

in the Python programming language that includes several tools and modules for setting up, modifying and analyze atomistic simulations [75]. It is in particular used together with the cloud service Computational Materials Repository (CMR).

Another commonly used module is the Python Materials Genomics (pymatgen) [76]. This is a well-documented open module with both introductory and advanced use case examples written in Jupyter Notebook for easy reproducibility, and is integrated with the Materials Project RESTful API.

Another exceedingly popular library is matminer [77], which is an open-source toolkit for material analysis written in Python. Matminer is powered by a group known as *Hacking Materials Research Group*¹. Matminer provides modules to extract data sets from many cloud-services and databases, with examples in table 5.2 and 5.3. Additionally, they provide the tools to extract possibly thousands of features from calculations based on DFT and more, and have modules for visualization and automatic machine learning. These tools will be exemplified in the next chapter.

TODO : Add paragraph about sklearn.

A full selection of python libraries used and their versions can be found in the Github page (TODO: Add github page.)

5.2 Databases and cloud services

Every database has its own speciality, and no two databases are the same. There exists entries that are fundamentally identical in several databases, but

¹Project's Github site: <https://github.com/hackingmaterials>.

with different properties as a consequence of parameters used, such as the functional utilised in VASP or the relaxation scheme. This section digs up what exactly is each respective database's claim to fame.

5.2.1 Novel Materials Discovery

The Novel Materials Discovery (NOMAD) [78] Repository is an open-access platform for sharing and utilizing computational materials science data. NOMAD also consists of several branches such as NOMAD Archive, which is the representation of the NOMAD repository parsified into a code-independent format, NOMAD Encyclopedia, which is a graphical user interface (GUI) for characterizing materials, and lastly NOMAD Analytics Toolkit, which includes early-development examples of artificial-intelligence tools [78].

Databases that are a part of NOMAD data collection includes Materials Project, the Open Quantum Materials Database and AFLOW. They are all based on the underlying quantum engine VASP.

5.2.2 Materials project

Materials project [71] is an open source project that offers a variety of properties of over one hundred thousand of inorganic crystalline materials. It is known as the initiator of materials genomics and has as its mission to accelerate the discovery of new technological materials, with an emphasis on batteries and electrodes, through advanced scientific computation and innovative design.

Every compound has an initial relaxation of cell and lattice parameters performed using a 1000k-point mesh to ensure that all properties calculated are representative of the idealized unit cell for each respective crystal structure. The functional GGA is used to calculate band structures, while for transition metals it is applied +U correction to correct for correlation effects in d- and f-orbital systems that are not addressed by GGA calculations [79]. The thermodynamic stability for each phase with respect to decomposition, is also calculated. This is denoted as E Above Hull, with a value of zero is defined as the most stable phase at a given composition, while larger positive values indicate increased instability.

Each material contains multiple computations for different purposes, resulting in different 'tasks'. The reason behind this is that each computation has a purpose, such as to calculate the band structure or energy. Therefore, it is possible to receive several tasks for one material which results in more features per material.

5.2.3 AFLOW

The AFLOW[80–82] repository is an automatic software framework for the calculations of a wide range of inorganic material properties. They utilise the GGA-PBE functional within VASP with projector-augmented wavefunction (PAW) potentials to relax twice and optimize the ICSD-sourced structure. They are using a 3000 – 6000 k-point mesh, indicating a more computationally expensive calculation compared to the Materials Project. Next, the band structure is calculated with an even higher k-point density, in addition to the +U correction term for most occupied d- and f-orbital systems, resulting in a standard band gap [83]. Furthermore, they apply a standard fit gathered from a study of DFT-computed versus experimentally measured band gap widths to the initial calculated value, obtaining a fitted band gap [84].

AFLOW-ML [85] is an API that uses machine learning to predict thermo-mechanical and electronic properties based on the chemical composition and atomic structure alone, which they denote as *fragment descriptors*. They start with applying a classification model to predict if a compound is either a metal or an insulator, where the latter is confirmed with an additional regression model to predict the band gap width. To be able to predict properties on an independent data set, they utilise a fivefold cross validation process for each model. They report a 93% prediction success rate of their initial binary classification model, whereas the majority of the wrongful predictions are narrow-gap semiconductors. The authors does not compare their predicted band gap to experimental values, but it is found that 93% of the machine-learning-derived values are within 25% of the DFT +U-calculated band gap width [86].

5.2.4 Open Quantum Materials Database

The Open Quantum Materials Database (OQDM) [69, 70] is a free and available database of DFT-calculations. It has included thermodynamic and structural properties of more than 600.000 materials, including all unique entries in the Inorganic Crystal Structure Database (ICSD) consisting of less than 34 atoms.

The DFT calculations are performed with the VASP software whereas the electron exchange and correlation are described with the GGA-PBE, while using the PAW potentials. They relax a structure using 4000 – 8000 k-point mesh, indicating an even increasing computational expensive calculation than AFLOW again. Several element-specific settings are included such as using the +U extension for various transition metals, lanthanides and actinides. In addition, any calculation containing 3d or actinide elements are spin-polarized with a ferromagnetic alignment of spins to capture possible magnetism. However, the authors note that this approach does not capture com-

plex magnetic, such as antiferromagnetism, which has been found to result in substantial errors for the formation energy [87].

5.2.5 JARVIS

Joint Automated Repository for Various Integrated Simulations (JARVIS) [88] - DFT is an open database based on the VASP software to perform a variety of material property calculations. It consists of roughly 40,000 3D and 1,000 2D materials using the vdW-DF-OptB88 van der Waals functional, which was originally designed to improve the approximation of properties of two-dimensional van der Waals materials, but has also shown to be effective for bulk materials [89, 90]. The functional has shown accurate predictions for lattice-parameters and energetics for both vdW and non-vdW bonded materials [91].

Structures included in the data set are originally taken from the materials project, and then re-optimized using the OPT-functional. Finally, the combination of the OPT and modified Becke-Johnson (mBJ) functionals are used to obtain a representative band gap of each structure, since both have shown unprecedented accuracy in the calculation of band gap compared to any other DFT-based calculation methods [92].

The JARVIS-DFT database is part of a bigger platform that includes JARVIS-FF, which is the evaluation of classical forcefield with respect to DFT-data, and JARVIS-ML, which consists of 25 machine learning to predict properties of materials. In addition, JARVIS-DFT also includes a data set of 1D-nanowire and 0D-molecular materials, yet not publically distributed.

Chapter 6

Extraction and featurization of data

The initial step for gathering and building features can be visualised through the flochart in figure 6.1. Initially, we start by extracting all entries in the Materials Project that matches a specific query. Thereafter, we apply Matminer's featurization tools to make thousands of features of the data. In a parallel step, entries that are deemed similar to the entries from the initial Materials Project query are extracted from AFLOW, AFLOW-ML, JARVIS-DFT, OQMD and Citrine Informatics. Finally, we combine the steps together as interim data and prepare the data for further analysis.

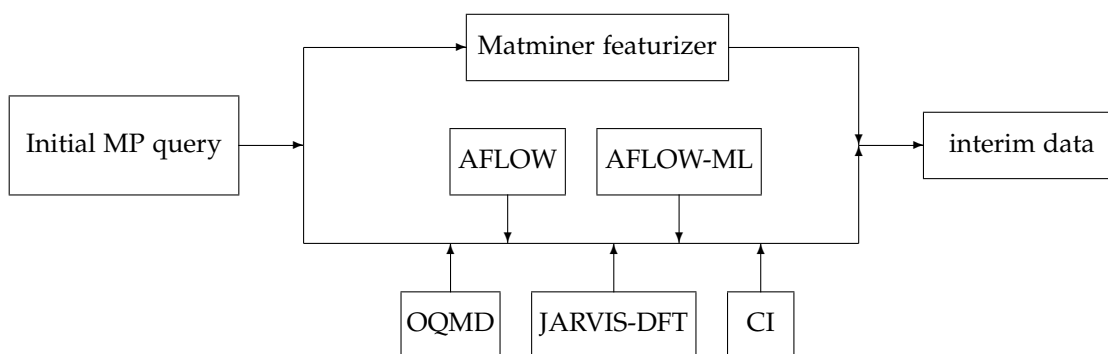


Figure 6.1: The data flow of the main project, starting from an initial MP-query, and ending with a featurized dataset with entries from several other databases. The matminer featurizer step is further visualized in-depth in figure 6.2.

The initial query has the requirement that all entries has to be derived from an experimental ICSD entry, and is reasoned by that we can identify equivalent entries in other databases. Furthermore, all entries in the Materials Project needs to have a band gap larger than 0.1eV. Recall that Materials Project applies the functional GGA in estimating the band gap, which is known to severely underestimate the given electronic property. Therefore, we have chosen a low value to not rule out any potential candidates but high enough to leave out all materials that can be considered metallic.

From figure 6.1 we notice that by using many databases we do not add additional entries that exist in some databases but is not to be found in Materials Project. This is by design since it preserves the versatility of choosing a database to work with. Therefore, one can completely ignore steps such as the initial query of Materials Project or the featurization process, and rather focus on e.g. all the 400.000 entries existing in OQMD. The examples that follows will illustrate the ease of extracting data from several different databases, and can serve as the starting point for other research projects in computational material science.

6.1 Practical data extraction with Python-examples

For this section, we will show practical examples of how to extract data that might fulfill the criteria for a material to host a qubit candidate given in the theory part. We will begin with the database of Materials Project, and then search for entries in other databases that match entries from MP. This process is reproducible as a jupyter notebook¹ and the databases in question are the ones referred to in the previous section.

Instead of building multiple HTTP-methods from scratch, we will here take a look at the easiest method at obtaining data from each database. This includes looking into the APIs that supports data-extraction and that are recommended by each respective database.

The range of data in a database can consist of data from a few entries up to an unlimited amount of entries with even further optional parameters, and has limitless use in applications. However, the amount of data in a database is irrelevant if the data is inaccessible. Therefore, we provide a toolbox in how to extract information in the easiest way possible.

Every data extraction class is based on an abstract parent class, which is listed in code listing 6.1. The advantages of using a base parent class are many, since it improves the readability during code reviews and reduce the main barrier for understanding the underlying structure of a project, while utilising reusable components. Yet, the main advantage of using a base parent

¹add and insert DOI for JN 01-generateDataset-notebook.ipynb

class is the fact that it can effortlessly be extended for further implementations since it provides a code skeleton.

The structure of extraction is centered around using the data extraction tools, and not understanding them. Therefore, we only show how to use them here, while the code is found in the Appendix.

```

1  import abc
2  import pandas as pd
3  from typing import Optional, Iterable, Tuple, Dict
4  import os
5  __all__ = ("data_base", )
6
7  class data_base(abc.ABC):
8      # TODO: ADD COMMENTS.
9      data_dir : Optional[str] = None
10     raw_data_path : Optional[str] = None
11     interim_data_path : Optional[str] = None
12
13     df : Optional[pd.DataFrame] = None
14
15     def _does_file_exist(self)-> bool:
16         if os.path.exists(self.raw_data_path):
17             print("Data path {} detected. Reading now...".
18                   format(self.raw_data_path))
19             return True
20         else:
21             print("Data for MP not detected. Applying query
22                   now... ")
23             return False #self.get_data()
24
25     def get_dataframe(self, sorted: Optional[bool] = True)->
26         pd.DataFrame:
27
28         if self._does_file_exist():
29             self.df = pd.read_pickle(self.raw_data_path)
30         else:
31             self.df = self._apply_query(sorted=sorted)
32         print("Done")
33         return(self.df)

```

Listing 6.1: Base parent class of all data extraction classes.

6.1.1 Materials Project

The most up-to-date version of Materials Project can be extracted using the python package pymatgen, which is integrated with Materials Project REST API. Other retrieval tools that is dependent on pymatgen includes matminer, with the added functionality of returning a pandas dataframe. Copies of

Materials Project are added frequently to cloud services such as Citrine Informatics, but the latest added entries to Materials Project cannot be guaranteed in such a query.

Entries in Materials Project are characterized using more than 60 features², some features being irrelevant for some materials while fundamental for others. The data is divided into three different branches, where the first can be described as basic properties of materials including over 30 features, while the second branch describes experimental thermochemical information. The last branch yields information about a particular calculation, in particular information that's relevant for running a DFT script.

To extract information from the database, we will be utilising the module `pymatgen`. This query supports MongoDB query and projection operators³, resulting in an almost instant query.

1. Register for an account⁴, and generate a secret API-key.
2. Set the required criteria.
3. Set the wanted properties.
4. Apply the query.

The code nipplet in code listing 6.2 resembles steps 2 – 4, and is filtered as the initial query.

```
1 from src.data.get_data_MP import data_MP
2
3 MAPI_KEY = 'very_secret_key_here'
4 MP = data_MP(API_KEY=MAPI_KEY)
5 df = MP.get_dataframe()
```

Listing 6.2: Practical example of extracting information from Materials Project using `pymatgen`, resulting in a Pandas DataFrame named `entries` that contains the properties given after performing a filter on the database. The criteria is given as a JSON, and supports MongoDB operators.

6.1.2 Citrine Informatics

Citrine Informatics is a cloud service, which means that the spectrum of stored information varies broadly. We will access research through open access for institutional and educational purposes. Information in Citrine can

²All features can be viewed in the documentation of the project: <https://github.com/materialsproject/mapidoc/master/materials>

³<https://docs.mongodb.com/manual/reference/operator/query/>

⁴<https://materialsproject.org>

be stored using a scheme that is broken down into two sections, with private properties for each entry in addition to common fields that are the same for all entries.

In this example, we will gather experimental data using the module `matminer`. The following steps are required to extract information from Citrine Informatics.

1. Register for an account⁵, and generate a secret API-key.
2. Set the required criteria.
3. Set the wanted properties and common fields.
4. Apply the query.

The code listed in code listing 6.3 gives an easy example to steps 2 – 4 with experimental data as filter.

```
1 from src.data.get_data_Citrine import data_Citrine
2
3 CAPI_KEY = 'very_secret_key_here'
4 citrine = data_Citrine(API_KEY=CAPI_KEY)
5 df = citrine.get_dataframe()
```

Listing 6.3: Practical example of extracting information from Citrine Informatics using `matminer`, resulting in a Pandas DataFrame named `experimental_entries` that contains the properties given after performing a filter on the database. The criteria is given as a JSON.

6.1.3 AFLOW

The query from AFLOW API [80] supports lazy formatting, which means that the query is just a search and does not return values but rather an object. This object is then used in the query when asking for values. For every object it is necessary to request the desired property, consequently making the query process significantly more time-demanding than similar queries using APIs such as `pymatgen` or `matminer` for Citrine Informatics. Hence, the accessibility is strictly limited to either searching for single compounds or if the user possess sufficient time.

`Matminer`'s data retrieval tool for AFLOW is currently an ongoing issue [93], thus we present in code listing 6.4 a function that extracts information from AFLOW and returns a Pandas DataFrame. In contrast to Materials Project and Citrine Informatics, AFLOW does not require an API-key for a

⁵<https://citrination.com>

query, which reduces the amount of steps to obtain data. The class searches for an stored AFLOW-data, and initialises a MP-query with the initial criteria if not successful. The resulting query will then be used as input to AFLOW.

```
1 from src.data.get_data_AFLOW import data_AFLOW
2
3 AFLOW = data_AFLOW()
4 df = AFLOW.get_dataframe()
```

Listing 6.4: Practical example of extracting information from AFLOW. The function can extract all information in AFLOW for a given list of compounds, however, it is a slow method and requires consistent internet connection.

6.1.4 AFLOW-ML

In this part, we will be using a machine learning algorithm named AFLOW-ML Property Labeled Material Fragments (PLMF) [85] to predict the band gap of structures. This algorithm is compatible with a POSCAR of a compound, which can be generated by the CIF (Crystallographic Information File) that describes a crystal's generic structure. It is possible to download a structure as a poscar by using Materials Project front-end API, but is a cumbersome process to do so individually if the task includes many structures. Extracting the feature of POSCAR is yet to be implemented in the RESful API of pymatgen, thus we demonstrate the versatility of pymatgen with a workaround.

We begin with extracting the desired compounds formula, its material_id for identification, and their respectful structure in CIF-format from Materials Project. In an iterative process, each CIF-structure is parsed to a pymatgen structure, where pymatgen can read and convert the structure to a POSCAR stored as a Python dictionary. Finally, we can use the POSCAR as input to AFLOW-ML, which will return the predicted band gap of the structure. This iterative process parsing and converting, but is an undemanding process. The function that handles this is presented in code listing 6.5. Similar to AFLOW-query, this code listing is dependent on MP-data and will apply for a query if the data is not present.

A significant portion of the process is tied up to obtaining the input-file for AFLOW-ML, and fewer structures will result in an easier process. Nevertheless, we present the following steps in order to receive data from AFLOW-ML.

1. Download AFLOWmlAPI⁶.
2. Getting POSCAR from MP.

⁶<http://aflow.org/src/aflow-ml/> to the same directory as code listing 6.5

- (a) Apply the query from Materials Project with "CIF", "material_id" and "full_formula" as properties.
- (b) Insert resulting DataFrame into function defined in code listing 6.5.

3. Insert POSCAR to AFLOW-ML.

```
1 from src.data.get_data_AFLOWML import data_AFLOWML
2
3 AFLOWML = data_AFLOWML()
4 df = AFLOWML.get_dataframe()
```

Listing 6.5: Practical example of extracting information from AFLOW-ML. The function will convert a CIF-file (from e.g. Materials Project) to a POSCAR, and will use it as input to AFLOW-ML. In return, one will get the structure's predicted band gap. It should be noted that this requires the AFLOW-ML library in the same directory.

6.1.5 JARVIS-DFT

The newest version of the JARVIS-DFT dataset can be obtained by requesting an account at the official webpage, but with the drawback that an administrator has to either accept or deny the request. Thus, the accessibility of the database is dependent on if there is an active administrator paying attention to the requests, which is a limitation experienced during this work. Another approach is to download the database through matminer, however with the limitation of not necessarily having the latest version of the database. A third approach is to download a version of JARVIS-DFT that have been made available for requests the 30.04.2020 at <http://figshare.com> by Choudhary *et al.* [88]. The author provides tools for extraction, yet not compatible with the latest version of Python (3.8) at the time writing (12.03.2021). Therefore, we provide a tool to extract this data through the use of our base class.

```
1 from src.data.get_data_JARVIS import data_JARVIS
2
3 JARVIS = data_JARVIS()
4 df = JARVIS.get_dataframe()
```

Listing 6.6: Practical example of extracting information from JARVIS-DFT. For this example, we exclude all metals by removing all non-measured band gaps.

We observe that there is no advanced search filter when loading the database from matminer. The author of matminer regards this as the user's task, and is indeed easily done through the use of the python library Pandas.

6.2 Matminer featurization

Before applying any machine learning algorithm, raw data needs to be transformed into a numerical representation that reflects the relationship between the input and output data. This transformation is known as generating descriptors or features, however, we will in this work adapt the name *featurization*. The open source library of Matminer provides many tools to featurize existing features extracted from Materials Project. In this section we will describe how to extract the features from an initial Materials Project query result (see subsection. 6.1.1), and the resulting features. It is beyond the scope of this work to go in-depth of each feature since the resulting dataset contains a quantity of several thousand features, but we will here take the liberty to serve a brief overview of the features and refer to each respective citation for more information. The respective table with information regarding 39 distinct matminer featurizers is situated in the Appendix, table A.1.

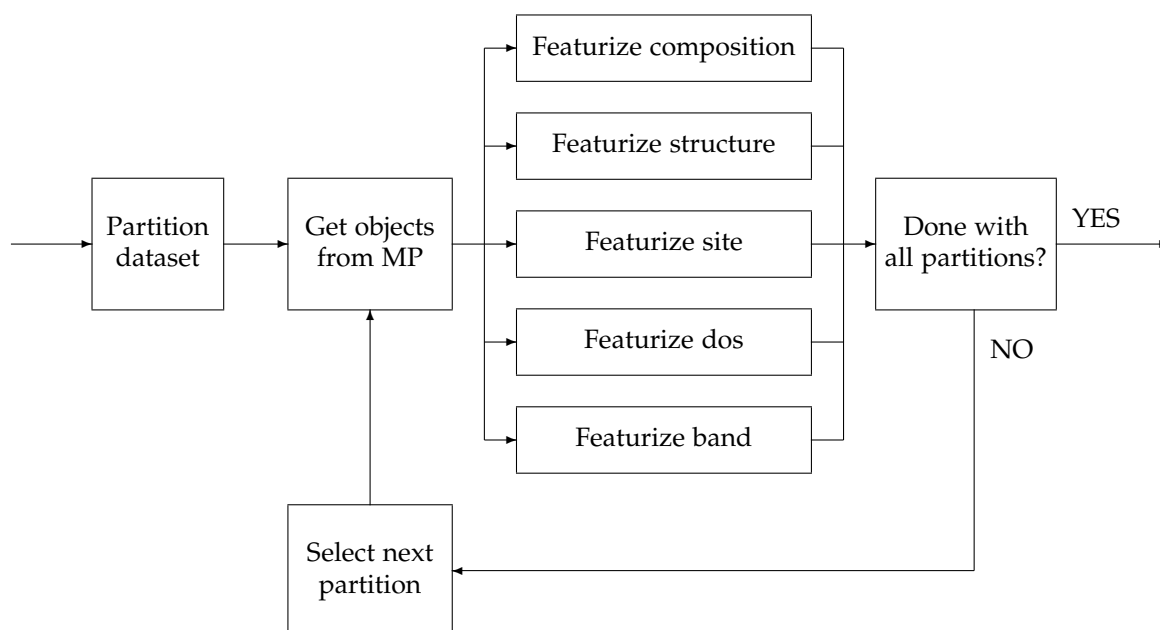


Figure 6.2: The process of the matminer featurizer step as seen in figure 6.1. To limit the memory and computational usage, the data is partitioned into smaller subsets where the respective pymatgen objects are obtained through a query to be used in the following featurization steps. This is iteratively done until all the data has been featurized.

To apply matminer’s featurization tools, we extend an existing implementation by De Breuck *et al.* [94] called the Materials Optimal Descriptor Network (MODNet). The author specifies that MODNet is a supervised machine learning framework for learning material properties based on either composition or crystal structure. To provide the training data for their model, MODNet featurizes (through matminer) structures either from Materials Project or in the form of a structure object made by pymatgen. Their current implementation provides featurization for compositions, structures and sites. However, matminer also provides featurization tools for density of states (DOS) and band structures, therefore we modify MODNet and extend it to facilitate such featurizations.

One immediate limitation of our extension is that Matminer’s tools is dependent on a pymatgen DOS- and bandstructure object. These objects contains information up to 5MB, and becomes a challenge when dealing with data containing several thousand such objects. This is solved by the required features for matminer’s featurization for a subsample of the data, followed by a featurization process of the same subsample. When the featurization is done, we store the new features and throw away the pymatgen features. This is done iteratively for the entire data set. Thus, a compromise between applying several queries and storing information has been done. The scheme can be visualised as the flow chart seen in figure 6.2.

In the extended version of the featurization process, we eliminate all columns that does not have any entries with physical meaning. This is beneficial for several reasons, such as to reduce memory allocated and to preprocess the data. If there are entries existing with both physical and non-physical for the same column, we replace the non-physical meanings with -1 for recognition in a later step. Additionally, we convert columns that are categorical or lacks a numerical representation into a categorical portrayal. Thus, we strive to limit the necessary steps for further processing of data into a machine learning algorithm.

Even if the first version of Matminer was released in 2016, many issues concerning daily operational use are still present. During the featurization process in this work, we manually identified 14 erroneous entries that are summarized in the Appendix, table A.2. These entries were excluded from the dataset.

6.3 Preprocessing

6.4 Screen procedure

kaffe

Part IV

Results

Chapter 7

Validation

A thorough testing procedure is important to find out if the code is working as intended. The procedure might reveal the presence or absence of bugs, and as a project grows, it can give an indication if a new implementation breaks the original project. Therefore, we present a test-case scenario to test if the two machine learning algorithms are able to predict the correct label. It is the same algorithm that will be used in the following chapters, and it will provide us the opportunity to understand how the algorithm works and to draw parallels between the separate works.

The validation process is a reproduction of Ref [95]. To be able to draw any parallel to their work, we use the exact same dataset in the beginning phase. It should be noted that even if the computational aspects of the validation is closely related to Ref. [95], the work eventually diverges in terms of focus. In their work they include a stability analysis using convex hull analysis in DFT calculations from OQMD, however, we will in this thesis not decide whether a compound is considered stable or not in an atomic configuration.

7.1 The perovskite dataset

The dataset in question contains 390 experimentally reported ABO_3 compounds. All compounds are charged balanced, and for every compound there is a feature explaining which structure the compound takes, either being a cubic perovskite, perovskite, or not a perovskite at all. Of the 390 compounds, there are 254 perovskites and 136 non-perovskites. Of the 254 perovskites, 232 takes a non-cubic perovskite structure while only 22 takes the cubic perovskite structure. Consequently, this will be visualized by two columns named Perovskite, which represents if a compound is either perovskite (1) or not perovskite (0), and Cubic, which represents if a compound is cubic perovskite (1), non-cubic perovskite (-1), or not perovskite(0).

7.1.1 Features

There are in total 9 features we can train a model on. Many of the features are based on the Shannon ionic radii [96], which are estimates of an element's ionic hard-sphere radii extracted from experiment. They are dimensionless numbers, and are frequently used in studies involving perovskite structures of materials since they can be a measurement of the ionic misfit of the B atom. This can be used to find the deviation of the structure from an ideal cubic geometry. The octahedral factor for an ABO_3 solid is known as

$$O = \frac{r_b}{r_O}, \quad (7.1)$$

where r_b and r_O are the Shannon radii for the B-atom and oxygen ($r_O = 1.4\text{\AA}$), respectively. If the octahedral factor is $O = 0.435$, it corresponds to a hard-sphere closed-packed arrangement where B and O ions are touching, while a six-fold coordination appear to require $0.414 < O < 0.732$ according to empirical studies [97]. O , r_A and r_b are represented as features in our data set. We can also compute the Goldschmidt tolerance factor, which is defined as

$$t = \frac{r_A + r_O}{\sqrt{2}(r_A + r_O)}. \quad (7.2)$$

The tolerance factor favors the following structures in the interval:

- $t > 1$: Hexagonal nonperovskite.
- $0.9 < t < 1.0$: Cubic perovskite.
- $0.75 < t < 0.9$: Orthorombic perovskite.
- $t < 0.75$: Not a perovskite.

If the tolerance factor is exactly $t = 1$, the structure is known as perfectly cubic and is free for any structural alterations.

Furthermore, the Shannon radii r_A and r_B can be directly correlated with the structure. Perovskites require $r_A > r_B$, and that A-atoms are in a 12-fold coordinated site if $r_A > 0.9\text{\AA}$. A-atoms also occur in a sixfold coordinated site if $r_A < 0.8\text{\AA}$ and $r_B > 0.7\text{\AA}$.

From bond valence theory we can find the valence of an ion to be the sum of valences, that is

$$V_i = \sum_i v_{ij} \quad (7.3)$$

$$= \sum_i \frac{\exp(d_o - d_{ij})}{b}, \quad (7.4)$$

where d_{ij} is the bond length while d_o and b are parameters from experimental data. The bond length can be found from 7.4 given the general value $b = 1.4\text{\AA}$ and d_o , that can be found from Zhang *et al.* database [97]. The valence of an ion is associated with its neighboring ions and the chemical bonds, and therefore the bond length d_{AO} and d_{BO} are included in the data set.

The two last features originates from the Mendeleev numbers of Villars *et al.* [98] for the A- and B atom. The given values positions the elements in structurally similar groups. This means that he groups the elements in the following interval.

- s-block $\in \{1, 10\}$.
- Sc = 11.
- Y = 12.
- f-block $\in \{13, 42\}$.
- d-block $\in \{43, 66\}$.
- p-block $\in \{67, 10\}$.

Part V

Appendices

Appendix A

Featurizaton

A.1 Table of featurizers

Table A.1: This thesis' chosen 39 featurizers from matminer. Descriptions are either found from Ref. [77] or from the project's Github page.

Features	Description	Original reference
Composition features		
		Continued on next page

Table A.1 – continued from previous page

Features	Description	Original reference
AtomicOrbitals	Highest occupied molecular orbital (HOMO) and lowest unoccupied molecular orbital (LUMO).	[99]
AtomicPacking-Efficiency	Packing efficiency.	[100]
BandCenter	Estimation of absolute position of band center using geometric mean of electronegativity.	[101]
ElementFraction	Fraction of each element in a composition.	-
ElementProperty	Statistics of various element properties.	[76, 102, 103]
IonProperty	Maximum and average ionic character.	[102]
Miedema	Formation enthalpies of intermetallic compounds, solid solutions, and amorphous phases using semi-empirical Miedema model.	[104]
Stoichiometry	L ^p norm-based stoichiometric attributes.	[102]
TMetalFraction	Fraction of magnetic transition metals.	[103]
ValenceOrbital	Valence orbital attributes such as the mean number of electrons in each shell.	[102]
YangSolid-Solution	Mixing thermochemistry and size mismatch terms.	[105]
Oxid composition features		
Electronegativity-Diff	Statistics on electronegativity difference between anions and cations.	[103]
OxidationStates	Statistics of oxidation states.	[103]
Structure features		
Continued on next page		

Table A.1 – continued from previous page

Features	Description	Original reference
DensityFeatures	Calculate density, volume per atom and packing fraction.	-
GlobalSymmetry-Features	Determines spacegroup number, crystal system (1-7) and inversion symmetry.	-
RadialDistribution-Function	Calculates the radial distribution function of a crystal system.	-
CoulombMatrix	Generate the Coulomb matrix, which is a representation of the nuclear coulombic interaction of the input structure.	[106]
PartialRadial-Distribution-Function	Compute the partial radial distribution function of a crystal structure	[107]
SineCoulomb-Matrix	Computes a variant of the coulomb matrix developed for periodic crystals.	[108]
EwaldEnergy	Computes the energy from Coulombic interactions based on charge states of each site.	[109]
BondFractions	Compute the fraction of each bond in a structure, based on nearest neighbours.	[110]
Structural-Heterogeneity	Calculates the variance in bond lengths and atomic volumes in a structure.	[111]
MaximumPacking-Efficiency	Calculates the maximum packing efficiency of a structure.	[111]
ChemicalOrdering	Computes how much the ordering of species differs from random in a structure.	[111]
XRDPowder-Pattern	1D array representing normalized powder diffraction of a structure as calculated by pymatgen.	[76]
Site features		
Continued on next page		

Table A.1 – continued from previous page

Features	Description	Original reference
AGNI-Fingerprints	Calculates the product integral of RDF and Gaussian window function	[112]
AverageBond-Angle	Determines the average bond angle of a specific site with its nearest neighbors using pymatgens implementation.	[113]
AverageBond-Length	Determines the average bond length between one specific site and all its nearest neighbors using pymatgens implementation.	[113]
BondOrientational-Paramater	Calculates the averages of spherical harmonics of local neighbors	[114, 115]
ChemEnvSite Fingerprint	Calculates the resemblance of given sites to ideal environment using pymatgens ChemEnv package.	[116, 117]
Coordination-Number	The number of first nearest neighbors of a site	[117]
CrystalNN-Fingerprint	A local order parameter fingerprint for periodic crystals.	-
GaussianSymm-Func	Calculates the gaussian radial and angular symmetry functions originally suggested for fitting machine learning potentials.	[118, 119]
GeneralizedRadial-Distribution-Function	Computes the general radial distribution function for a site	[114]
LocalProperty-Difference	Computes the difference in elemental properties between a site and its neighboring sites.	[111, 113]
OPSite-Fingerprint	Computes the local structure order parameters from a site's neighbor environment.	[117]
Continued on next page		

Table A.1 – continued from previous page

Features	Description	Original reference
Voronoi-Fingerprint	Calculates the Voronoi tessellation-based features around a target site.	[120, 121]
Density of state features		
DOSFeaturizer	Computes top contributors to the density of states at the valence and conduction band edges. Thus includes chemical specie, orbital character, and orbital location information.	[122]
Band structure features		
BandFeaturizer	Converts a complex electronic band structure into quantities such as band gap and the norm of k point coordinates at which the conduction band minimum and valence band maximum occur.	-

A.2 Erroneous entries

MPID	Full formula	Reference
mp-555563	$\text{PH}_6\text{C}_2\text{S}_2\text{NCl}_2\text{O}_4$	[123]
mp-583476	$\text{Nb}_7\text{S}_2\text{I}_{19}$	[124]
mp-600205	$\text{H}_{10}\text{C}_5\text{SeS}_2\text{N}_3\text{Cl}$	-
mp-600217	$\text{H}_{80}\text{C}_{40}\text{Se}_8\text{S}_{16}\text{Br}_8\text{N}_{24}$	-
mp-1195290	$\text{Ga}_3\text{Si}_5\text{P}_{10}\text{H}_{36}\text{C}_{12}\text{N}_4\text{Cl}_{11}$	-
mp-1196358	$\text{P}_4\text{H}_{120}\text{Pt}_8\text{C}_{40}\text{I}_8\text{N}_4\text{Cl}_8$	-
mp-1196439	$\text{Sn}_8\text{P}_4\text{H}_{128}\text{C}_{44}\text{N}_{12}\text{Cl}_8\text{O}_4$	-
mp-1198652	$\text{Te}_4\text{H}_{72}\text{C}_{36}\text{S}_{24}\text{N}_{12}\text{Cl}_4$	-
mp-1198926	$\text{Re}_8\text{H}_{96}\text{C}_{24}\text{S}_{24}\text{N}_{48}\text{Cl}_{48}$	-
mp-1199490	$\text{Mn}_4\text{H}_{64}\text{C}_{16}\text{S}_{16}\text{N}_{32}\text{Cl}_8$	-
mp-1199686	$\text{Mo}_4\text{P}_{16}\text{H}_{152}\text{C}_{52}\text{N}_{16}\text{Cl}_{16}$	-
mp-1203403	$\text{C}_{121}\text{S}_2\text{Cl}_{20}$	-
mp-1204279	$\text{Si}_{16}\text{Te}_8\text{H}_{176}\text{Pd}_8\text{C}_{64}\text{Cl}_{16}$	-
mp-1204629	$\text{P}_{16}\text{H}_{216}\text{C}_{80}\text{N}_{32}\text{Cl}_8$	-

Table A.2: A table of manually identified entries from Materials Project that experience issues concerning Matminer’s featurization tools. These were excluded from the dataset.

Bibliography

1. Griffiths, D. *Introduction to quantum mechanics* ISBN: 9781107179868 (Cambridge University Press, Cambridge, 2017).
2. Turing, A. M. On computable numbers, with an application to the Entscheidungsproblem. *Proceedings of the London mathematical society* **2**, 230–265 (1937).
3. Moore, G. Cramming More Components Onto Integrated Circuits. *Proceedings of the IEEE* **86**, 82–85 (Jan. 1965).
4. Pavičić, M. *Quantum computation and quantum communication : theory and experiments* ISBN: 9786610743704 (Springer, New York, 2006).
5. Gwennap, L. Apple's 5 Nanometer Chip Is Another Signpost That Moore's Law Is Running Out. *Forbes*. <<https://www.forbes.com/sites/linleygwennap/2020/10/12/apple-moores-law-is-running-out/>> (Oct. 12, 2020).
6. Weber, J. R. *et al.* Quantum computing with defects. *Proceedings of the National Academy of Sciences* **107**, 8513–8518 (Apr. 2010).
7. DiVincenzo, D. P. The Physical Implementation of Quantum Computation. *Fortschritte der Physik* **48**, 771–783 (Sept. 2000).
8. Ladd, T. D. *et al.* Quantum computers. *Nature* **464**, 45–53 (Mar. 2010).
9. Mizel, A., Lidar, D. A. & Mitchell, M. Simple Proof of Equivalence between Adiabatic Quantum Computation and the Circuit Model. *Physical Review Letters* **99**. doi:[10.1103/physrevlett.99.070502](https://doi.org/10.1103/physrevlett.99.070502) (Aug. 2007).
10. Grover, L. K. A framework for fast quantum mechanical algorithms. arXiv: [quant-ph/9711043v2](https://arxiv.org/abs/quant-ph/9711043v2) [[quant-ph](https://arxiv.org/abs/quant-ph)] (Nov. 20, 1997).
11. Shor, P. *Algorithms for quantum computation: discrete logarithms and factoring* in *Proceedings 35th Annual Symposium on Foundations of Computer Science* (IEEE Comput. Soc. Press, 1994). doi:[10.1109/sfcs.1994.365700](https://doi.org/10.1109/sfcs.1994.365700).
12. Martinis, J. M. *et al.* Quantum supremacy using a programmable superconducting processor en. 2019. doi:[10.5061/DRYAD.K6T1RJ8](https://doi.org/10.5061/DRYAD.K6T1RJ8).
13. Griffiths, R. B. Nature and location of quantum information. *Physical Review A* **66**. doi:[10.1103/physreva.66.012311](https://doi.org/10.1103/physreva.66.012311) (July 2002).

14. Gisin, N., Ribordy, G., Tittel, W. & Zbinden, H. Quantum cryptography. *Reviews of Modern Physics* **74**, 145–195 (Mar. 2002).
15. Gisin, N. & Thew, R. Quantum communication. *Nature Photonics* **1**, 165–171 (Mar. 2007).
16. Acín, A. *et al.* The quantum technologies roadmap: a European community view. *New Journal of Physics* **20**, 080201 (Aug. 2018).
17. Boaron, A. *et al.* Secure Quantum Key Distribution over 421 km of Optical Fiber. *Physical Review Letters* **121**. doi:[10.1103/physrevlett.121.190502](https://doi.org/10.1103/physrevlett.121.190502) (Nov. 2018).
18. Degen, C. L., Reinhard, F. & Cappellaro, P. Quantum sensing. *Reviews of Modern Physics* **89**. doi:[10.1103/revmodphys.89.035002](https://doi.org/10.1103/revmodphys.89.035002) (July 2017).
19. Georgescu, I. The DiVincenzo criteria 20 years on. *Nature Reviews Physics* **2**, 666–666 (Nov. 2020).
20. Kristian Fossheim, A. S. *Superconductivity: Physics and Applications* 442 pp. ISBN: 0470844523. <https://www.ebook.de/de/product/3608091/kristian_fossheim_asle_sudboe_superconductivity_physics_and_applications.html> (WILEY, 2004).
21. Ben Streetman, S. B. *Solid State Electronic Devices, Global Edition* 632 pp. ISBN: 1292060557. <https://www.ebook.de/de/product/30394493/ben_streetman_sanjay_banerjee_solid_state_electronic_devices_global_edition.html> (Pearson Education Limited, 2015).
22. Renganathan, G., Tanneru, N. & Madurai, S. L. in *Fundamental Biomaterials: Metals* 211–241 (Elsevier, 2018). doi:[10.1016/b978-0-08-102205-4.00010-6](https://doi.org/10.1016/b978-0-08-102205-4.00010-6).
23. Lufaso, M. W. & Woodward, P. M. Prediction of the crystal structures of perovskites using the software program SPuDS. *Acta Crystallographica Section B Structural Science* **57**, 725–738 (Nov. 2001).
24. Bednorz, J. G. & Müller, K. A. Perovskite-type oxides—The new approach to high-Tc superconductivity. *Reviews of Modern Physics* **60**, 585–600 (July 1988).
25. Boivin, J. C. & Mairesse, G. Recent Material Developments in Fast Oxide Ion Conductors. *Chemistry of Materials* **10**, 2870–2888 (Oct. 1998).
26. Cheong, S.-W. & Mostovoy, M. Multiferroics: a magnetic twist for ferroelectricity. *Nature Materials* **6**, 13–20 (Jan. 2007).
27. Ibn-Mohammed, T. *et al.* Perovskite solar cells: An integrated hybrid lifecycle assessment and review in comparison with other photovoltaic technologies. *Renewable and Sustainable Energy Reviews* **80**, 1321–1344 (Dec. 2017).

28. Chen, P.-Y. *et al.* Environmentally responsible fabrication of efficient perovskite solar cells from recycled car batteries. *Energy Environ. Sci.* **7**, 3659–3665 (2014).
29. Pauli, W. Über den Zusammenhang des Abschlusses der Elektronengruppen im Atom mit der Komplexstruktur der Spektren. *Zeitschrift für Physik* **31**, 765–783 (Feb. 1925).
30. Martienssen, W. *Springer handbook of condensed matter and materials data* ISBN: 9786610625949 (Springer, Heidelberg New York, 2005).
31. Pelant, I. *Luminescence spectroscopy of semiconductors* ISBN: 0191738549 (Oxford University Press, Oxford, 2012).
32. Kun Huang, A. R. Theory of light absorption and non-radiative transitions in F-centres. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* **204**, 406–423 (Dec. 1950).
33. Gordon, L. *et al.* Quantum computing with defects. *MRS Bulletin* **38**, 802–807 (Oct. 2013).
34. Bernien, H. *et al.* Heralded entanglement between solid-state qubits separated by three metres. *Nature* **497**, 86–90 (Apr. 2013).
35. Taylor, J. M. *et al.* High-sensitivity diamond magnetometer with nanoscale resolution. *Nature Physics* **4**, 810–816 (Sept. 2008).
36. Barclay, P. E., Fu, K.-M. C., Santori, C., Faraon, A. & Beausoleil, R. G. Hybrid Nanocavity Resonant Enhancement of Color Center Emission in Diamond. *Physical Review X* **1**. doi:[10.1103/physrevx.1.011007](https://doi.org/10.1103/physrevx.1.011007) (Sept. 2011).
37. Bassett, L. C., Alkauskas, A., Exarhos, A. L. & Fu, K.-M. C. Quantum defects by design. *Nanophotonics* **8**, 1867–1888 (Oct. 2019).
38. James, W. J. Theory of defects in solids by A. M. Stoneham. *Acta Crystallographica Section A* **32**, 527–527 (May 1976).
39. Togan, E. *et al.* Quantum entanglement between an optical photon and a solid-state spin qubit. *Nature* **466**, 730–734 (Aug. 2010).
40. Neudeck, P. G. Progress in silicon carbide semiconductor electronics technology. *Journal of Electronic Materials* **24**, 283–288 (Apr. 1995).
41. Silveira, E., Freitas, J. A., Glembocki, O. J., Slack, G. A. & Schowalter, L. J. Excitonic structure of bulk AlN from optical reflectivity and cathodoluminescence measurements. *Physical Review B* **71**. doi:[10.1103/physrevb.71.041201](https://doi.org/10.1103/physrevb.71.041201) (Jan. 2005).
42. Lawaetz, P. Valence-Band Parameters in Cubic Semiconductors. *Physical Review B* **4**, 3460–3467 (Nov. 1971).

43. Beckers, L. *et al.* Structural and optical characterization of epitaxial waveguiding BaTiO₃ thin films on MgO. *Journal of Applied Physics* **83**, 3305–3310 (Mar. 1998).
44. Kumbhojkar, N., Nikesh, V. V., Kshirsagar, A. & Mahamuni, S. Photo-physical properties of ZnS nanoclusters. *Journal of Applied Physics* **88**, 6260–6264 (Dec. 2000).
45. Persson, C. *Brief Introduction to the density functional theory* 2020.
46. Top500. SUPERCOMPUTER FUGAKU June 2020. <<https://www.top500.org/system/179807/>> (visited on 10/02/2020).
47. Hohenberg, P. & Kohn, W. Inhomogeneous Electron Gas. *Physical Review* **136**, B864–B871 (Nov. 1964).
48. Toulouse, J. *Introduction to density-functional theory* Sept. 2019. <http://www.lct.jussieu.fr/pagesperso/toulouse/enseignement/introduction_dft.pdf> (visited on 10/25/2020).
49. Kohn, W. & Sham, L. J. Self-Consistent Equations Including Exchange and Correlation Effects. *Physical Review* **140**, A1133–A1138 (Nov. 1965).
50. Allen, J. P. & Watson, G. W. Occupation matrix control of d- and f-electron localisations using DFT+U. *Phys. Chem. Chem. Phys.* **16**, 21016–21031 (2014).
51. David Sholl, J. A. S. *Density Functional Theory: A Practical Introduction* 238 pp. ISBN: 0470373172. <https://www.ebook.de/de/product/7207845/david_sholl_janice_a_steckel_density_functional_theory_a_practical_introduction.html> (WILEY, 2009).
52. Perdew, J. P. & Wang, Y. Accurate and simple analytic representation of the electron-gas correlation energy. *Physical Review B* **45**, 13244–13249 (June 1992).
53. Perdew, J. P., Burke, K. & Ernzerhof, M. Generalized Gradient Approximation Made Simple. *Physical Review Letters* **77**, 3865–3868 (Oct. 1996).
54. Freysoldt, C. *et al.* First-principles calculations for point defects in solids. *Reviews of Modern Physics* **86**, 253–305 (Mar. 2014).
55. Becke, A. D. A new mixing of Hartree–Fock and local density-functional theories. *The Journal of Chemical Physics* **98**, 1372–1377 (Jan. 1993).
56. Perdew, J. P., Ernzerhof, M. & Burke, K. Rationale for mixing exact exchange with density functional approximations. *The Journal of Chemical Physics* **105**, 9982–9985 (Dec. 1996).
57. Aryasetiawan, F. & Gunnarsson, O. TheGWmethod. *Reports on Progress in Physics* **61**, 237–312 (Mar. 1998).

58. Heyd, J., Scuseria, G. E. & Ernzerhof, M. Hybrid functionals based on a screened Coulomb potential. *The Journal of Chemical Physics* **118**, 8207–8215 (May 2003).
59. Krukau, A. V., Vydrov, O. A., Izmaylov, A. F. & Scuseria, G. E. Influence of the exchange screening parameter on the performance of screened hybrid functionals. *The Journal of Chemical Physics* **125**, 224106 (Dec. 2006).
60. Freitas, L. C. G. Prêmio Nobel de Química em 1998: Walter Kohn e John A. Pople. *Química Nova* **22**, 293–298 (Apr. 1999).
61. Murphy, K. *Machine learning : a probabilistic perspective* ISBN: 9780262018029 (MIT Press, Cambridge, Mass, 2012).
62. Guido, S. *Introduction to Machine Learning with Python* 400 pp. ISBN: 1449369413. <https://www.ebook.de/de/product/23308778/sarah_guido_introduction_to_machine_learning_with_python.html> (O'Reilly UK Ltd., 2016).
63. Caruana, R. & Niculescu-Mizil, A. *An empirical comparison of supervised learning algorithms in Proceedings of the 23rd international conference on Machine learning - ICML '06* (ACM Press, 2006). doi:[10.1145/1143844.1143865](https://doi.org/10.1145/1143844.1143865).
64. Friedman, J., Hastie, T. & Tibshirani, R. Additive logistic regression: a statistical view of boosting (With discussion and a rejoinder by the authors). *The Annals of Statistics* **28**, 337–407 (Apr. 2000).
65. Friedman, J. H. Greedy function approximation: A gradient boosting machine. *Ann. Statist.* **29**, 1189–1232 (Oct. 2001).
66. James, G., Witten, D., Hastie, T. & Tibshirani, R. *An Introduction to Statistical Learning* ISBN: 1461471370. <https://www.ebook.de/de/product/20292548/gareth_james_daniela_witten_trevor_hastie_robert_tibshirani_an_introduction_to_statistical_learning.html> (Springer-Verlag GmbH, 2017).
67. Kresse, G. & Furthmüller, J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B* **54**, 11169–11186 (Oct. 1996).
68. Battle, R. & Benson, E. Bridging the semantic Web and Web 2.0 with Representational State Transfer (REST). *Journal of Web Semantics* **6**, 61–69 (Feb. 2008).
69. Kirklin, S. *et al.* The Open Quantum Materials Database (OQMD): assessing the accuracy of DFT formation energies. *npj Computational Materials* **1**. doi:[10.1038/npjcompumats.2015.10](https://doi.org/10.1038/npjcompumats.2015.10) (Dec. 2015).

70. Saal, J. E., Kirklin, S., Aykol, M., Meredig, B. & Wolverton, C. Materials Design and Discovery with High-Throughput Density Functional Theory: The Open Quantum Materials Database (OQMD). *JOM* **65**, 1501–1509 (Sept. 2013).
71. Jain, A. *et al.* Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials* **1**, 011002 (July 2013).
72. Ong, S. P. *et al.* The Materials Application Programming Interface (API): A simple, flexible and efficient API for materials data based on Representational State Transfer (REST) principles. *Computational Materials Science* **97**, 209–215 (Feb. 2015).
73. Levin, I. *NIST Inorganic Crystal Structure Database (ICSD)* en. 2020. doi:[10.18434/M32147](https://doi.org/10.18434/M32147).
74. Landis, D. D. *et al.* The Computational Materials Repository. *Computing in Science & Engineering* **14**, 51–57 (Nov. 2012).
75. Larsen, A. H. *et al.* The atomic simulation environment—a Python library for working with atoms. *Journal of Physics: Condensed Matter* **29**, 273002 (June 2017).
76. Ong, S. P. *et al.* Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science* **68**, 314–319 (Feb. 2013).
77. Ward, L. *et al.* Matminer: An open source toolkit for materials data mining. *Computational Materials Science* **152**, 60–69 (Sept. 2018).
78. Draxl, C. & Scheffler, M. The NOMAD laboratory: from data sharing to artificial intelligence. *Journal of Physics: Materials* **2**, 036001 (May 2019).
79. Wang, L., Maxisch, T. & Ceder, G. Oxidation energies of transition metal oxides within the GGA+U framework. *Physical Review B* **73**. doi:[10.1103/physrevb.73.195107](https://doi.org/10.1103/physrevb.73.195107) (May 2006).
80. Curtarolo, S. *et al.* AFLOWLIB.ORG: A distributed materials properties repository from high-throughput ab initio calculations. *Computational Materials Science* **58**, 227–235 (June 2012).
81. Curtarolo, S. *et al.* AFLOW: An automatic framework for high-throughput materials discovery. *Computational Materials Science* **58**, 218–226 (June 2012).
82. Calderon, C. E. *et al.* The AFLOW standard for high-throughput materials science calculations. *Computational Materials Science* **108**, 233–238 (Oct. 2015).

83. Setyawan, W. & Curtarolo, S. High-throughput electronic band structure calculations: Challenges and tools. *Computational Materials Science* **49**, 299–312 (Aug. 2010).
84. Setyawan, W., Gaume, R. M., Lam, S., Feigelson, R. S. & Curtarolo, S. High-Throughput Combinatorial Database of Electronic Band Structures for Inorganic Scintillator Materials. *ACS Combinatorial Science* **13**, 382–390 (June 2011).
85. Isayev, O. *et al.* Universal fragment descriptors for predicting properties of inorganic crystals. *Nature Communications* **8**. doi:[10.1038/ncomms15679](https://doi.org/10.1038/ncomms15679) (June 2017).
86. Ferrenti, A. M., de Leon, N. P., Thompson, J. D. & Cava, R. J. Identifying candidate hosts for quantum defects via data mining. *npj Computational Materials* **6**. doi:[10.1038/s41524-020-00391-7](https://doi.org/10.1038/s41524-020-00391-7) (Aug. 2020).
87. Stevanović, V., Lany, S., Zhang, X. & Zunger, A. Correcting density functional theory for accurate predictions of compound enthalpies of formation: Fitted elemental-phase reference energies. *Physical Review B* **85**. doi:[10.1103/physrevb.85.115104](https://doi.org/10.1103/physrevb.85.115104) (Mar. 2012).
88. Choudhary, K. *et al.* JARVIS: An Integrated Infrastructure for Data-driven Materials Design. arXiv: [2007.01831v1](https://arxiv.org/abs/2007.01831) [[cond-mat.mtrl-sci](#)] (July 3, 2020).
89. Thonhauser, T. *et al.* Van der Waals density functional: Self-consistent potential and the nature of the van der Waals bond. *Physical Review B* **76**. doi:[10.1103/physrevb.76.125112](https://doi.org/10.1103/physrevb.76.125112) (Sept. 2007).
90. Klimeš, J., Bowler, D. R. & Michaelides, A. Van der Waals density functionals applied to solids. *Physical Review B* **83**. doi:[10.1103/physrevb.83.195131](https://doi.org/10.1103/physrevb.83.195131) (May 2011).
91. Choudhary, K., Cheon, G., Reed, E. & Tavazza, F. Elastic properties of bulk and low-dimensional materials using van der Waals density functional. *Physical Review B* **98**. doi:[10.1103/physrevb.98.014107](https://doi.org/10.1103/physrevb.98.014107) (July 2018).
92. Choudhary, K. *et al.* Computational screening of high-performance optoelectronic materials using OptB88vdW and TB-mBJ formalisms. *Scientific Data* **5**. doi:[10.1038/sdata.2018.82](https://doi.org/10.1038/sdata.2018.82) (May 2018).
93. Rosenbrock, C. W. A Practical Python API for Querying AFLOWLIB. arXiv: [1710.00813v1](https://arxiv.org/abs/1710.00813) [[cs.DB](#)] (Sept. 28, 2017).
94. Breuck, P.-P. D., Evans, M. L. & Rignanese, G.-M. Robust model benchmarking and bias-imbalance in data-driven materials science: a case study on MODNet. arXiv: [2102.02263v1](https://arxiv.org/abs/2102.02263) [[cond-mat.mtrl-sci](#)] (Feb. 3, 2021).

95. Balachandran, P. V. *et al.* Predictions of new ABO_3 perovskite compounds by combining machine learning and density functional theory. *Physical Review Materials* **2**. doi:[10.1103/physrevmaterials.2.043802](https://doi.org/10.1103/physrevmaterials.2.043802) (Apr. 2018).
96. Shannon, R. D. Revised effective ionic radii and systematic studies of interatomic distances in halides and chalcogenides. *Acta Crystallographica Section A* **32**, 751–767 (Sept. 1976).
97. Zhang, H., Li, N., Li, K. & Xue, D. Structural stability and formability of ABO_3 -type perovskite compounds. *Acta Crystallographica Section B Structural Science* **63**, 812–818 (Nov. 2007).
98. Villars, P., Cenzual, K., Daams, J., Chen, Y. & Iwata, S. Data-driven atomic environment prediction for binaries using the Mendeleev number. *Journal of Alloys and Compounds* **367**, 167–175 (Mar. 2004).
99. Kotochigova, S., Levine, Z. H., Shirley, E. L., Stiles, M. D. & Clark, C. W. Local-density-functional calculations of the energy of atoms. *Physical Review A* **55**, 191–199 (Jan. 1997).
100. Laws, K. J., Miracle, D. B. & Ferry, M. A predictive structural model for bulk metallic glasses. *Nature Communications* **6**. doi:[10.1038/ncomms9123](https://doi.org/10.1038/ncomms9123) (Sept. 2015).
101. Butler, M. A. & Ginley, D. S. Prediction of Flatband Potentials at Semiconductor-Electrolyte Interfaces from Atomic Electronegativities. *Journal of The Electrochemical Society* **125**, 228–232 (Feb. 1978).
102. Ward, L., Agrawal, A., Choudhary, A. & Wolverton, C. A general-purpose machine learning framework for predicting properties of inorganic materials. *npj Computational Materials* **2**. doi:[10.1038/npjcompumats.2016.28](https://doi.org/10.1038/npjcompumats.2016.28) (Aug. 2016).
103. Deml, A. M., O’Hayre, R., Wolverton, C. & Stevanović, V. Predicting density functional theory total energies and enthalpies of formation of metal-nonmetal compounds by linear regression. *Physical Review B* **93**. doi:[10.1103/physrevb.93.085142](https://doi.org/10.1103/physrevb.93.085142) (Feb. 2016).
104. Weeber, A. W. Application of the Miedema model to formation enthalpies and crystallisation temperatures of amorphous alloys. *Journal of Physics F: Metal Physics* **17**, 809–813 (Apr. 1987).
105. Yang, X. & Zhang, Y. Prediction of high-entropy stabilized solid-solution in multi-component alloys. *Materials Chemistry and Physics* **132**, 233–238 (Feb. 2012).

106. Rupp, M., Tkatchenko, A., Müller, K.-R. & von Lilienfeld, O. A. Fast and Accurate Modeling of Molecular Atomization Energies with Machine Learning. *Physical Review Letters* **108**. doi:[10.1103/physrevlett.108.058301](https://doi.org/10.1103/physrevlett.108.058301) (Jan. 2012).
107. Schütt, K. T. *et al.* How to represent crystal structures for machine learning: Towards fast prediction of electronic properties. *Physical Review B* **89**. doi:[10.1103/physrevb.89.205118](https://doi.org/10.1103/physrevb.89.205118) (May 2014).
108. Faber, F., Lindmaa, A., von Lilienfeld, O. A. & Armiento, R. Crystal structure representations for machine learning models of formation energies. *International Journal of Quantum Chemistry* **115**, 1094–1101 (Apr. 2015).
109. Ewald, P. P. Die Berechnung optischer und elektrostatischer Gitterpotentiale. *Annalen der Physik* **369**, 253–287 (1921).
110. Hansen, K. *et al.* Machine Learning Predictions of Molecular Properties: Accurate Many-Body Potentials and Nonlocality in Chemical Space. *The Journal of Physical Chemistry Letters* **6**, 2326–2331 (June 2015).
111. Ward, L. *et al.* Including crystal structure attributes in machine learning models of formation energies via Voronoi tessellations. *Physical Review B* **96**. doi:[10.1103/physrevb.96.024104](https://doi.org/10.1103/physrevb.96.024104) (July 2017).
112. Botu, V. & Ramprasad, R. Adaptive machine learning framework to accelerate ab initio molecular dynamics. *International Journal of Quantum Chemistry* **115**, 1074–1083 (Dec. 2014).
113. De Jong, M. *et al.* A Statistical Learning Framework for Materials Science: Application to Elastic Moduli of k-nary Inorganic Polycrystalline Compounds. *Scientific Reports* **6**. doi:[10.1038/srep34256](https://doi.org/10.1038/srep34256) (Oct. 2016).
114. Seko, A., Hayashi, H., Nakayama, K., Takahashi, A. & Tanaka, I. Representation of compounds for machine-learning prediction of physical properties. *Physical Review B* **95**. doi:[10.1103/physrevb.95.144110](https://doi.org/10.1103/physrevb.95.144110) (Apr. 2017).
115. Steinhardt, P. J., Nelson, D. R. & Ronchetti, M. Bond-orientational order in liquids and glasses. *Physical Review B* **28**, 784–805 (July 1983).
116. Waroquiers, D. *et al.* Statistical Analysis of Coordination Environments in Oxides. *Chemistry of Materials* **29**, 8346–8360 (Sept. 2017).
117. Zimmermann, N. E. R., Horton, M. K., Jain, A. & Haranczyk, M. Assessing Local Structure Motifs Using Order Parameters for Motif Recognition, Interstitial Identification, and Diffusion Path Characterization. *Frontiers in Materials* **4**. doi:[10.3389/fmats.2017.00034](https://doi.org/10.3389/fmats.2017.00034) (Nov. 2017).

118. Behler, J. Atom-centered symmetry functions for constructing high-dimensional neural network potentials. *The Journal of Chemical Physics* **134**, 074106 (Feb. 2011).
119. Khorshidi, A. & Peterson, A. A. Amp: A modular approach to machine learning in atomistic simulations. *Computer Physics Communications* **207**, 310–324 (Oct. 2016).
120. Peng, H. L., Li, M. Z. & Wang, W. H. Structural Signature of Plastic Deformation in Metallic Glasses. *Physical Review Letters* **106**. doi:[10.1103/physrevlett.106.135503](https://doi.org/10.1103/physrevlett.106.135503) (Mar. 2011).
121. Wang, Q. & Jain, A. A transferable machine-learning framework linking interstice distribution and plastic heterogeneity in metallic glasses. *Nature Communications* **10**. doi:[10.1038/s41467-019-13511-9](https://doi.org/10.1038/s41467-019-13511-9) (Dec. 2019).
122. Dylla, M. T., Dunn, A., Anand, S., Jain, A. & Snyder, G. J. Machine Learning Chemical Guidelines for Engineering Electronic Structures in Half-Heusler Thermoelectric Materials. *Research* **2020**, 1–8 (Apr. 2020).
123. None Available. *Materials Data on PH6C2S2N(ClO2)2 by Materials Project* en. 2020. doi:[10.17188/1268877](https://doi.org/10.17188/1268877).
124. None Available. *Materials Data on Nb7S2I19 by Materials Project* en. 2014. doi:[10.17188/1277059](https://doi.org/10.17188/1277059).