



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

Bachelorthesis

Title

vorgelegt von

TJARK SIEVERS

Fakultät: Mathematik, Informatik und Naturwissenschaften

Fachbereich: Physik

Studiengang: Physik

Matrikelnummer: 7147558

Erstgutachter:

Zweitgutachter:



# Kurzzusammenfassung



# Inhaltsverzeichnis

<b>Kurzzusammenfassung</b>	<b>i</b>
<b>Danksagung</b>	<b>vii</b>
<b>Motivation</b>	<b>ix</b>
<b>I Many-body physics</b>	<b>1</b>
I.1 The electronic structure problem . . . . .	1
I.2 Density Functional Theory . . . . .	1
I.2.1 Hohenberg-Kohn theorems . . . . .	1
I.2.2 Kohn-Sham equations . . . . .	2
<b>II Computational Basics</b>	<b>3</b>
II.1 Parallel computing . . . . .	3
II.1.1 On scalability . . . . .	3
II.2 QUANTUM ESPRESSO . . . . .	4
<b>III Parallelisation of self-consistent calculations of electronic-structure properties</b>	<b>5</b>
III.0.1 First scaling tests . . . . .	5
III.0.2 Testing different compilers and mathematical libraries . . . . .	6
III.0.3 Using the parallelisation parameters of QUANTUM ESPRESSO . . . . .	6
III.0.4 Comparison with calculations on the HLRN cluster . . . . .	8
III.0.5 Conclusion: Parameters for optimal scaling . . . . .	8
<b>IV Parallelization of DFPT calculations</b>	<b>11</b>
<b>Literaturverzeichnis</b>	<b>13</b>
<b>Listings</b>	<b>15</b>
Abbildungsverzeichnis . . . . .	15
Tabellenverzeichnis . . . . .	15



# Akronyme

**DFT** Density Functional Theory. [1](#)





# Danksagung



# Motivation

## Conventions

Throughout the text of this thesis scalars are written in italic  $s$ , vectors in bold italic  $\mathbf{v}$  and matrices in bold  $\mathbf{M}$  fonts. The summation/multiplication over nearest neighbour sites  $i$  and  $j$  is  $\langle ij \rangle$  as a subscript to  $\sum/\prod$ . Furthermore, Hartree atomic units are used in general and only in selected instances it is deviated from this:  $\hbar = m_e = e = 4\pi/\varepsilon_0 = 1$ .



# I Many-body physics

## I.1 The electronic structure problem

In solid state physics, one general problem we are concerned with is finding the properties of the ground state of an isolated system of  $N$  interacting electrons in an external potential. The system is described by the Schrödinger equation

$$\hat{H}\Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) = E\Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) \quad (\text{I.1})$$

with the Hamiltonian

$$\hat{H} = \hat{T}_e + \hat{V}_{n-e} + \hat{V}_{e-e} + \hat{V}_{n-n} \quad (\text{I.2})$$

$$= -\sum_i \frac{1}{2} \nabla_i^2 - \sum_i \sum_{\alpha} \frac{Z_{\alpha}}{|\mathbf{r}_i - \mathbf{R}_{\alpha}|} + \frac{1}{2} \sum_{i \neq j} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} + \frac{1}{2} \sum_{\alpha \beta} \frac{Z_{\alpha} Z_{\beta}}{R_{\alpha \beta}} \quad (\text{I.3})$$

where:

- $\hat{T}_e$  is the kinetic energy of the electrons
- $\hat{V}_{n-e}$  is the potential energy of the electrons in the field of the nuclei
- $\hat{V}_{e-e}$  is the Coulomb interaction between the electrons and
- $\hat{V}_{n-n}$  is the Coulomb interaction between the nuclei

## I.2 Density Functional Theory

A direct solution to the electronic structure problem, this meaning obtaining the ground-state many-body wavefunction  $\Psi(\mathbf{r}_1, \dots, \mathbf{r}_N)$  for a given potential is analytically impossible even for a small number of electrons. As such, the need for good approximations to obtain results for real world systems is high. One particularly successful approach is [Density Functional Theory \(DFT\)](#). In the following section, the theoretical framework of

### I.2.1 Hohenberg-Kohn theorems

The basis for DFT lies in the exact reformulation of the outlined electronic structure problem by Hohenberg and Kohn [1]. This reformulation uses the ground state density of the electronic system as the basic variable.

Hohenberg-Kohn theorems:

stuff missing

- I The external potential (and by extension the ground state wave function and the ground state energy) are unique functionals of the ground state density (except for an additive constant).
- II The ground state energy minimizes the energy functional,

$$E[n(r)] > E_0 \quad \forall n(r) \neq n_0(r)$$

.

### **I.2.2 Kohn-Sham equations**

One way of approximating the functional  $F[n]$  was given by Kohn and Sham [2]. The idea is to use a non-interacting auxiliary system of electrons and introduce a correction potential.

## II Computational Basics

### II.1 Parallel computing

The following section will give an overview of the technical aspects of running computer code (such as QUANTUM ESPRESSO ) on massively parallel computing environments (such as the PHYSnet compute cluster). The information presented can be found in any textbook on parallel or high-performance computing [3].

#### II.1.1 On scalability

In scientific computing, one can identify two distinct reasons to distribute workloads to multiple processors:

- The execution time on a single core is not sufficient. The definition of sufficient is dependent on the specific task and can range from „over lunch“ to „multiple weeks“
- The memory requirements grow outside the capabilities of a single core

In order to judge how well a task can be parallelized, usually some sort of scalability metric is employed, for example:

- How fast can a problem be solved with  $N$  processors instead of one?
- What kind of bigger problem (finer resolution, more particles, etc.) can be solved with  $N$  processors?
- How much of the resources is used for solving the problem?

The speedup by using  $N$  workers to solve a problem instead of one is defined as  $S = \frac{T_1}{T_N}$ , where  $T_1$  is the execution time on a single processor and  $T_N$  is the execution time on  $N$  processors. In the ideal case, where all the work can be perfectly distributed among the processors, all processors need the same time for their respective workloads and don't have to wait for others processors to finish their workload to continue, the execution time on  $N$  processors would be  $\frac{T_1}{N}$ , so the speedup would be  $S = \frac{T_1}{\frac{T_1}{N}} = N$ .

In reality, there are many factors either limiting or in some cases supporting parallel code scalability. Limiting factors include:

- *Algorithmic limitations*: when parts of a calculation are mutually dependent on each other, the calculation cannot be fully parallelized
- *Bottlenecks*: in any computer system exist resources which are shared between processor cores with limitations on parallel access. This serializes the execution by requiring cores to wait for others to complete the task which uses the shared resources in question

- *Startup Overhead*: introducing parallelization into a program necessarily introduces an overhead, e.g. for distributing data across all the processors
- *Communication*: often solving a problem requires communication between different cores (e.g. exchange of interim results after a step of the calculation). Communication can be implemented very effectively, but can still introduce a big prize in computation time

On the other hand, faster parallel code execution can come from:

- *Better caching*: when the data the program is working with is distributed among processors (assuming constant problem size), it may enable the data to be stored in faster cache memory. Modern computers typically have three layers of cache memory, with level 1 cache being the smallest and fastest and level 3 being the largest and slowest, so smaller data chunks per processor can lead to the data not being stored in main memory, but completely in cache or in a faster cache level

## II.2 Quantum ESPRESSO



# III Parallelisation of self-consistent calculations of electronic-structure properties

## III.0.1 First scaling tests

The first step in analysing the scaling of QUANTUM ESPRESSO is to perform a baseline scaling test without any optimisations yet applied. In Fig. ?? two scaling tests on the earlier mentioned benchmarking systems Si and TaS<sub>2</sub> are pictured. The QUANTUM ESPRESSO version used is compiled using OpenMPI 4.1.0 , without any further compilation or runtime optimisation parameters used.

What exactly?  
Which compiler?

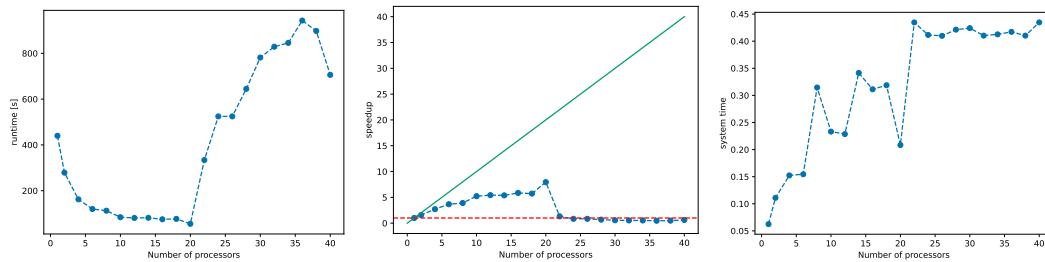


Abbildung III.1: XXX

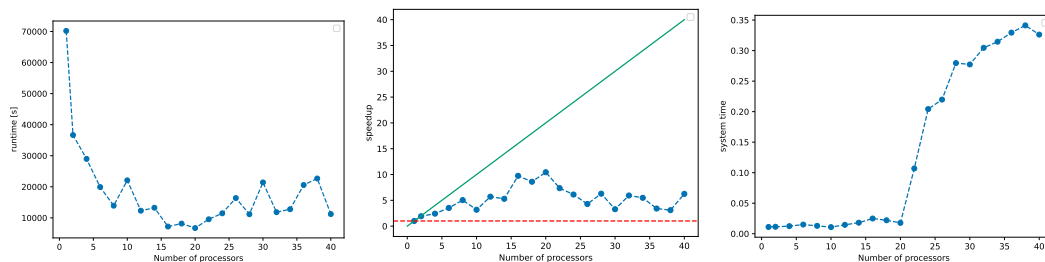


Abbildung III.2: XXX

caption

Three different metrics of scalability are pictured in ??.

- runtime: absolute runtime of the compute job
- speedup: runtime divided by runtime of the job on a single core

- system time: percentage of wall time used by system tasks, e.g. writing to disk, etc. (calculated as  $\text{walltime} - \text{cputime}$ )

For further analysis mainly speedup will be used as a metric of scalability, because it lends itself to easy interpretation: optimal scalability is achieved when the speedup scales linearly with the number of processors (with a slope of one), as discussed in ch. II.1.

On a single node, both the Si and TaS<sub>2</sub> calculations show good, but not perfect scaling behavior: the speedup does approximately scale linearly with the number of processors, but the slope is round about  $\frac{1}{2}$ . Even though the scaling behavior is not perfect, there is just a small, almost constant amount of runtime used by system calls, this speaks for good parallelisation: as discussed in ch. II.1, startup time is part of every

missing

When using more than one node, not only does the scaling get worse, the execution needs longer than on a single core for the Si system, with a marginally better performance for the TaS<sub>2</sub> system. This is also seen in the plots of system time. The percentage of time used for tasks not directly related to calculations (mostly exchange of data in this case, which induces long waiting time when the connection between processors is not as fast as on one motherboard) goes from a near constant value for under 20 processors to 50% of the execution time for the Si system and 35% for the TaS<sub>2</sub> system.

These scaling tests pose now two questions to be answered:

- Is better scaling on a single node possible?
- How can scaling over more than one node be achieved?

### III.0.2 Testing different compilers and mathematical libraries

A first strategy for solving issues with parallelization is trying different compilers and mathematical libraries. In the PHYSnet cluster a variety of software packages is available. For the compilation of QUANTUM ESPRESSO a

missing

For testing QUANTUM ESPRESSO will be compiled using the following software combinations:

- OpenMPI 4.1.0 and OpenBLAS
- OpenMPI 4.1.0 and Scalapack
- Intel oneAPI 2021.4 (includes Intel MPI, Fortran and C compilers as well as Intel MKL, a scalable mathematical library)

caption

### III.0.3 Using the parallelisation parameters of Quantum ESPRESSO

As detailed in section II.2, QUANTUM ESPRESSO offers ways to manage how the workload is distributed among the processors. In `pw.x` the default plane wave parallelization, k-point-parallelization and linear-algebra parallelization are implemented.

The benchmark pictured in III.4 is set up as follows: for a given number of processors  $N_p$ , the parameter  $N_k$  splits the  $N_p$  processors into  $N_k$  processors pools. As the number of processors in one pool has to be a whole number, only certain combinations of  $N_p$  and  $N_k$  are possible, for example  $N_p = 32$  could be split into processor pools of size 2 with  $N_k = 16$ , size 8 with  $N_k = 4$  or size 16 with  $N_k = 2$ . This leads to choosing the size of the processor pools as a variable,

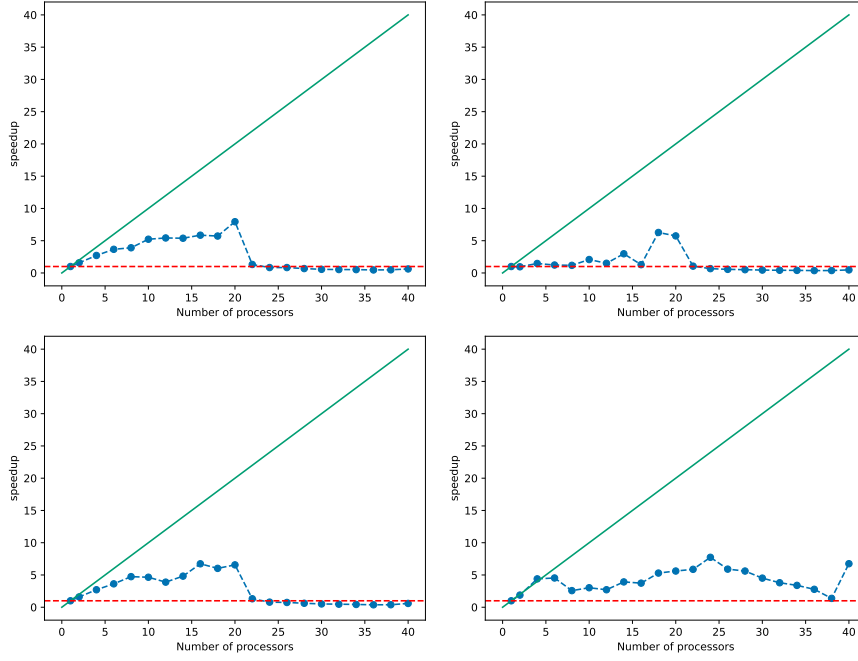


Abbildung III.3: XXX

not the parameter `nk`. Fig. III.4 shows the scaling for pool sizes 2, 8 and 16 for QUANTUM ESPRESSO being compiled with OpenMPI/Scalapack and Intel oneAPI. This choice of pool sizes showcases the smallest pool size possibly (namely 2), as well as a bigger pool size with 16, that still gives rise to a few data points over the chosen range of processors.

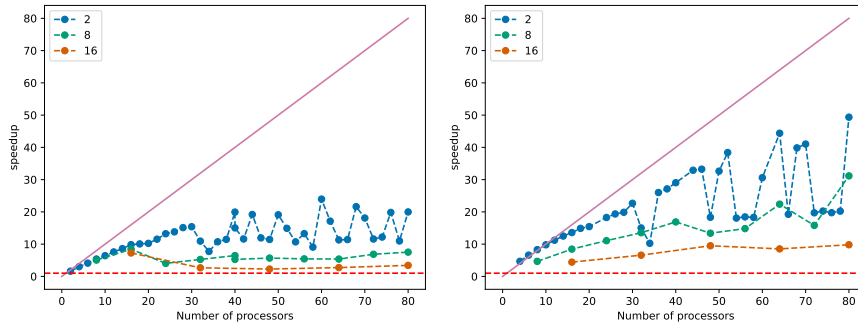


Abbildung III.4: XXX

Fig. III.4 shows

The same scaling test is applied to the TaS2 system in fig. III.5, with the same list of pool sizes, but over a wider range of processors.

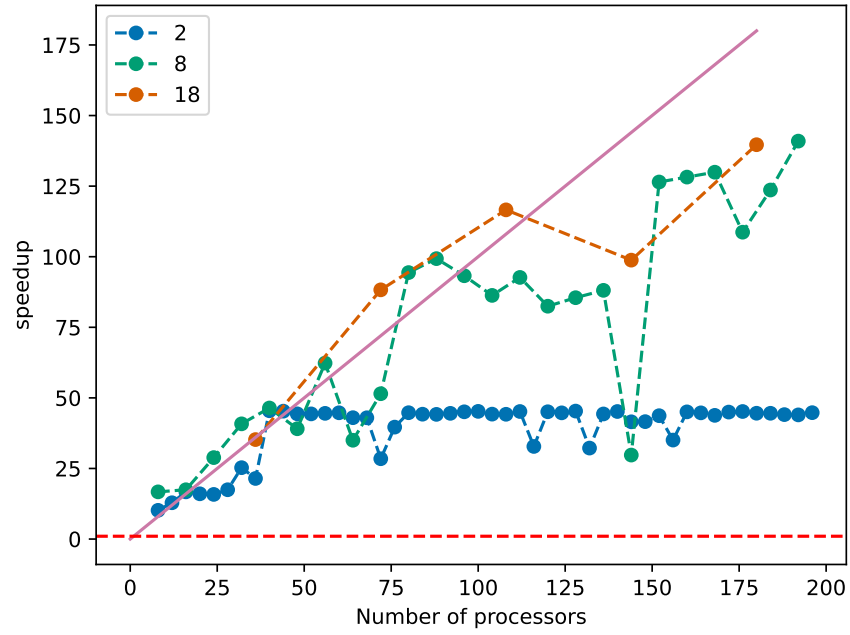


Abbildung III.5: XXX

Remarkably, the scaling behavior is swapped in comparison to III.4, as the pool size 2 saturates fast and the bigger pool sizes showing way better scaling behavior.

It can also be instructive to look at the idle time for this benchmark to judge the quality of parallelization.

Fig. III.6 shows a distribution of idle times between 1 % and 4 % of the whole wall time, without any kind of systemic increase over any range of processors. This means the parallelization is as good as possible for these

#### III.0.4 Comparison with calculations on the HLRN cluster

#### III.0.5 Conclusion: Parameters for optimal scaling

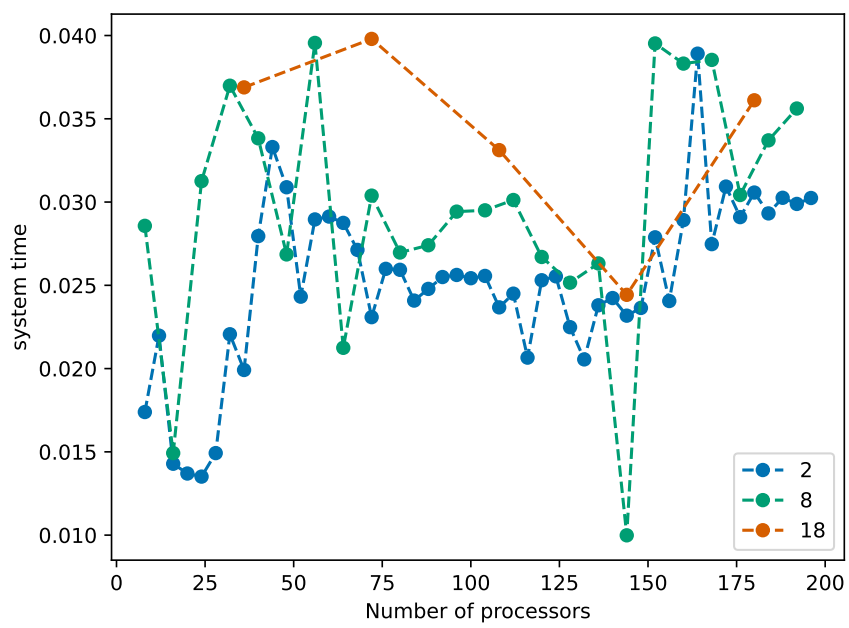


Abbildung III.6: XXX



## IV Parallelization of DFPT calculations





# Literatur

- [1] P. Hohenberg und W. Kohn. „Inhomogeneous Electron Gas“. In: *Phys. Rev.* 136.3 (Nov. 1964). Publisher: American Physical Society, B864–B871. DOI: [10.1103/PhysRev.136.B864](https://doi.org/10.1103/PhysRev.136.B864).
- [2] W. Kohn und L. J. Sham. „Self-Consistent Equations Including Exchange and Correlation Effects“. In: *Phys. Rev.* 140.4 (Nov. 1965). Publisher: American Physical Society, A1133–A1138. DOI: [10.1103/PhysRev.140.A1133](https://doi.org/10.1103/PhysRev.140.A1133).
- [3] G. Hager und G. Wellein. *Introduction to High Performance Computing for Scientists and Engineers*. 0. Aufl. CRC Press, 2. Juli 2010. ISBN: 978-1-4398-1193-1. DOI: [10.1201/EBK1439811924](https://doi.org/10.1201/EBK1439811924).



# Listings

## Abbildungsverzeichnis

III.1 XXX . . . . .	5
III.2 XXX . . . . .	5
III.3 XXX . . . . .	7
III.4 XXX . . . . .	7
III.5 XXX . . . . .	8
III.6 XXX . . . . .	9

## Tabellenverzeichnis