

Master Thesis

Performance Modelling and Analysis of the openQxD Lattice QCD Application

Roman Gruber

ETH Zürich, TODO:date, TODO: supervisor(s)

Abstract

TODO

This work is licensed under a [Creative Commons](#) “Attribution-ShareAlike 4.0 International” license.



Contents

1	Introduction	2
2	Conventions	2
3	Non-Abelian gauge theories	2
4	Lattice Gauge Theories	6
4.1	Euclidean theory	7
4.2	Lattice Discretisation	10
5	Performance Models	11
6	Software: openQxD	11
7	Real number formats	11
7.1	IEEE Standard for Floating-Point Arithmetic	11
7.2	Posits	14
7.3	Floating point numbers in openQxD	16
8	Conjugate Gradient algorithm	19
8.1	Derivation	19
8.2	CG kernel in openQxD	26
8.3	Simulating CG with different datatypes	28
8.3.1	Discussion of figures 8 - 11	28
8.3.2	8 ⁴ lattice	33
8.3.3	Conclusion	33
9	SAP preconditioned GCR algorithm	39
9.1	Even-Odd Preconditioning	39
9.2	Schwarz Alternating Procedure	40
9.3	SAP as a Preconditioner	42
9.4	Generalised Conjugate Residual algorithm	43
9.5	GCR in openQxD	45

9.6	Simulating SAP_GCR	47
9.6.1	Discussion of figures 18 - 21	48
9.6.2	Conclusion	51
10	Deflated SAP preconditioned GCR algorithm	51
10.1	Deflation	52
10.2	Choosing the deflation subspace	56
11	Multi-shift Conjugate Gradient algorithm	56
12	Dirac operator	56
13	GPU Implementation	58
14	Algorithm-independent considerations	59
15	Summary	60
16	Future	60
17	References	61
	Appendices	63
A	Proofs	63
B	Code	63
C	List of Proposals	63
	Acronyms	64
	Glossary	64

1 Introduction

TODO

Proposal 1.1: Example proposal

Reference here with pp:one.

In QCD blabla see proposal 1.1. orange, yellow, blue, brown, pink, red, green, purple, turquoise, lightblue, lightgreen, lightpink, darkblue, lightblue, lightpink, lightgreen, linkcolor

The result of integrating $\int \sqrt{1+x} \, dx$ is given by $\frac{2(x+1)^{\frac{3}{2}}}{3}$
Python says "Hello!"

2 Conventions

3 Non-Abelian gauge theories

The goal of this section is to derive a Lagrangian descibing fundamental fermions that is Lorentz-invariant (coming from special relativity) as well as invariant under local phase-transformations (coming from basic quantum mechanics).

Let's consider a set of N complex independent Dirac spinors

$$\psi(x) = \begin{pmatrix} \psi_1(x) \\ \vdots \\ \psi_N(x) \end{pmatrix},$$

where every $\psi_a(x)$, $a \in \{1, \dots, N\}$, has 4 components accessed via spinor index $\alpha \in \{1, \dots, 4\}$. We have a set of complex Grassmann-valued fields $\psi_a^\alpha(x)$, a is called the **color index**, α is called the **spinor index**. It makes sense to demand the theory to be invariant under $SU(N)$ -transformations, because basic quantum mechanics tells us that these phases are unobservable. Such a field theory that is invariant under the gauge group $SU(N)$ has to introduce (massless) vectorial fields $A_\mu(x)$, called the **gauge fields**. From special relativity, we also demand the theory to be Lorentz-invariant. The goal is to construct a Lagrangian for fermionic fields that satisfies these symmetries.

Let's introduce the local¹ gauge transformation under which we want the theory to be invariant

$$\psi(x) \longrightarrow V(x)\psi(x),$$

with $V(x) \in SU(N)$ unitary. Since we want our fermions to possess a non-zero mass, the Lagrangian will contain a term quadratic in the field ψ . A first approach of a $SU(N)$ -invariant expression would be a term proportional to $\psi^\dagger\psi$. Unfortunately ψ is Grassmann valued, thus anti-commutes with itself, meaning the mentioned term is equal to zero². Special relativity also demands the term to be a Lorentz-scalar, which $\psi^\dagger\psi$ is not. Under Lorentz-transformation Λ , ψ and $\psi(x)^\dagger$ transform as

$$\begin{aligned}\psi(x) &\longrightarrow \Lambda\psi(x), \\ \psi(x)^\dagger &\longrightarrow \psi(x)^\dagger\Lambda^\dagger.\end{aligned}$$

We can use a property of the γ -matrices here. The γ -matrices are defined such that they obey the **Clifford-algebra**

$$\{\gamma^\mu, \gamma^\nu\} = 2\eta^{\mu\nu} \cdot id, \quad (3.1)$$

with $\mu, \nu \in \{0, 1, \dots, D-1\}$, where D is the spacetime dimension and id is the identity operator in spinor space. The needed property of the γ -matrices is

$$\Lambda^\dagger\gamma^\mu = \gamma^\mu\Lambda^{-1}.$$

With this, it makes sense to define the **Dirac-adjoint** as $\bar{\psi} := \psi^\dagger\gamma^0$ and construct a Lorentz- and $SU(N)$ -invariant quadratic expression $c\bar{\psi}\psi$, with $c \in \mathbb{C}$.

The Lagrangian of the theory inevitably contains derivatives of ψ and since the transformation is local (different for every spacetime point x), we have to redefine a derivative that compensates for this. The **directional derivative** along direction n is

$$n^\mu\partial_\mu\psi(x) = \lim_{\epsilon \rightarrow 0} \frac{\psi(x + n\epsilon) - \psi(x)}{\epsilon}.$$

The two fields appearing in this expression are evaluated at different spacetime points and thus transform differently under $V(x)$. In order for the kinetic expression in the Lagrangian $n^\mu\partial_\mu$ to be invariant under $SU(N)$, we introduce a compensator for the shifts in the derivative.

Definition 3.1 (Compensator field). *The **compensator field** $U(x, y)$ is a non-local matrix quantity that transforms under $V(x) \in SU(N)$ as*

$$U(x, y) \longrightarrow V(x)U(x, y)V^\dagger(y).$$

$U(x, y)$ is an element of $SU(N)$ for all x, y , with $U(x, x) = id$.

¹The phase depends on the spacetime coordinate x .

²To be pedantic, 0 is $SU(N)$ -invariant.

We then redefine the derivative as

$$n^\mu D_\mu \psi(x) = \lim_{\epsilon \rightarrow 0} \frac{\psi(x + n\epsilon) - U(x + n\epsilon, x)\psi(x)}{\epsilon}, \quad (3.2)$$

to compensate for the shift in x . Now, the derivative transforms as

$$n^\mu D_\mu \psi(x) \longrightarrow n^\mu V(x) D_\mu \psi(x).$$

We can Taylor-expand $U(x + n\epsilon, x)$ around $\epsilon = 0$ and find

$$\begin{aligned} U(x, x + n\epsilon) &\approx U(x, x) + \frac{1}{1!} \frac{\partial U}{\partial x^\mu} \frac{\partial(x^\mu + n^\mu \epsilon)}{\partial \epsilon} \Big|_{\epsilon=0} \cdot \epsilon + O(\epsilon^2) \\ &= id + \frac{\partial U}{\partial x^\mu} \Big|_{\epsilon=0} n^\mu \epsilon + O(\epsilon^2) \\ &= id - ig\epsilon n^\mu A_\mu^a(x) T^a + O(\epsilon^2), \end{aligned} \quad (3.3)$$

where we introduced real-valued bosonic vector-fields $A_\mu^a(x)$ ³, an arbitrary constant $g \in \mathbb{R}_{>0}$ and the generators $T^a \in su(N)$, the Lie-algebra of $SU(N)$. The fields $A_\mu^a(x)$ are defined by

$$\frac{\partial U}{\partial x^\mu} \Big|_{\epsilon=0} =: -igA_\mu^a(x)T^a.$$

Since $C_{x,n}(t) := U(x + nt, x)$ for every x and n is a curve in $SU(N)$ going through id at $t = 0$, its derivative evaluated at $t = 0$ is therefore an element of the Lie-algebra $su(N)$. Inserting this into (3.2)

$$\begin{aligned} n^\mu D_\mu \psi(x) &= \lim_{\epsilon \rightarrow 0} \frac{\psi(x + n\epsilon) - \psi(x) + ig\epsilon n^\mu A_\mu^a(x) T^a \psi(x)}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} \frac{\psi(x + n\epsilon) - \psi(x)}{\epsilon} + ign^\mu A_\mu^a(x) T^a \psi(x) \\ &= n^\mu (\partial_\mu + igA_\mu^a(x) T^a) \psi(x), \end{aligned}$$

which leads directly to the ***gauge covariant derivative***

$$D_\mu := \partial_\mu + igA_\mu^a(x) T^a. \quad (3.4)$$

Since we introduced a new field and by this a corresponding particle, in order for the particle to have a propagator, we also have to implement a kinetic term that is quadratic in A_μ^a or its derivatives. Obviously the kinetic term should be invariant under $SU(N)$ transformations. For this we need the plaquette⁴.

Definition 3.2 (Plaquette). *Let $n_1 \neq n_2$ be two 4-vectors and $\epsilon > 0$. The **plaquette** (see figure 1) in the (n_1, n_2) -subspace is defined as*

$$\hat{U}_{n_1, n_2}(\epsilon, x) := U(x, x + \epsilon n_2) U(x + \epsilon n_2, x + \epsilon n_2 + \epsilon n_1) U(x + \epsilon n_2 + \epsilon n_1, x + \epsilon n_1) U(x + \epsilon n_1, x). \quad (3.5)$$

³ $A_\mu^a(x)$ is an auxiliary field that is the infinitesimal limit of a compensator field, also called **connection**

⁴There are other (probably simpler and more elegant) methods to derive a kinetic term for the gauge fields, but the plaquette will arise again later in the context of lattice gauge theories. And as a physicist sometimes one has to go through some pain and suffer a little bit here and there - this is part of the game.

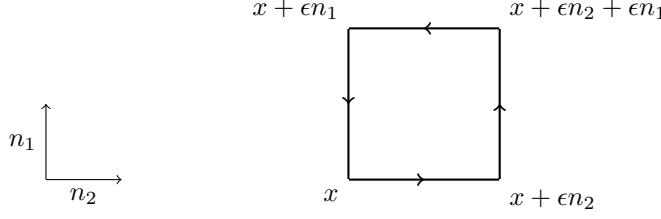


Figure 1: Scheme of the plaquette $\hat{U}_{n_1, n_2}(\epsilon, x)$ in the (n_1, n_2) -subspace.

The plaquette $\hat{U}_{n_1, n_2}(\epsilon, x)$ is not invariant under $SU(N)$, but its trace over color space $\text{tr}(\hat{U}_{n_1, n_2}(\epsilon, x))$ is, because of the cyclicity of the trace. Using equation (3.3), we can write

$$\begin{aligned} U(x, y) &= \exp\left(-ig(x-y)^\mu A_\mu^a\left(\frac{x+y}{2}\right)T^a + O((x-y)^2)\right) \\ U(x + \epsilon n_1, x + \epsilon n_2) &= \exp\left(-ig\epsilon(n_1 - n_2)^\mu A_\mu^a\left(x + \epsilon\frac{n_1 + n_2}{2}\right)T^a + O(\epsilon^2)\right). \end{aligned} \quad (3.6)$$

If ϵ is small we can expand

$$A_\mu^a(x + \epsilon n) = A_\mu^a(x) + \epsilon n^\nu \partial_\nu A_\mu^a(x) + O(\epsilon^2)$$

and insert it in equation (3.6),

$$\begin{aligned} U(x + \epsilon n_1, x + \epsilon n_2) &= \\ \exp\left(-ig\epsilon(n_1 - n_2)^\mu A_\mu^a(x)T^a - ig\frac{\epsilon^2}{2}(n_1 + n_2)^\nu \partial_\nu(n_1 - n_2)^\mu A_\mu^a(x)T^a + O(\epsilon^3)\right). \end{aligned}$$

Using this formula as well as the Baker–Campbell–Hausdorff formula⁵ we obtain

$$\begin{aligned} \hat{U}_{n_1, n_2}(\epsilon, x) &= \exp\left(+ig\epsilon n_2^\mu A_\mu^a(x)T^a + ig\frac{\epsilon^2}{2}n_2^\nu \partial_\nu n_2^\mu A_\mu^a(x)T^a \right. \\ &\quad + ig\epsilon n_1^\mu A_\mu^a(x)T^a + ig\frac{\epsilon^2}{2}(n_1 + 2n_2)^\nu \partial_\nu n_1^\mu A_\mu^a(x)T^a \\ &\quad - ig\epsilon n_2^\mu A_\mu^a(x)T^a - ig\frac{\epsilon^2}{2}(2n_1 + n_2)^\nu \partial_\nu n_2^\mu A_\mu^a(x)T^a \\ &\quad - ig\epsilon n_1^\mu A_\mu^a(x)T^a - ig\frac{\epsilon^2}{2}n_1^\nu \partial_\nu n_1^\mu A_\mu^a(x)T^a \\ &\quad - \frac{1}{2}g^2\epsilon^2 n_2^\mu A_\mu^a(x)n_1^\nu A_\nu^b(x)[T^a, T^b] + \frac{1}{2}g^2\epsilon^2 n_2^\mu A_\mu^a(x)n_1^\nu A_\nu^b(x)[T^a, T^b] \\ &\quad \left. + \frac{1}{2}g^2\epsilon^2 n_1^\mu A_\mu^a(x)n_2^\nu A_\nu^b(x)[T^a, T^b] - \frac{1}{2}g^2\epsilon^2 n_2^\mu A_\mu^a(x)n_1^\nu A_\nu^b(x)[T^a, T^b] + O(\epsilon^3) \right) \end{aligned}$$

By staring long enough at this expression, we see that there are a lot of terms cancelling each other; the first terms in the first 4 lines cancel exactly, the second terms in the first 4 lines all cancel

⁵In this case we used the BCH-formula on steroids,

$$e^{\epsilon A} e^{\epsilon B} e^{\epsilon C} e^{\epsilon D} = e^{\epsilon A + \epsilon B + \epsilon C + \epsilon D + \frac{\epsilon^2}{2}[A, B] + \frac{\epsilon^2}{2}[A, C] + \frac{\epsilon^2}{2}[A, D] + \frac{\epsilon^2}{2}[B, C] + \frac{\epsilon^2}{2}[B, D] + \frac{\epsilon^2}{2}[C, D] + O(\epsilon^3)}.$$

except the terms involving $2n_1$ and $2n_2$, and the two terms in the fifth line cancel, and the two terms in the sixth line are actually the same. Thus, we have 3 terms surviving corresponding to the yellow color. Using the directed fields and derivatives $n_i^\mu A_\mu^a(x) =: A_i^a(x)$ and $n_j^\nu \partial_\nu n_i^\mu A_\mu^a(x) =: \partial_j A_i^a(x)$, we end up in

$$\hat{U}_{n_1, n_2}(\epsilon, x) = \exp \left(-ig\epsilon^2 \underbrace{[\partial_1 A_2^a(x) - \partial_2 A_1^a(x) + gf^{abc} A_1^b(x) A_2^c(x)]}_{=: F_{12}^a} T^a + O(\epsilon^3) \right) \quad (3.7)$$

where we used the totally anti-symmetric **structure constants** of the Lie-algebra defined via the commutation relations $[T^a, T^c] = if^{abc} T^b$.

The expression in the square bracket of the last line is what we define as the **Yang-Mills field strength tensor**⁶

$$F_{\mu\nu}^a = \partial_\mu A_\nu^a(x) - \partial_\nu A_\mu^a(x) + gf^{abc} A_\mu^b(x) A_\nu^c(x). \quad (3.8)$$

The field strength tensor transforms in the same way as the plaquette $\hat{U}_{n_1, n_2}(\epsilon, x)$, thus the trace over $F_{\mu\nu}^a T^a$ in color space is invariant under $SU(N)$ transformations. In order to be Lorentz-invariant as well, we need a Lorentz-scalar. Thus we need to contract the indices and obtain⁷

$$tr(F_{\mu\nu}^a T^a F^{\mu\nu, b} T^b) = F_{\mu\nu}^a F^{\mu\nu, b} tr(T^a T^b) \quad (3.9)$$

$$= \frac{1}{2} F_{\mu\nu}^a F^{\mu\nu, a}. \quad (3.10)$$

The quadratic terms in the Lagrangian are always of the form $-\frac{1}{2}(\text{field})^2$. To honour this convention, we finally end up in a Lorentz- and $SU(N)$ -invariant Lagrangian including the fermion part, the **Yang-Mills Lagrangian** [24], of the form

$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu}^a F^{\mu\nu, a} + \bar{\psi} (i\not{D} - m) \psi. \quad (3.11)$$

Where $\not{D} = \gamma^\mu D_\mu$ is the **Feynman slash notation**. Just as in the Abelian theory, there is no term quadratic in A proportional to $m_A A_\mu^a A^{\mu, a}$, a mass-term for the gauge fields, because this would violate the $SU(N)$ -invariance. The fundamental particle which is represented by the field must therefore be massless $m_A = 0$. The constant g can be interpreted as **coupling constant**. The terms in the above Lagrangian are actually the only ones that can appear if we demand the spacetime to have $D = 4$ dimensions and parity- (P) and time-reversal- (T) invariance. The gauge part of the Lagrangian employs interactions among the gauge fields, namely the theory is equipped with pure 3- and 4-vertices of gauge bosons.

4 Lattice Gauge Theories

Non-Abelian gauge theories with a certain number of fundamental fermionic particles can have the property of asymptotic freedom, meaning that the strength of the interaction becomes asymptotically weak as the distance between elementary particles decreases and the energy scale increases. Perturbatively such theories can only be treated in the high energy scale, where the running coupling constant is small to allow perturbative expansion of the problem. Consequently at low energies, the interaction is strong leading to confinement. Lattice gauge theory is a non-perturbative approach to deal with aforementioned theories in the low energy regime by discretising the problem on a finite spacetime lattice. The finiteness of the lattice volume results in a momentum cutoff at $1/a$ curing IR-divergencies, where a is the **lattice spacing** (also called the **lattice constant**). Also, the finiteness of the lattice spacing on the other hand results in a cutoff at a curing UV-divergencies. Lattice discretisation of a field theory therefore acts as a regularisation scheme.

⁶When replacing n_1 and n_2 with unit vectors in arbitrary direction μ and ν respectively, we can set $F_{12}^a = F_{\mu\nu}^a$.

⁷The (arbitrary) normalisation convention used is $tr(T^a T^b) = \frac{1}{2} \delta^{ab}$.

4.1 Euclidean theory

Starting from the continuum Yang-Mills Lagrangian in Minkowski D -dimensional spacetime (the superscript M stands for Minkowski, see equation (3.11))

$$\mathcal{L}_{YM}^M = \mathcal{L}_G^M + \mathcal{L}_F^M, \quad (4.1)$$

with fermion- (F) and gauge-part (G)

$$\begin{aligned} \mathcal{L}_F^M &= \bar{\psi} (i \not{D} - m) \psi, \\ \mathcal{L}_G^M &= -\frac{1}{4} F_{\mu\nu}^a F_a^{\mu\nu}, \end{aligned}$$

where - for simplicity - there is only one fundamental Dirac spinor field with mass m , $D_\mu = \partial_\mu + igA_\mu^a(x)T^a$ is the ***gauge covariant derivative***⁸ (see equation (3.4)), T^a are the generators of the Lie-algebra of the gauge group and $A_\mu^a(x)$ are the (massless) gauge fields introduced in the previous section 3. The (color) index a runs from 1 to $N^2 - 1$, where N is degree of the special unitary symmetry group $SU(N)$. The ***field strength tensor*** is defined as (see equation (3.8))

$$F_{\mu\nu}^a = \partial_\mu A_\nu^a - \partial_\nu A_\mu^a - gf_{abc}A_\mu^b A_\nu^c.$$

The (Minkowski) action is as usual defined as the integral over spacetime of the Lagrangian

$$\mathcal{S}_{YM}^M = \int d^4x \mathcal{L}_{YM}^M.$$

We perform a ***Wick-rotation*** to obtain the Euclidean Lagrangian and action. This is done, because the Wick rotation in path integral formulation translates as $e^{i\mathcal{S}^M} \rightarrow e^{-\mathcal{S}^E}$, where \mathcal{S}^E is a positive real number. The Euclidean path integral is then in the form of a classical statistical mechanics model, enabling us to interpret $e^{-\mathcal{S}^E}$ as probability density.

The Minkowski metric $\eta^{\mu\nu}$ becomes Euclidean if - through analytic continuation - we restrict the time coordinate to take imaginary values. The substitution is (for covariant and contravariant vectors)

$$\begin{aligned} t &\longrightarrow -i\tau, \\ x^0 &\longrightarrow -ix^4, \\ x_0 &\longrightarrow +ix_4, \end{aligned} \quad (4.2)$$

where (real) t is the Minkowski time coordinate and the real number τ is the Euclidean time coordinate. Equation (4.2) only holds in signature $(+, -, \dots, -)$, else the signs in front of the i would be opposite. The fields transform as well, and the transformed Euclidean fields take τ instead of t as time coordinate. We have to take care when transforming the fields and derivatives to Euclidean spacetime. The spinor fields transform as

$$\begin{aligned} \psi(\vec{x}, t) &\longrightarrow S\psi_E(\vec{x}, \tau) \\ \psi(\vec{x}, t)^\dagger &\longrightarrow \psi_E(\vec{x}, \tau)^\dagger S, \end{aligned}$$

where S is a (invertible) matrix in spinor space and still has to be determined. Since the gauge fields are vector quantities, they transform under Wick-rotation just as the coordinates $x^0 \rightarrow -ix^4$

⁸We use the particle physics convention of the metric tensor $\eta_{\mu\nu}$ with signature $(+, -, \dots, -)$. In $D = 4$ dimensions, $\eta_{\mu\nu} = \text{diag}(+1, -1, -1, -1)_{\mu\nu}$

and $x^k \rightarrow x^k$ with $k \in \{1, 2, 3\}$, but the fields appear with lower indices and are therefore covariant. In analogy to the spacetime components $x_0 \rightarrow ix_4$,

$$\begin{aligned} A^{0,a}(\vec{x}, t) &\longrightarrow -i(A_E)^{4,a}(\vec{x}, \tau) \\ A_0^a(\vec{x}, t) &\longrightarrow +i(A_E)_4^a(\vec{x}, \tau) \\ A_k^a(\vec{x}, t) &\longrightarrow (A_E)_k^a(\vec{x}, \tau). \end{aligned}$$

Notice that when in Minkowski space μ takes values $0, 1, 2, 3$, $\mu = 0$ being the time component, but when in Euclidean space μ takes values $1, 2, 3, 4$, where $\mu = 4$ is the time component. Directly from equation (4.2), we obtain the rules for derivatives and integral measures

$$\begin{aligned} dt = dx^0 &\longrightarrow -id x^4 = -id\tau, \\ dx^k &\longrightarrow dx^k, \\ \partial_t = \partial_0 &\longrightarrow i\partial_4 = i\partial_\tau, \\ \partial^0 &\longrightarrow -i\partial^4, \\ \partial_k &\longrightarrow \partial_k. \end{aligned}$$

Let's first transform the fermion and interaction part of \mathcal{L}_{YM}^M (we write the space and time components explicitly)

$$\begin{aligned} \mathcal{L}_F^M &= \bar{\psi} (i\gamma^\mu (\partial_\mu + igA_\mu^a T^a) - m) \psi \\ &= \bar{\psi}(\vec{x}, t) (i\gamma^\mu \partial_\mu - g\gamma^\mu A_\mu^a(\vec{x}, t) T^a - m) \psi(\vec{x}, t) \\ &\xrightarrow{\text{WR}} \bar{\psi}_E^\dagger(\vec{x}, \tau) S i\gamma^0 (i\gamma^0 i\partial_4 + i\gamma^k \partial_k - g\gamma^0 i(A_E)_4^a(\vec{x}, \tau) T^a - g\gamma^k (A_E)_k^a(\vec{x}, \tau) T^a - m) S \psi_E(\vec{x}, \tau). \end{aligned}$$

We want in the Euclidean Lagrangian the term $\gamma_E^\mu \partial_\mu$ to appear with Euclidean versions of the γ -matrices. To fulfill these requirements, we need the following theorem.

Theorem 4.1 (Constructing the Euclidean Clifford algebra). *Let γ^μ obey the Clifford algebra, equation (3.1). Let S be an invertible operator in spinor space. Then the Euclidean γ -matrices defined as*

$$\begin{aligned} \gamma_E^4 &:= S^{-1} \gamma^0 S \\ \gamma_E^k &:= i S^{-1} \gamma^k S. \end{aligned}$$

satisfy the Euclidean Clifford algebra

$$\{\gamma_E^\mu, \gamma_E^\nu\} = 2\delta^{\mu\nu} \cdot id,$$

where $\delta^{\mu\nu}$ has signature $(+, +, \dots, +)$.

Proof. It's straight forward to check the properties

$$\begin{aligned} \{\gamma_E^4, \gamma_E^4\} &= \{S^{-1} \gamma^0 S, S^{-1} \gamma^0 S\} \\ &= 2S^{-1} \gamma^0 S S^{-1} \gamma^0 S \\ &= 2S^{-1} \gamma^0 \gamma^0 S \\ &= S^{-1} \{\gamma^0, \gamma^0\} S \\ &= 2\eta^{00} \cdot id \\ &= 2\delta^{44} \cdot id. \end{aligned}$$

Let $k, l \in \{1, 2, 3\}$

$$\begin{aligned}
\{\gamma_E^k, \gamma_E^l\} &= \{iS^{-1}\gamma^k S, iS^{-1}\gamma^l S\} \\
&= -\{S^{-1}\gamma^k S, S^{-1}\gamma^l S\} \\
&= -S^{-1}\{\gamma^k, \gamma^l\}S \\
&= -S^{-1}2\eta^{kl}S \\
&= 2\delta^{kl} \cdot id.
\end{aligned}$$

And finally, let $k \in \{1, 2, 3\}$

$$\begin{aligned}
\{\gamma_E^4, \gamma_E^k\} &= \{S^{-1}\gamma^0 S, iS^{-1}\gamma^k S\} \\
&= iS^{-1}\{\gamma^0, \gamma^k\}S \\
&= 0.
\end{aligned}$$

□

It remains to determine the operator S . Since the Wick rotation does not affect space coordinates x^k , it should also not rotate the space components of the γ^k [21]. It therefore makes sense to demand that the spatial γ -matrices commute with S and the temporal one satisfies

$$\begin{aligned}
[S, \gamma^k] &= 0, \\
S\gamma^0 &= \gamma^0 S^{-1}.
\end{aligned}$$

Using the above restriction, we end up in

$$\begin{aligned}
S &= e^{\frac{\theta}{2}\gamma^4\gamma^5}, \\
\gamma^4 &:= i\gamma^0, \\
\gamma^5 &:= \gamma^1\gamma^2\gamma^3\gamma^4,
\end{aligned}$$

Where θ is an arbitrary angle. Now we can replace the Minkowski γ -matrices with the Euclidean ones.

$$\begin{aligned}
\mathcal{L}_F^M &\xrightarrow{\text{WR}} \psi_E^\dagger(\vec{x}, \tau) i\gamma^0 \left[-S^{-1}\gamma^0 S \partial_4 + iS^{-1}\gamma^k S \partial_k - m \right. \\
&\quad \left. - igS^{-1}\gamma^0 S (A_E)_4^a(\vec{x}, \tau) T^a - gS^{-1}\gamma^k S (A_E)_k^a(\vec{x}, \tau) T^a \right] \psi_E(\vec{x}, \tau) \\
&= -\psi_E^\dagger(\vec{x}, \tau) i\gamma^0 \left[(S^{-1}\gamma^0 S) \partial_4 - (iS^{-1}\gamma^k S) \partial_k + m \right. \\
&\quad \left. + ig(S^{-1}\gamma^0 S) (A_E)_4^a(\vec{x}, \tau) T^a - ig(iS^{-1}\gamma^k S) (A_E)_k^a(\vec{x}, \tau) T^a \right] \psi_E(\vec{x}, \tau) \\
&= -\bar{\psi}_E(\vec{x}, \tau) \left[\gamma_E^4 \partial_4 + \gamma_E^k \partial_k + m \right. \\
&\quad \left. + ig\gamma_E^4 (A_E)_4^a(\vec{x}, \tau) T^a + ig\gamma_E^k (A_E)_k^a(\vec{x}, \tau) T^a \right] \psi_E(\vec{x}, \tau) \\
&= -\bar{\psi}_E(\vec{x}, \tau) \left[\gamma_E^\mu \partial_\mu + m + ig\gamma_E^\mu (A_E)_\mu^a(\vec{x}, \tau) T^a \right] \psi_E(\vec{x}, \tau). \\
&= -\bar{\psi}_E \left[\gamma_E^\mu \partial_\mu + m + ig\gamma_E^\mu (A_E)_\mu^a T^a \right] \psi_E = -\mathcal{L}_E.
\end{aligned}$$

And we obtain the fermion and interaction part of the Euclidean Lagrangian \mathcal{L}_E . The action transforms as

$$i\mathcal{S}_F^M = i \int d^4x \mathcal{L}_F^M = i \int d^3\vec{x} dt \mathcal{L}_F^M$$

$$\xrightarrow{\text{WR}} -i \int d^3 \vec{x} (-id\tau) \mathcal{L}_F^E = - \int d^4 x \mathcal{L}_F^E = -\mathcal{S}_F^E,$$

as desired to model a probability density in the path integral. Next we look at the pure gauge Lagrangian \mathcal{L}_G^M . The trace over the field strength tensor transforms just as it is

$$\begin{aligned} \mathcal{L}_G^M &= -\frac{1}{4} F_{\mu\nu}^a F_a^{\mu\nu} \\ &\xrightarrow{\text{WR}} -\frac{1}{4} (F_E)_{\mu\nu}^a (F_E)_a^{\mu\nu} = -\mathcal{L}_G^E, \end{aligned}$$

because the Euclidean field strength tensor transforms trivially. It derives as

$$(F_E)_{\mu\nu}^a = \partial_\mu (A_E)_\nu^a - \partial_\nu (A_E)_\mu^a - gf_{abc} (A_E)_\mu^b (A_E)_\nu^c.$$

4.2 Lattice Discretisation

To discretise the theory, from now on we work in $D = 4$ dimensional Euclidean spacetime and discretise the spacetime in terms of a lattice.

Definition 4.1 (4D lattice). *Let $L_0, L_1, L_2, L_3 \in \mathbb{N}$ be natural numbers, define the **full lattice volume** as $V := L_0 \cdot L_1 \cdot L_2 \cdot L_3$ and the **full 4D lattice** of volume V as*

$$\Lambda := \{n = (n_0, n_1, n_2, n_3) \mid n_i \in \{0, 1, \dots, L_i - 1\}, i \in \{0, 1, 2, 3\}\}.$$

We define the **lattice spacing** $a > 0$ (also called **lattice constant**), which is a real number of mass dimension 1⁹. Using this we can span the (finite) 4D spacetime lattice.

We start by discretising the gauge-part of the theory. Here our effort in deriving the gauge-part using the plaquette will bear fruit. We note the gauge-part of the Lagrangian (we drop the subscript E for "Euclidean" from now on)

$$\begin{aligned} \mathcal{L}_G &= \frac{1}{4} F_{\mu\nu}^a F^{\mu\nu, a} \\ &= \frac{1}{2} \sum_{\mu, \nu} \text{tr}(F_{\mu\nu}^a F_{\mu\nu}^b T^a T^b). \end{aligned}$$

Notice that in Euclidean spacetime upper and lower Lorentz-indices do not differ, $F_{\mu\nu}^a = F^{\mu\nu, a}$. The Lagrangian was actually constructed as the trace over the sum of plaquettes in all possible (μ, ν) -subspaces. To see this, we reuse the result from equation (3.7) with $n_1 = e_\mu$ and $n_2 = e_\nu$, where e_μ is the unit vector in direction μ ,

$$\begin{aligned} \hat{U}_{\mu\nu}(\epsilon, x) &= \exp(-ig\epsilon^2 F_{\mu\nu}^a T^a + O(\epsilon^3)) \\ &= \delta_{\mu\nu} \cdot id - ig\epsilon^2 F_{\mu\nu}^a T^a + iO(\epsilon^3) + \frac{1}{2!} (-ig\epsilon^2)^2 F_{\mu\nu}^a F_{\mu\nu}^b T^a T^b + \delta_{\mu\nu} O(\epsilon^6). \end{aligned}$$

This is the plaquette in subspace (e_μ, e_ν) . In the last line, the exponential was written in its series form, where it has to be noted that there is no sum over the Lorentz-indices μ, ν in the expression above. Also when tracking down the $O(\epsilon^3)$ -term, it turns out to be imaginary. The yellow term is actually the same that is traced and summed over in the Lagrangian. Solving the equation for the yellow part and taking the real-part on both sides,

$$F_{\mu\nu}^a F_{\mu\nu}^b T^a T^b = \frac{2}{g^2 \epsilon^4} \text{Re} \left[\delta_{\mu\nu} \cdot id - \hat{U}_{\mu\nu}(\epsilon, x) \right] + \delta_{\mu\nu} O(\epsilon^2)$$

⁹We use particle physics units $c = \hbar = 1$ in which length and time have the same unit. It is therefore legal to use the same a as lattice spacing in all 4 spacetime dimensions.

and inserting this expression into the Lagrangian gives

$$\begin{aligned}\mathcal{L}_G &= \frac{1}{2} \sum_{\mu, \nu} \text{tr}(F_{\mu\nu}^a F_{\mu\nu}^b T^a T^b) \\ &= \frac{1}{g^2 \epsilon^4} \sum_{\mu, \nu} \text{Re tr} \left[\delta_{\mu\nu} \cdot id - \hat{U}_{\mu\nu}(\epsilon, x) \right] + O(\epsilon^2).\end{aligned}$$

One usually uses the same lattice spacing a in all spatial and time directions ($\hbar = 1$). Of importance is only that the value of a is chosen small, because it determines to what length and time scale the system can be resolved. The limit $a \rightarrow 0$ is called the continuum limit.

The lattice is structured as follows.

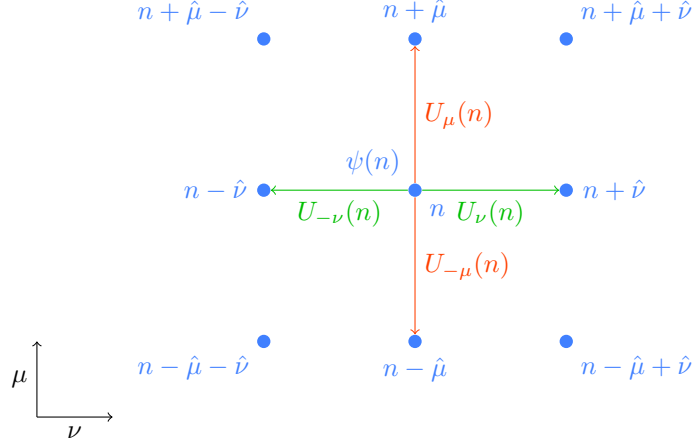


Figure 2: TODO.

Since the link variables $U_\mu(n)$ connect neighbouring spinors $\psi(n)$ and $\psi(n + \hat{\mu})$ they reside visually "in-between" the lattice points. The index μ decorating the link variable is a directed one; it can take 8 values in a 4D lattice, the possible values are therefore $\mu \in \{-4, -3, -2, -1, 1, 2, 3, 4\}$. Using this notation, we have $U_\mu(n) = U_{-\mu}(n + \hat{\mu})^\dagger$, because the connection between lattice site n and $n + \mu$ is equal in both directions. The spinor fields $\psi(n)_{\alpha,a}$ residing on the lattice points carry color- (a) and spinor indices (α).

5 Performance Models

TODO: why are they important? semi-analytical, analytical vs. empirical models

6 Software: openQxD

the software package openQxD: description * importance of CG in openQxD and what it does / how it's used in the software / why 90% computation time

7 Real number formats

7.1 IEEE Standard for Floating-Point Arithmetic

Floating point numbers are omnipresent in the scientific applications. In the **Conjugate Gradient (CG)** kernel of openQxD[5], there are large scalar products over vectors of very high dimensionality over multiple ranks. The components of these vectors are single precision floating point numbers (I call them **binary32** from here on). The precision was degraded from **binary64** to **binary32** already and a speedup of a factor of 2 was achieved. This motivates to explore even smaller floating point

formats with encoding lengths of 16 bits. Since scalar products as well as matrix-vector products are memory-bound operations, going to a smaller bit-length will increase the throughput of the calculation. Therefore, a 16 bits floating point format with a smaller exponent could lead to a double of performance if the new operation is still memory-bound.

Definition 7.1 (IEEE 754 Floating point format). *The **IEEE 754 floating point format** [17] is defined using the **number of exponent bits** e and the **number of mantissa bits** m respectively. A binary floating point number is illustrated in Figure 3.*

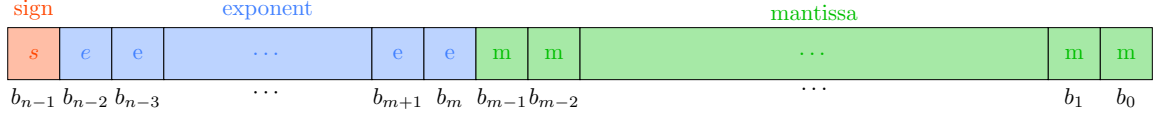


Figure 3: Binary representation of a IEEE 754 n -bit precision floating-point number. The **orange** bit represents the **sign bit**, the **blue** bits represent the fixed-length **e exponent bits** and the **green** bits represent the fixed-length **m mantissa bits**. Notice that $n = 1 + e + m$.

The resulting floating point number is then calculated as

$$f = (-1)^s \cdot M \cdot 2^E,$$

where $E = E' - B$ denotes the biased exponent, B is the exponent bias, M the mantissa and s the sign bit. The (unbiased) exponent E' is calculated as follows

$$E' = \sum_{i=0}^{(e-1)} b_{m+i} 2^i, \quad (7.1)$$

where B is the exponent bias.

Definition 7.2 (Exponent bias).

$$B = 2^{(e-1)} - 1,$$

The calculation of the mantissa is a bit more involved, since it depends on the number being normal or subnormal.

Definition 7.3 (Subnormal numbers). *The IEEE 754 standard introduces so called **subnormal numbers**. If all the exponent bits are 0, meaning the unbiased exponent $E' = 0$, and the mantissa bits are not all 0, then the number is called subnormal. The exponent being zero causes the implicit bit to flip to 0, instead of 1.*

Remark. Subnormal numbers have a variable-length mantissa and exponent, because some of the mantissa bits are used as additional exponent bits, making the numbers less precise the lower they get (see the smooth cutoff in Figure 5).

Therefore the mantissa of a regular (non-subnormal) number is (when the exponent $0 < E < B$, this implies that the implicit bit is 1)

$$M = \underset{\substack{\uparrow \\ \text{implicit bit}}}{1} + \sum_{i=1}^m b_{m-i} 2^{-i},$$

whereas the mantissa of a subnormal number (when the exponent $E = 0$) is

$$M = \underset{\substack{\uparrow \\ \text{implicit bit}}}{0} + \sum_{i=1}^m b_{m-i} 2^{-i},$$

The 1 or 0 in the front of the summand is the leading **implicit bit**, sometimes also called the $(m+1)$ -th **mantissa bit** that tells us whether the number is subnormal or not.

Floating point formats				
name	s	e	m	comment
binary64	1	11	52	double precision, IEEE 754 [18]
binary32	1	8	23	single precision, IEEE 754 [18]
binary16	1	5	10	half precision, IEEE 754 [18]
bfloat16	1	8	7	Googles Brain Float [23]
tensorfloat32	1	8	10	NVIDIAs TensorFloat-32 [14] ¹⁰
binary24	1	7	16	AMDs fp24 [3]
binary128	1	15	112	IEEE 754 [18]
binary256	1	19	236	IEEE 754 [18]

Table 1: Commonly used floating point formats, where s is the number of sign bits, e the number of exponent bits and m the number of mantissa bits.

Remark. The mantissa range of a regular floating point number is $M \in [1, 2)$, whereas the mantissa range of a subnormal floating point number is $M \in (0, 1)$. The number zero is not considered subnormal.

Usual floating point formats are summarised in Table 1.

The format of interest is the **binary16** half precision IEEE 754 floating point format. The highest representable number is when the exponent is highest. This is not the case when all e exponent bits are 1, because then - according to the specification [17] - the number is either $\pm\infty$ or **not a number** (**NaN**), depending on the mantissa. The maximal unbiased exponent is therefore the next smaller number,

$$E'_{max} = \underbrace{1 \dots 1}_e 0.$$

Using equation (7.1), we find

$$\begin{aligned} E'_{max} &= \sum_{i=1}^{(e-1)} 2^i \\ &= 2^e - 2. \end{aligned}$$

The mantissa on the other hand is maximal when all mantissa bits are 1 (including the implicit bit),

$$\begin{aligned} M_{max} &= 1 + \sum_{i=1}^m 2^{-i} \\ &= 2 - 2^{-m}. \end{aligned}$$

Using these two formulas we can define the

Definition 7.4 (highest representable number). *The highest representable number in any floating point format is*

$$\begin{aligned} f_{max} &= (-1)^0 \cdot M_{max} \cdot 2^{(E'_{max}-B)} \\ &= (2 - 2^{-m}) \cdot 2^{(2^e - 2^{e-1} - 1)} \\ &= (2 - 2^{-m}) \cdot 2^{(2^{e-1} - 1)}. \end{aligned}$$

¹⁰Allocates 32 bits, but only 19 bits are actually used.

Floating point format limits				
name	f_{max}	f_{min}	f_{smin}	sign. digits 11
binary64	1.8×10^{308}	2.2×10^{-308}	4.9×10^{-324}	≤ 15.9
binary32	3.4×10^{38}	1.2×10^{-38}	1.4×10^{-45}	≤ 7.2
binary16	6.6×10^4	6.1×10^{-5}	6.0×10^{-8}	≤ 3.3
bfloat16	3.4×10^{38}	1.2×10^{-38}	9.2×10^{-41}	≤ 2.4
tensorfloat32	3.4×10^{38}	1.2×10^{-38}	1.1×10^{-41}	≤ 7.2
binary24	1.8×10^{19}	2.2×10^{-19}	3.3×10^{-24}	≤ 5.1
binary128	1.2×10^{4932}	3.4×10^{-4932}	6.5×10^{-4966}	≤ 34
binary256	$1.6 \times 10^{78,913}$	$1 \times 10^{-78,912}$	$1 \times 10^{-78,983}$	≤ 71.3

Table 2: Summary of highest representable numbers, minimal subnormal and non-subnormal representable numbers above 0 in any IEEE 754 floating point format together with their approximated precision.

The minimal number above 0 can be found similarly, using minimal unbiased exponent (when all exponent bits are 0, except the last one, therefore $E'_{min} = 1$) and the minimal mantissa ($M_{min} = 1$).

Definition 7.5 (*minimal (non-subnormal) representable number above 0*). The *minimal (non-subnormal) representable number above 0 in any floating point format is*

$$\begin{aligned} f_{min} &= (-1)^0 \cdot M_{min} \cdot 2^{(E'_{min}-B)} \\ &= 2^{(2-2^{e-1})}. \end{aligned}$$

The minimal subnormal number can be found, when the unbiased exponent consists of only zeros ($E'_{smin} = 0$) and for the mantissa, only the rightmost bit is one ($M_{smin} = 2^{1-m}$).

Definition 7.6 (*minimal subnormal representable number above 0*). The *minimal subnormal representable number above 0 in any floating point format is*

$$\begin{aligned} f_{min} &= (-1)^0 \cdot M_{smin} \cdot 2^{(E'_{smin}-B)} \\ &= 2^{1-m} \cdot 2^{(1-2^{e-1})} \\ &= 2^{(2-m-2^{e-1})}. \end{aligned}$$

See Table 2 for these limiting numbers in the different floating point formats.

7.2 Posits

The posit datatype is designed to be a replacement for the IEEE floating point format, fixing its various quirks. Some of the more entertaining are:

- The appearance of **NaNs**. They are considered unnatural, because a specific bit pattern describing a number that is not a number is a contradiction.
- The **NaNs** and the fact that floats have two different bit-representations for the number zero (0 and -0) lead to very complicated and slow comparison units as well as (funny) theoretical contradictions¹².
- Floats may under- or overflow, because the standard employs the round to nearest even rounding rule ($\pm\infty$ and 0 are considered even).

¹¹Number of significant digits in decimal; $-\log_{10}(\text{MACHINE_EPSILON}) = \log_{10}(2^{m+1})$.

¹²According to IEEE 754 floating point arithmetic, it holds $\frac{1}{\infty} = 0$ and $\frac{1}{-\infty} = -0$, but we have $0 = -0$. This implies $\infty = -\infty$, but (also according to IEEE 754) $\infty > -\infty$.

- Floats are non-associative and non-distributive¹³ leading to rounding errors that have to be taken into account, specially in scientific computing.
- The standard gives no guarantee of bit-identical results across systems.

The goal is to utilise the number of bits more efficiently and remove these inconsistencies. The key idea is to place half of all numbers between 0 and 1 and the other half are the reciprocals (the reciprocal of 0 being $\pm\infty$). The number can then be drawn on a projective real number circle [13]. The structure of a binary posit number is illustrated in Figure 4.

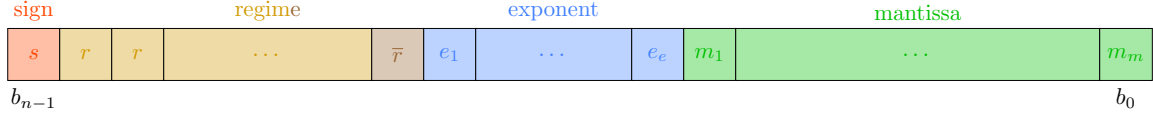


Figure 4: Binary representation of a n -bit posit number. As with regular floats the **orange** bit represents the **sign bit**, the **yellow** bit(s) represent the variable length **regime bit(s)** terminated by the **brown** bit that is the **opposite regime bit**, the **blue** bit(s) represent the variable-length **exponent bit(s)** and the **green** bit(s) represent the variable-length **mantissa bit(s)**.

The actual value of the number is calculated as follows. The yellow and brown bits determine the regime of the number. They either start with a **row of all 0 or all 1** terminated by the **opposite bit** indicating the end of the row. The number of bits in the row are counted as m and if they are all 0 they get a minus sign, the regime being $k = -m$. If they are all 1 the regime is calculated as $k = m - 1$. After the regime is decoded, the remaining bits contain the exponent with at most es bits depending on how much bits remain. If no bits remain the exponent is 0. The exponent and the mantissa are both of variable length. Both can have 0 bits, in this case the number consists of only regime bits. This is the reason why posits have a larger number range than floats. The exponent is encoded as unsigned integer, so there is no bias and no bit pattern denoting special numbers such as subnormals or **NaNs**. Therefore n -bit posits have more numbers than n -bit floats, because they have no **NaNs**. After the exponent - if there are still bits remaining - the fraction follows, else the fraction is just 1.0. Since there are no subnormals the implicit bit is always 1. There are two special numbers that do not follow the above encoding scheme; zero which has the bit pattern of all 0 and $\pm\infty$ with a 1 followed by all 0. These two numbers are reciprocals of each other. A general posit number can therefore be written as

$$p = (-1)^s \cdot useed^k \cdot M \cdot 2^E,$$

where s is the sign bit, $useed$ is defined to be $useed = 2^{2^{es}}$, with es the number of predefined exponent bits, M is the mantissa and E the exponent.

The mantissa is calculated as

$$M = 1 + \sum_{i=1}^m m_i 2^{m-i},$$

where m is the variable number of mantissa bits and the implicit bit in front of the sum is always 1. The exponent is

$$E = \sum_{i=1}^e e_i 2^{e-i},$$

where e is the variable number of exponent bits satisfying $e \leq es$.

Using these two equations, we are now able to calculate the highest representable number and the minimal representable number above 0 in posit format.

¹³There was even a system using IEEE 754 that had non-commutative floating point operations [7].

Posit format limits				
name	<i>es</i>	p_{max}	p_{min}	sign. digits 14
posit64	3	2.0×10^{149}	4.9×10^{-150}	≤ 17.7
posit32	2	1.3×10^{36}	7.5×10^{-37}	≤ 8.1
posit16	1	2.7×10^8	3.7×10^{-9}	≤ 3.6
posit8	0	64	1.6×10^{-2}	≤ 1.5

Table 3: Summary of highest representable numbers, minimal representable numbers above 0 in any posit format together with their approximated precision.

Definition 7.7 (highest representable number). *The highest representable number in any posit format is*

$$\begin{aligned}
p_{max} &= (-1)^0 \cdot useed^{n-2} \\
&= 2^{2^{es}(n-2)}.
\end{aligned}$$

Definition 7.8 (minimal representable number above 0). *The minimal representable number above 0 in any posit format is the reciprocal of the highest representable number p_{max}*

$$\begin{aligned}
p_{min} &= \frac{1}{p_{max}} \\
&= 2^{2^{es}(2-n)}.
\end{aligned}$$

See Table 3 for these limiting numbers in the different posit formats.

Posits employ a feature called the **quire**, which is the generalised answer to the **fused multiply-add** operation that recently found its way into [18] in 2008, where the rounding is deferred to the very end of the operation.

7.3 Floating point numbers in openQxD

To explore how the conjugate gradient kernel in openQxD would perform when using smaller bit lengths, one can look at the exponentials of the numbers in the matrix and vectors, see Figure 6. The plot shows all exponents appearing together with their overall occurrence in percent. The number zero was taken from the plot, because it has biased exponent $E = -127$. The occurrences for zero are given in the legend.

The highest exponent in all 4 runs was $E = 4$, whereas the lowest exponent decreased when the number of lattice points increased. The range of exponents that is representable in **binary16** spans from -24 to $+16$ and is indicated by the **solid orange line** and the **solid pink line**. Between -24 and -14 is the regime of subnormal numbers in **binary16**, with the lowest regular (non-subnormal) exponent indicated by the **solid blue line**. When using half precision instead of single precision, all numbers with exponents below -24 , will be converted to zero, whereas exponents above $+16$ will be cast to $\pm\infty$ depending on the sign of the number. It can be seen, that when calculating the norm of these numbers, only numbers between the **dashed blue line** and the **dashed pink line** will participate. If there is a number above the dashed pink line in the **unsafe region** this number will - after squaring - be cast to ∞ and therefore the norm will be ∞ as well¹⁵. In this case the variable representing the norm $x = \|\vec{v}\|$ should be of higher precision than **binary16**. The plot shows that the Dirac matrix `Dop()` is confined in a narrow exponent regime and a representation in 16-bit floats would suffice. Notice the sparsity the Dirac matrix.

¹⁴Number of significant digits in decimal; $-\log_{10}(\text{MACHINE_EPSILON})$. Notice that posits have **tapered accuracy**; numbers near 1 have much more precision than numbers at the borders of the regime. The precision of floats decreases as well with very large and small numbers, but posit precision decreases faster, see Figure 5.

¹⁵A method to circumvent this is to scale the vector entries during the calculation and scale the result back, exploiting homogeneity of the norm, $\|\vec{v}\| = \frac{1}{s} \|s\vec{v}\|$ for $s \in \mathbb{R}_{>0}$.

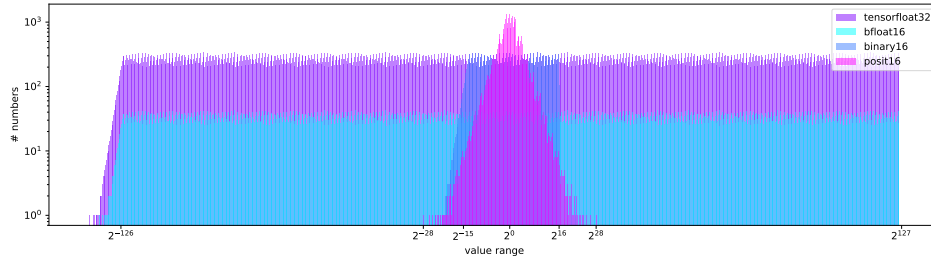


Figure 5: Density or distribution of numbers for `tensorflow32`, `binary16`, `posit16` and `bfloat16`. The number of bins was chosen to be 1024 of logarithmic width. The IEEE conformant floats `tensorflow32`, `binary16` and `bfloat16` exhibit a similar shape, namely the distribution of numbers is exponential decreasing for higher and smaller numbers. The high numbers undergo a rough cutoff at the highest representable number. Numbers above that value will be cast to infinity. Compared to this, the small numbers show a smooth cutoff, because of the existence subnormal numbers. The range of `posit16` is bigger than the range of `binary16`, but specially in the very small numbers this difference in range is negligible. Some features of posits can be observed: First, their distribution is symmetric around 1, because posits have no subnormals. Second, more numbers are closer to 1 than in case of floats; the closer to 1, the better the number resolution. Closest to 1, the number resolution becomes better than `binary16` resolution. Third, posits have no fixed-length mantissa nor exponent. That's the reason why the height of the posit shape depends on the number regime, which happens for floats only in the subnormal regime, where the exponent and mantissa are indeed of variable length. For all formats, the amount of numbers decreases exponentially when going away from 1, but posits decrease faster. This suggests that when calculating in the number regime close to 1 posits might be the better choice, but when numbers span the whole number range equally, floats might be superior. But in that case one has to take care about over- and underflows. Notice that the height of the shape is determined by the number of mantissa bits, therefore giving the precision, whereas the width is determined by the number of exponent bits, therefore giving the number range. For example `tensorflow32` and `binary16` have a very different number range, but exhibit the same precision for numbers in their intersection, meaning that `binary16` is a subset of `tensorflow32`. On the other hand comparing `tensorflow32` and `bfloat16` they have approximately the same number range, but different precisions in them, meaning that `bfloat16` is as well a subset of `tensorflow32`, which itself is a subset of `binary32`. Notice that when plotting `binary32` and `posit32` in such a plot, they would look very similar to `binary16` versus `posit16`.

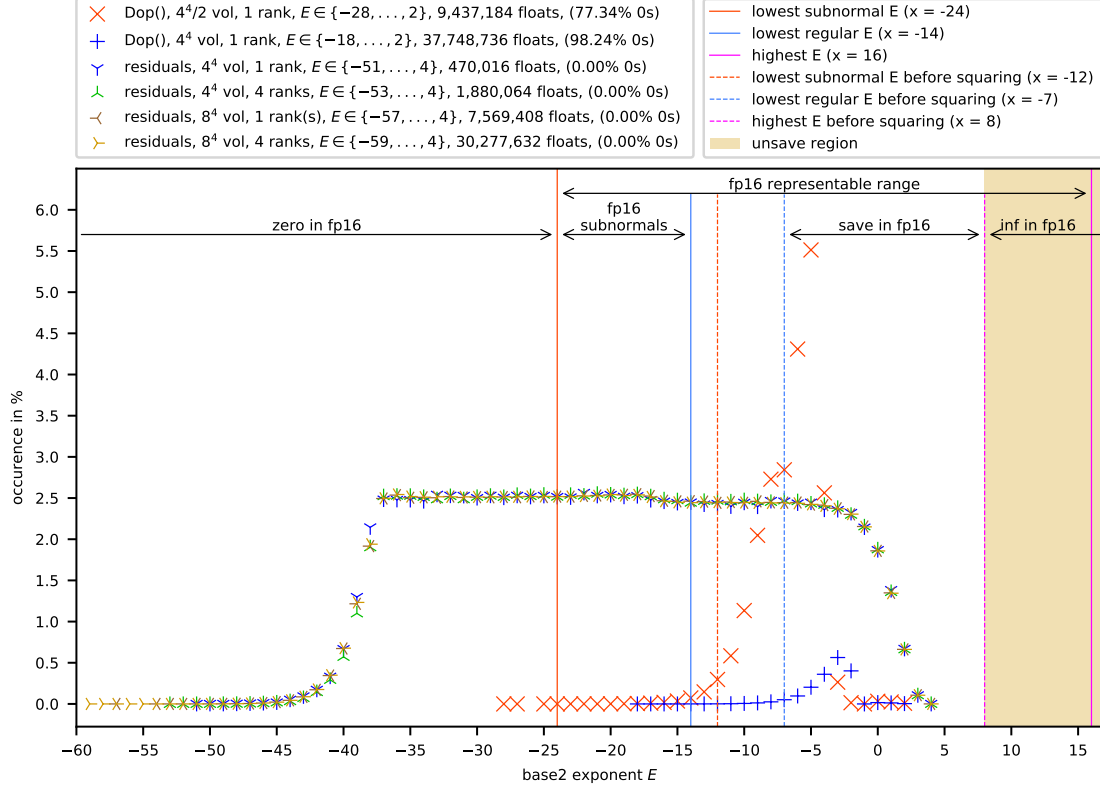


Figure 6: Exponent distribution of **binary32** single precision floats in the residual vectors of all steps in a conjugate gradient run in openQxD as well as entries of the Dirac operator. 4 runs were made, with a lattice size of 4^4 and 8^4 on one single rank and 4 ranks respectively. The number is normalised to $(-1)^s \cdot M \cdot 2^E$, where $M \in [1, 2)$.

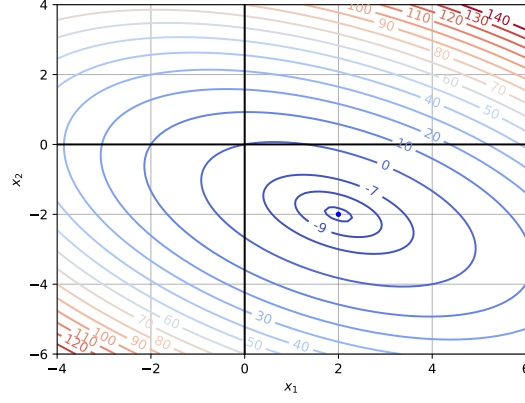


Figure 7: Quadratic form TODO

8 Conjugate Gradient algorithm

In many scientific computations large systems of linear equations need to be solved. Usually these systems are huge and the matrices and vectors are distributed among many **ranks**. The method to solve such systems should therefore be iterative. The problem can be formulated mathematically in the following way.

8.1 Derivation

Let $n \in \mathbb{N}$ and let A be a $n \times n$ -matrix with components in \mathbb{C} , Hermitian, positive definite and **sparse**

$$\begin{aligned} A^\dagger &= A, & (\text{Hermitian}) \\ \forall \vec{x} \in \mathbb{C}^n \setminus \{0\} &: \quad \vec{x}^\dagger A \vec{x} > 0, & (\text{positive definite}) \end{aligned}$$

as well as $\vec{b} \in \mathbb{C}^n$ be given, then the **system of linear equations** can be described as

$$A\vec{x} = \vec{b}. \tag{8.1}$$

We are interested in the **solution** vector \vec{x} , that is the one that satisfies the above equation, n is called the **problem size**. First let us define a function that will be helpful in the next sections.

Definition 8.1 (Quadratic form). *The **quadratic form** depends on the problem matrix A as well as on the **source** vector \vec{b} and is defined as*

$$f(\vec{x}) = \frac{1}{2} \vec{x}^\dagger A \vec{x} - \vec{b}^\dagger \vec{x} + c,$$

where $c \in \mathbb{C}$. When taking the derivative of this function with respect to \vec{x} , we find that

$$f'(\vec{x}) = A\vec{x} - \vec{b}.$$

Therefore finding the extrema of $f(\vec{x})$ is equivalent to solving the linear system of equations (8.1). The question whether the solution \vec{x} is unique remains.

Lemma 8.1 (Uniqueness of the solution). *The solution \vec{x} in equation (8.1) is unique and the global minimum of $f(\vec{x})$ if A is Hermitian and positive definite ¹⁶.*

Proof. Let us rewrite $f(\vec{p})$ at an arbitrary point $\vec{p} \in \mathbb{C}$ in terms of the solution vector \vec{x} :

$$f(\vec{p}) = f(\vec{x}) + \frac{1}{2}(\vec{p} - \vec{x})^\dagger A(\vec{p} - \vec{x}). \quad (8.2)$$

This is indeed the same as $f(\vec{p})$ (inserting $A\vec{x} = \vec{b}$ and using $A^\dagger = A$ and of $\vec{a}^\dagger \vec{b} = \vec{b}^\dagger \vec{a}$),

$$\begin{aligned} f(\vec{x}) + \frac{1}{2}(\vec{p} - \vec{x})^\dagger A(\vec{p} - \vec{x}) &= \frac{1}{2}\vec{x}^\dagger A\vec{x} - \vec{b}^\dagger \vec{x} + c + \frac{1}{2}\vec{p}^\dagger A\vec{p} - \frac{1}{2}\vec{p}^\dagger A\vec{x} - \frac{1}{2}\vec{x}^\dagger A\vec{p} + \frac{1}{2}\vec{x}^\dagger A\vec{x} \\ &= \frac{1}{2}\vec{p}^\dagger A\vec{p} + c + \vec{x}^\dagger \vec{b} - \vec{b}^\dagger \vec{x} - \vec{b}^\dagger \vec{p} \\ &= \frac{1}{2}\vec{p}^\dagger A\vec{p} - \vec{b}^\dagger \vec{p} + c \\ &= f(\vec{p}). \end{aligned}$$

In the new form of $f(\vec{p})$, one can directly see that if A is positive definite, \vec{x} must minimise the function:

$$f(\vec{p}) = f(\vec{x}) + \frac{1}{2} \underbrace{(\vec{p} - \vec{x})^\dagger A(\vec{p} - \vec{x})}_{> 0 \text{ if } A \text{ pos. def.}}$$

Therefore \vec{x} is the global unique minimum. □

TODO: figure of a pos/neg definite quadratic form.

Before deriving the conjugate gradient method, we look at a related method called the **method of steepest descent**. We are interested in a method that iteratively solves equation (8.1) starting at a **initial guess** \vec{x}_0 until the series is interrupted, because the approximate solution \vec{x}_i might be close to the real solution by a certain tolerance or the solution was found exactly,

$$\vec{x}_0 \longrightarrow \vec{x}_1 \longrightarrow \cdots \longrightarrow \vec{x}_i \longrightarrow \vec{x}_{i+1} \longrightarrow \cdots$$

For each step, we can define the **error** and **residual** of the current step i .

Definition 8.2 (Error and Residual). Define the **error** \vec{e}_i and the **residual** \vec{r}_i as

$$\vec{e}_i = \vec{x}_i - \vec{x}, \quad (8.3a)$$

$$\vec{r}_i = \vec{b} - A\vec{x}_i. \quad (8.3b)$$

The residual is the vector of discrepancies and the same as $\vec{r}_i = -f'(\vec{x}_i) = -A\vec{e}_i$, the negative derivative of the quadratic form. The derivative point in direction of the maximum increase, thus the residual points in direction of the steepest descent seen from the position of point \vec{x}_i .

Definition 8.3 (Method of Steepest Descent). The iteration step equation of the **method of steepest descent** is defined as

$$\vec{x}_{i+1} = \vec{x}_i + \alpha_i \vec{r}_i, \quad (8.4)$$

where the $\alpha_i \in \mathbb{C}$ are the amounts to go in direction \vec{r}_i . The α_i are determined by minimising the parabola with respect to α_i , $\frac{d}{d\alpha_i} f(\vec{x}_{i+1}) \stackrel{!}{=} 0$.

TODO: figure of steepest descent zigzag.

¹⁶Notice that, negative definiteness is sufficient as well and \vec{x} would be the global maximum instead - just define $A' = -A$ which is positive definite and all of the argumentation that follows will hold as well. Indefinite matrices on the other hand might have local minima and maxima.

Remark (Convergence). As seen in figure [TODO], the method of steepest descent converges very slowly to the actual solution, when starting at an unfavourable starting point \vec{x}_0 . The speed of convergence also heavily depends on the condition number of matrix A . We see that the iteration goes in the same direction multiple times. How about, when we only go *once* in each direction i , but by the perfect amount α_i ? Then we would be done after at most n steps.

This gives motivation for a enhanced method. Let's define a new **step equation** as

$$\vec{x}_{i+1} = \vec{x}_i + \alpha_i \vec{p}_i, \quad (8.5)$$

with **directions** \vec{p}_i and **amounts** α_i that have to be determined. But this time, we will impose the condition to go in every direction only once at most. This will lead us to the **method of conjugate gradient**.

Using the step equation (8.5), we can update the error and residuals,

$$\vec{e}_{i+1} = \vec{x}_{i+1} - \vec{x} \quad (8.6a)$$

$$= \vec{e}_i + \alpha_i \vec{p}_i \quad (8.6b)$$

$$= \vec{e}_0 + \sum_{j=0}^i \alpha_j \vec{p}_j, \quad (8.6c)$$

$$\vec{r}_{i+1} = \vec{b} - A\vec{x}_{i+1} \quad (8.7a)$$

$$= \vec{r}_i - \alpha_i A\vec{p}_i \quad (8.7b)$$

$$= -A\vec{e}_{i+1}. \quad (8.7c)$$

The $\{\vec{p}_i\}$ need to form a basis of \mathbb{C}^n , because the method should succeed with any arbitrary initial guess \vec{x}_0 . Since we move in the vector space \mathbb{C}^n from an arbitrary point \vec{x}_0 to the solution \vec{x} , the n direction vectors need cover all possible directions in the space, therefore need to be linear independent.

To be done after at most n steps, we need that the n -th error is zero, $\vec{e}_n = 0$. Since the directions form a basis, we can write \vec{e}_0 as a linear combination of the $\{\vec{p}_i\}$,

$$\vec{e}_0 = \sum_{j=0}^{n-1} \delta_j \vec{p}_j.$$

Using this we can rewrite \vec{e}_n ,

$$\begin{aligned} \vec{e}_n &= \vec{e}_0 + \sum_{j=0}^{n-1} \alpha_j \vec{p}_j \\ &= \sum_{j=0}^{n-1} \delta_j \vec{p}_j + \sum_{j=0}^{n-1} \alpha_j \vec{p}_j \\ &= \sum_{j=0}^{n-1} (\delta_j + \alpha_j) \vec{p}_j. \end{aligned}$$

In order for this to be zero, all coefficients need to be zero, thus $\delta_j = -\alpha_j$. Then the i -th error can be written in a different way

$$\vec{e}_i = \vec{e}_0 + \sum_{j=0}^{i-1} \alpha_j \vec{p}_j$$

$$\begin{aligned}
&= \sum_{j=0}^{n-1} \delta_j \vec{p}_j - \sum_{j=0}^{i-1} \delta_j \vec{p}_j \\
&= \sum_{j=i}^{n-1} \delta_j \vec{p}_j.
\end{aligned} \tag{8.8}$$

In the last row, we can see that after every step in the iteration, we shave off the contribution of one direction \vec{p}_i to the initial error \vec{e}_0 (or phrased differently: \vec{e}_{i+1} has no contribution from direction \vec{p}_i). But we still need to find these directions. We could for example impose that the $(i+1)$ -th error should be orthogonal to the i -th direction, because we never want to go in that direction again,

$$\begin{aligned}
0 &\stackrel{!}{=} \vec{p}_i^\dagger \vec{e}_{i+1} \\
&= \vec{p}_i^\dagger (\vec{e}_i + \alpha_i \vec{p}_i).
\end{aligned}$$

This gives us a expression for the amount α_i ,

$$\alpha_i = -\frac{\vec{p}_i^\dagger \vec{e}_i}{\vec{p}_i^\dagger \vec{p}_i}.$$

The problem with this expression is that we don't know the value of \vec{e}_i - if we would, we could just subtract it from the current \vec{x}_i and obtain \vec{x} exactly. So, we do not know \vec{e}_i , but what we actually know is something similar, namely $-A\vec{e}_i$, with is the residual. So if we manage to sandwich an A in the expression above, we are save. It turns out that imposing A -orthogonality instead of regular orthogonality between \vec{e}_{i+1} and \vec{p}_i achieves what we're up to by the exact same steps¹⁷,

$$\begin{aligned}
0 &\stackrel{!}{=} \vec{p}_i^\dagger A \vec{e}_{i+1} \\
&= \vec{p}_i^\dagger A (\vec{e}_i + \alpha_i \vec{p}_i)
\end{aligned}$$

Solving for α_i gives the (almost) final expression for the amounts,

$$\implies \alpha_i = -\frac{\vec{p}_i^\dagger A \vec{e}_i}{\vec{p}_i^\dagger A \vec{p}_i} = \frac{\vec{p}_i^\dagger \vec{r}_i}{\vec{p}_i^\dagger A \vec{p}_i}. \tag{8.9}$$

Notice that the denominator is never zero, because A is positive definite. Let us continue with the expression for A -orthogonality, but insert the derived expression (8.8) for \vec{e}_{i+1} this time,

$$\begin{aligned}
0 &\stackrel{!}{=} \vec{p}_i^\dagger A \vec{e}_{i+1} \\
&= \vec{p}_i^\dagger A \left[\sum_{j=i+1}^{n-1} \delta_j \vec{p}_j \right] \\
&= \sum_{j=i+1}^{n-1} \underbrace{\delta_j}_{\neq 0} \vec{p}_i^\dagger A \vec{p}_j.
\end{aligned}$$

This implies that for $j > i$ and $i \in \{0, \dots, n-1\}$, we have

$$\vec{p}_i^\dagger A \vec{p}_j = 0.$$

¹⁷This is equivalent to imposing $0 \stackrel{!}{=} \vec{r}_{i+1}^\dagger \vec{p}_i$ which is done in most literature, but in the opinion of the author this is less intuitive.

But since A is Hermitian, we can Hermitian conjugate the whole expression above and obtain

$$0 = \left(\vec{p}_i^\dagger A \vec{p}_j \right)^\dagger = \vec{p}_j^\dagger A \vec{p}_i.$$

So the expression holds for $i > j$ as well, which implies that the $\{\vec{p}_i\}$ are ***A-orthogonal***,

$$\vec{p}_i^\dagger A \vec{p}_j = 0 \quad \forall i \neq j.$$

So the problem has reduced to finding a set of A -orthogonal vectors in an iterative way. Luckily there is a well know method to find orthogonal vectors from a set of linear independent vectors: ***Gram-Schmidt orthogonalization***. The procedure can be altered to find A -orthogonal vectors instead.

Definition 8.4 (Gram-Schmidt Orthogonalization). *Let $\{\vec{u}_0, \dots, \vec{u}_{n-1}\} \subset \mathbb{C}^n$ be a set of n linear independent vectors. The iterative Gram-Schmidt procedure is*

$$\begin{aligned} \vec{p}_0 &= \vec{u}_0 \\ \vec{p}_i &= \vec{u}_i + \sum_{k=0}^{i-1} \beta_{ik} \vec{p}_k, \end{aligned} \tag{8.10}$$

where the $\beta_{ik} \in \mathbb{C}$ are (to be determined) coefficients. In the regular procedure, the β_{ik} are just normalised projections of \vec{u}_i to \vec{p}_k that are subtracted from \vec{u}_i , leading to a vector \vec{p}_i that is orthogonal to all previously calculated \vec{p}_k .

In our problem, we need a set of vectors that are A -orthogonal. By imposing this condition we find a different expression for the β_{ik} ,

$$\begin{aligned} 0 &\stackrel{!}{=} \vec{p}_i^\dagger A \vec{p}_j \\ &= \vec{u}_i^\dagger A \vec{p}_j + \sum_{k=0}^{i-1} \beta_{ik} \vec{p}_k^\dagger A \vec{p}_j \\ &= \vec{u}_i^\dagger A \vec{p}_j + \beta_{ij} \vec{p}_j^\dagger A \vec{p}_j, \end{aligned}$$

where in the last step, we assumed $i > j$ (else we would not find an expression for β_{ij}) and therefore only the j -th term in the sum remains, because of the A -orthonormality of the directions. Solving this for β_{ij} gives

$$\beta_{ij} = - \frac{\vec{u}_i^\dagger A \vec{p}_j}{\vec{p}_j^\dagger A \vec{p}_j}. \tag{8.11}$$

In principle we are done here, we only need a set of linearly independent vectors $\{\vec{u}_i\}$. Since the conjugate gradient method is iterative and often dealing with huge problem sizes n , we need to store all previous directions \vec{p}_k in order to calculate the current direction (see equation (8.10)). This becomes a problem in limited memory situations. We want that the current step only depends on the previous one. By imposing this condition, we need the sum in equation (8.10) to collapse; the β_{ik} should only be non-zero for $k = i - 1$. If we manage to satisfy this, the orthogonalization procedure would simplify to

$$\begin{aligned} \beta_i &:= \beta_{i,i-1}, \\ \vec{p}_i &= \vec{u}_i + \beta_i \vec{p}_{i-1}, \end{aligned}$$

where in the second equation, the current \vec{p}_i only depends on the previous \vec{p}_{i-1} . For this to hold, all other β_{ij} need to be zero. For such a β_{ij} the numerator needs to be zero. Let therefore $j < i - 1$

$$\vec{u}_i^\dagger A \vec{p}_j \stackrel{!}{=} 0.$$

To find a different expression for the left hand side, consider

$$\begin{aligned} \vec{u}_i^\dagger \vec{r}_{j+1} &= \vec{u}_i^\dagger (\vec{r}_j + \alpha_j A \vec{p}_j) \\ &= \vec{u}_i^\dagger \vec{r}_j + \alpha_j \vec{u}_i^\dagger A \vec{p}_j, \\ \implies \vec{u}_i^\dagger A \vec{p}_j &= \frac{1}{\alpha_j} \left[\vec{u}_i^\dagger \vec{r}_{j+1} - \vec{u}_i^\dagger \vec{r}_j \right], \end{aligned} \tag{8.12}$$

where we inserted the recursive relation of the residuals (8.7b) and the yellow part is the expression we want to be zero for $j < i - 1$. We therefore find a condition for the linear independent set $\{\vec{u}_i\}$, namely that the scalar product of \vec{u}_i with \vec{r}_{j+1} and \vec{r}_j must be the same. But we can apply the same equation over and over again and obtain

$$\vec{u}_i^\dagger \vec{r}_{j+1} = \vec{u}_i^\dagger \vec{r}_j = \dots = \vec{u}_i^\dagger \vec{r}_0, \quad j < i - 1$$

We have to find $\{\vec{u}_i\}$ that satisfy the above equation. It is sufficient to find a set of $\{\vec{u}_i\}$ that are orthogonal to all the residuals and the equation would be obeyed.

Lemma 8.2. *The residuals are orthogonal, thus for all $i \neq j$, it holds*

$$\vec{r}_i^\dagger \vec{r}_j = 0.$$

Proof. The proof consists of 2 steps.

1) Let $i < j$,

$$\begin{aligned} \vec{p}_i^\dagger \vec{r}_j &= -\vec{p}_i^\dagger A \vec{e}_j \\ &= -\sum_{k=j}^{n-1} \delta_j \vec{p}_i^\dagger A \vec{p}_k \\ &= 0, \end{aligned}$$

where the yellow expression is zero, because $i < j \leq k$.

2) Let $i < j$. By step 1), we have

$$\begin{aligned} 0 &= \vec{p}_i^\dagger \vec{r}_j \\ &= \vec{r}_i^\dagger \vec{r}_j + \sum_{k=0}^{i-1} \beta_{ik} \vec{p}_k^\dagger \vec{r}_j \\ &= \vec{r}_i^\dagger \vec{r}_j. \end{aligned}$$

The yellow expression is again zero by step 1). Using the symmetry of the scalar product, the above equation also holds for i and j interchanged ($i > j$), therefore holds for all $i \neq j$.

□

From now on we set $\vec{u}_i = \vec{r}_i$. What remains to find is the final expression for the β_i .

$$\begin{aligned}
\beta_i &:= \beta_{i,i-1} = -\frac{\vec{u}_i^\dagger A \vec{p}_{i-1}}{\vec{p}_{i-1}^\dagger A \vec{p}_{i-1}} \\
&= -\frac{1}{\vec{p}_{i-1}^\dagger A \vec{p}_{i-1}} \frac{1}{\alpha_{i-1}} \left[\vec{r}_i^\dagger \vec{r}_i - \vec{r}_i^\dagger \vec{r}_{i-1} \right] \\
&= -\frac{\vec{r}_i^\dagger \vec{r}_i}{\alpha_{i-1} \vec{p}_{i-1}^\dagger A \vec{p}_{i-1}} \\
&= -\frac{\vec{r}_i^\dagger \vec{r}_i}{\vec{p}_{i-1}^\dagger \vec{r}_{i-1}},
\end{aligned}$$

where in the first row we used the definition (8.11), in the second row we have used equation (8.12) and the yellow expression is zero by the orthogonality of the residuals lemma 8.2. In the last line we used the expression for the α_j equation (8.9)

To obtain the final form of the α_i and the β_i , we can use a leftover of the proof of lemma 8.2, namely

$$\begin{aligned}
\vec{p}_i^\dagger \vec{r}_i &= \vec{r}_i^\dagger \vec{r}_i + \beta_i \underbrace{\vec{p}_{i-1}^\dagger \vec{r}_i}_{= 0 \text{ by lemma 8.2 step 1)}} \\
&= \vec{r}_i^\dagger \vec{r}_i.
\end{aligned}$$

Using this we find the final form of the α_i and the β_i as well as the **method of conjugate gradient**.

Definition 8.5 (Method of conjugate gradient). *The iteration step equation of the **method of conjugate gradient** is defined as*

$$\vec{x}_{i+1} = \vec{x}_i + \alpha_i \vec{p}_i,$$

with

$$\vec{r}_{i+1} = \vec{r}_i - \alpha_i A \vec{p}_i, \quad \alpha_i = \frac{\vec{r}_i^\dagger \vec{r}_i}{\vec{p}_i^\dagger A \vec{p}_i}, \quad (8.13)$$

$$\vec{p}_{i+1} = \vec{r}_{i+1} + \beta_{i+1} \vec{p}_i, \quad \beta_{i+1} = -\frac{\vec{r}_{i+1}^\dagger \vec{r}_{i+1}}{\vec{r}_i^\dagger \vec{r}_i}, \quad (8.14)$$

and initial starting vectors

$$\begin{aligned}
\vec{x}_0 &= \text{arbitrary starting point}, \\
\vec{p}_0 &= \vec{r}_0 = \vec{b} - A \vec{x}_0.
\end{aligned} \quad (8.15)$$

There are some remarks to note about the method of conjugate gradient.

Remark. The β_{i+1} of the current iteration depends on the norm of the current residual as well as the last one. This means that we can store the result of the last iteration and reuse it in the current, the norm may not be calculated twice.

Remark. In the source code of openQxD (see [5]) the matrix A is the Dirac matrix applied twice $A = D^\dagger D$. This means that the denominator of α_i is a regular inner product as well; $\vec{p}_i^\dagger A \vec{p}_i = \vec{p}_i^\dagger D^\dagger D \vec{p}_i = (D \vec{p}_i)^\dagger (D \vec{p}_i) = \|D \vec{p}_i\|^2$

Remark. Therefore in each iteration, we have:

- 2 times the norm of a vector,

- 2 matrix-vector multiplications,
- 3 times axpy.¹⁸

Remark (Floating point errors). Since the method contains recursive steps, floating point round-off accumulation is an issue. This causes the residuals to loose their A -orthogonality. It can be resolved by calculating the residual from time to time using its (computationally more expensive) definition $\vec{r}_i = \vec{b} - A\vec{x}_i$, which involves one matrix vector multiplication. One can for example do this every m -th step. The same problem applies to the directions \vec{p}_i that loose their A -orthogonality.

Remark (Problem size). The method of conjugate gradient is suitable for problems of very huge size n . The algorithm is done after n steps, but there might be problems such that even n steps are out of reach for an exact solution.

Remark (Complexity). The time complexity of the conjugate gradient method is $O(m\sqrt{\kappa})$, where m is the number of non-zero entries in A and κ is its **condition number**. The space complexity is $O(m)$.

Remark (Starting). The **starting vector** \vec{x}_0 can be chosen at wish. If there is already a rough estimate of the solution one can take that vector. But usually just $\vec{x}_0 = 0$ is chosen. Since the minimum is global, there is no issue in choosing a starting point. The method will always converge towards the real solution.

Remark (Stopping). If the problem size does not allow to run n steps, one can stop when the norm of the residual falls below a certain **threshold** value. Usually this threshold is a fraction of the initial residual $\|\vec{r}_i\| < \epsilon\|\vec{r}_0\|$ [20].

Remark (Initialisation). The very first step of the method is equivalent to a step in the method of steepest descent, see equation (8.4).

Remark (Speed of convergence). TODO: cg is quicker if there are duplicated eigenvalues. number of iterations for exact solution is at most the number of distinct eigenvalues.

Remark (Preconditioning). The linear system of equations can be transformed using a matrix M to

$$M^{-1}A\vec{x} = M^{-1}\vec{b}.$$

It is assumed M is such that is is easy to insert and it approximates A in some way, resulting in $M^{-1}A$ to be better conditioned than was A . An examples of a particular preconditioner M would be a diagonal matrix, with diagonal entries of D . It is indeed easy to invert and it approximates A quite well if A has non-zero diagonal entries and most off-diagonal entries are zero.

Remark (Conjugate Gradient on the normal equations (CGNE)). The algorithm can be used even if A is not symmetric nor Hermitian nor positive definite. The linear system of equations to be solved is then

$$A^\dagger A\vec{x} = A^\dagger \vec{b}.$$

If A is square and invertible, solving the above equation is equivalent to solving $A\vec{x} = \vec{b}$. Conjugate gradient can be applied, because $A^\dagger A$ is Hermitian and positive ($\vec{x}^\dagger A^\dagger A\vec{x} = \|A\vec{x}\|^2 \geq 0$). Notice that $A^\dagger A$ is less sparse than A , and often $A^\dagger A$ is badly conditioned.

8.2 CG kernel in openQxD

The conjugate gradient kernel `cgne()` in `modules/linsolv/cgne.c` in [5] implements the algorithm, see Listing 1. The algorithm is already implemented in mixed precision using `binary32` in most of the computations and `binary64` in correction steps¹⁹.

The function expects the Dirac matrix `Dop()` in `binary32`, `Dop_double()` in `binary64` format and the source vector `eta` (\vec{b}) in `binary64` only. In the initialisation the starting vector `psi` (\vec{x}_0) is set

¹⁸This stands for $a\vec{x} + \vec{y}$, scalar times vector plus vector, "a x plus y" (to resemble the BLAS level 1 routine call of the same name).

¹⁹The method is also referred to as *mixed precision defect-correction*, see ref. [11]

```

429 double cgne(int vol,int icom,void (*Dop)(spinor *s,spinor *r),
430             void (*Dop_double)(spinor_double *s,spinor_double *r),
431             spinor **ws,spinor_double **wsd,int nmix,double res,
432             spinor_double *eta,spinor_double *psi,int *status)
433 {

```

Listing 1: The conjugate gradient kernel in `modules/linsolv/cgne.c` line 429ff.

```

490 if ((rn<=tol)|| (rn<=(PRECISION_LIMIT*xn))||(ncg>=100)||
491     ((*status)>=nmix))
492     break;

```

Listing 2: break condition in `modules/linsolv/cgne.c` line 490ff, `rn` is the norm of the current residual, `xn` is the norm of the current solution vector, both in `binary32`.

to zero. The algorithm stops when the desired maximal relative residue $\mathbf{res} (= \frac{\|\mathbf{eta} - D^\dagger D \mathbf{psi}\|}{\|\mathbf{eta}\|})$ is reached, where `psi` is the calculated approximate solution of the Dirac equation $D^\dagger D \mathbf{psi} = \mathbf{eta}$ in `binary64`. For this, the tolerance `tol` is calculated using $\mathbf{tol} = \|\mathbf{eta}\| * \mathbf{res}$. The parameter `nmix` is the maximal number of iterations that may be applied and `status` reports the total number of iterations that were required, or a negative value if the algorithm failed. `icom` is a control parameter and `ws` and `wsd` are work space allocations. The volume of the lattice should be given in `vol`.

Since the Dirac matrix is given in two precisions, the algorithm in the code bails out of the main conjugate gradient loop, when some particular conditions where met, see Listing 2.

This may happen in 4 cases:

1. if the recursively calculated residual is below the tolerance,
2. if the precision of `binary32` is reached²⁰,
3. after a hard coded number of 100 steps,
4. if the maximal number of steps is reached.

Point 2 is the most interesting condition, because lets imagine that this condition is met, but the algorithm does not break out of the main loop. Therefore the norm of the current residual compared to the norm of the current solution vector differ in their orders of magnitude by the precision limit of the datatype (`binary32` in this case). This means that the solution vector \vec{x}_i contains large numbers compared to the residual vector \vec{r}_i . Therefore the changing in residual from iteration to iteration is small compared to numbers in \vec{x}_i as well. Since \vec{r}_i contains small numbers, the amounts α_i are small as well. This causes $\vec{x}_{i+1} = \vec{x}_i + \alpha_i \vec{d}_i$ to not change anymore, because adding very large and very small numbers in floating point arithmetic will return the larger number unchanged if the two numbers differ in magnitude by the precision limit of the datatype. The algorithm stalls in that case and breaking out of the main loop is the emergency brake.

So when one of the above conditions are met, the algorithm performs a *reset step*. A reset step consists of calculating the residual not in the recursive way, instead calculating it in it's definition $\vec{r}_i = \vec{b} - A\vec{x}_i$ in double precision. This involves 2 invocations of each `Dop_double()` as well as `Dop()` which is very expensive. The algorithm is resetting in the sense that the solution vector is set back to $\vec{x}_i = 0$, but before resetting, the solution vector in `binary32` is added to the real solution vector `psi` in `binary64` which was initialised to zero at the start of the algorithm as well. It looks like a restart of the whole calculation, but the direction for the next iteration $\vec{d}_i = \vec{r}_i$ is set to the just calculated, very accurate residual. Therefore the the algorithm now continues in a new direction A -orthogonal to all previous directions and progression is kept. The step is meant to remove the accumulated round-off errors due to the recursive calculation of the residuals and directions. The first step following a reset step is a step in the direction of steepest descent just like the very first step of the algorithm. The less precise the datatype, the more reset steps need to be taken, because the precision limit is reached earlier.

²⁰The constant `PRECISION_LIMIT` is defined to be $100 * \text{MACHINE_EPSILON}$, where the `MACHINE_EPSILON` is the difference between 1 and the lowest value above 1 depending on the datatype. In case of `binary32` the `MACHINE_EPSILON` takes a value of $1.192,092,9 \times 10^{-7}$.

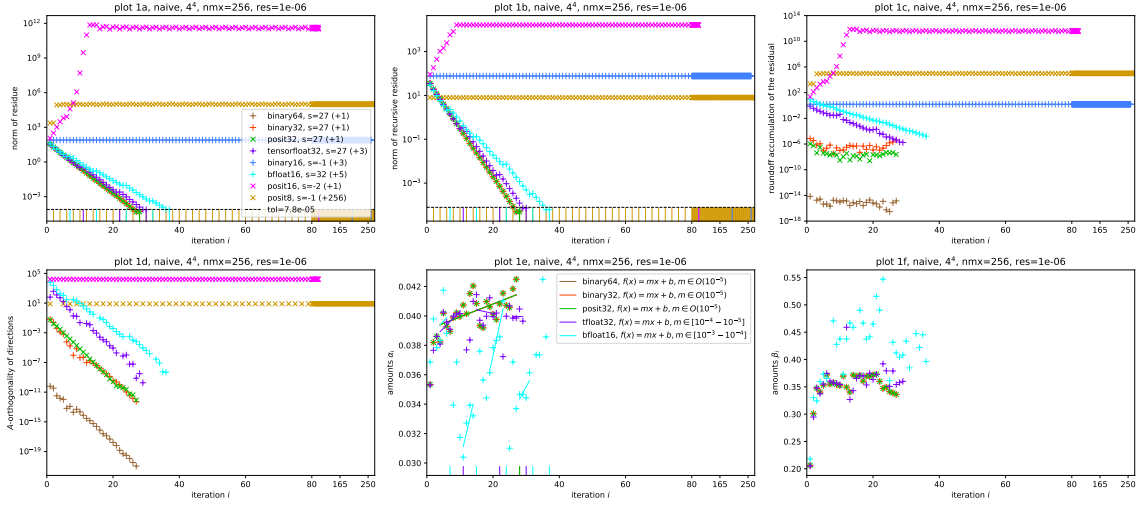


Figure 8: Convergence analysis of a conjugate gradient run, where **binary32** was replaced by one of the simulated datatypes. The number **s** describes the number of normal steps needed (the value of **status**), whereas the numbers in the brackets indicate the number of reset steps. All reset steps are indicated by ticks at the dashed black line denoting the tolerance limit. The iterations will always go up to **nmx=256**, but the range 80-256 is compressed since the most interesting behaviour happens before step 80 for most of the simulated datatypes. The 6 plots show the naive replacement of the **binary32** datatype with the simulated one. This means that every single variable containing a **binary32** was replaced with a variable of the simulated datatype. Plot *1a* shows the exact residue (8.7a) calculated in every iteration using the Dirac matrix and the source vector both in **binary64**, whereas plot *1b* shows the norm of the recursively calculated residue (8.7b) (cast from the simulated datatype to **binary64**). The relative residue suffers round-off accumulation because of the recursive calculation; this is the difference between plots *1a* and *1b*, which is plotted in plot *1c*. Plot *1d* shows the A -orthogonality of the current direction to the last direction, namely the value of $\vec{p}_i^\dagger A \vec{p}_{i+1}$. The last 2 plots, *1e* and *1f*, show the values of the amounts α_i and β_i (see equations (8.13) and (8.14)) in every iteration, but only of the datatypes that converged (**status**>0). The lines in plot *1e* are linearly fitted to the data points ($f(x) = mx + b$). The number range of the slope m is given in the plot legend.

8.3 Simulating CG with different datatypes

Some operations such as norms and scalar products are memory-bandwidth-bound, which means the on-chip memory bandwidth determines how much time is spent computing the output. Storing input data in a format with lower bit-length reduces the amount of data to be transferred, thus improving the speed of calculation.

The complete conjugate gradient kernel was simulated in different datatypes, floats as well as posits. In order to produce the plots, the Dirac matrix `Dop_double()` and the source vector `eta` were extracted in **binary64** format from the original code running a simulation of a 4^4 lattice, **Schrödinger functional** (SF) boundary conditions (**type** 1), no C* boundary conditions (**cstar** 0) and 1 rank. The first 2000 trajectories were considered of thermalisation. The matrix was extracted in trajectory 2001. A python script mimicking the exact behaviour of the `cgne()` kernel from the source code²¹, was implemented to cope with arbitrary datatypes. The simulated datatypes were **binary64**, **binary32**, **tensorfloat32**, **binary16**, **bfloat16**, **posit32**, **posit16**, and **posit8**. The Dirac matrix had approximately 2% non-zero value. The results are plotted in figures 8, 9, 10 and 11.

8.3.1 Discussion of figures 8 - 11

Figures 8, 9, 10 and 11 contain all relevant data. It is expected in general that the plots show datatypes of the same bit-length in clusters and exhibit a hierarchy in precision and exponent

²¹See line 429ff in `modules/linsolv/cgne.c` in [5].

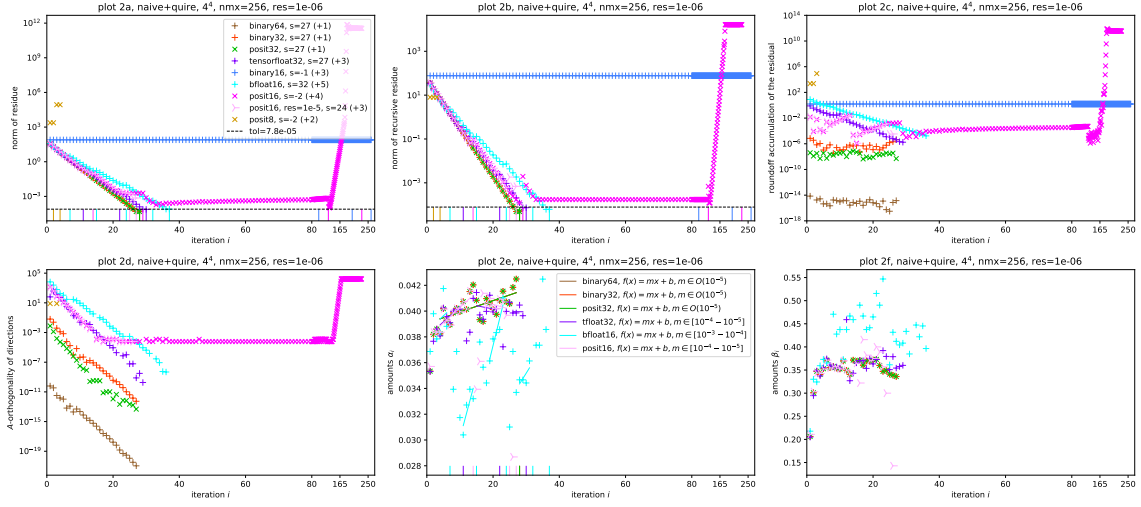


Figure 9: In these plots, the posits were utilising **quires** as their collective variables, the remaining setup was the same as for figure 8, therefore the floating point datatypes show exactly the same values, only posits changed their behaviour.

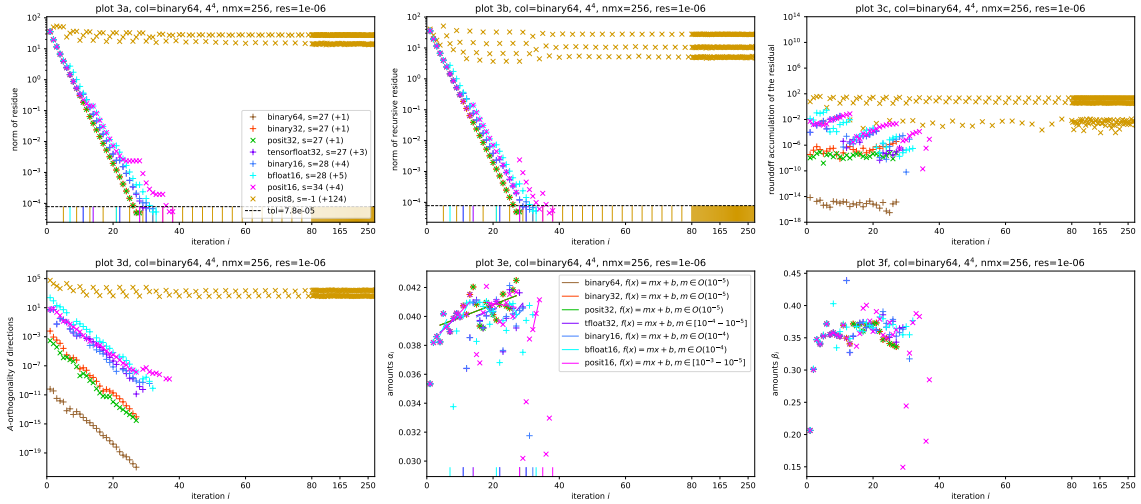


Figure 10: The 6 plots introduce a slightly smarter replacement. All collective variables such as norms were calculated in **binary64**, such that a datatype with a small number range such as **binary16** may not over- or underflow when calculating the norm of a vector full of said datatype. This replacement resembles the **quire** for posits. Using this replacement, even heavily reduced datatypes like **binary16** and **posit16** converged and threw a result of equal quality as the one simulated with **binary64**.

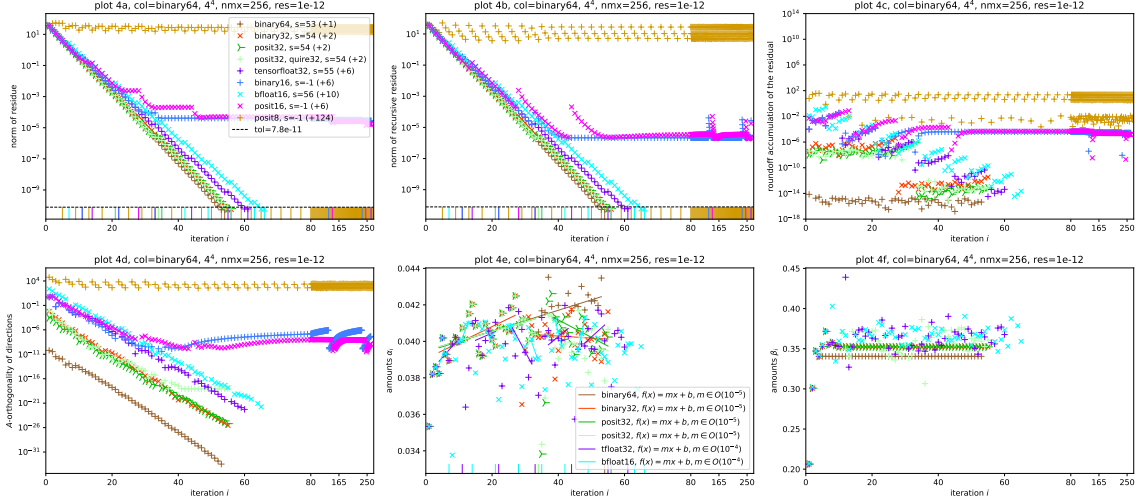


Figure 11: The configuration in this series of plots is equal to Figure 10, besides the value of **res** - the desired relative residue of the calculated solution - is set to 10^{-12} instead of 10^{-6} . Notice that 10^{-12} is outside the representable number range of the datatypes that did not converge; **binary16**, **posit16** and **posit8**.

range; more precision and larger exponent range should end up in faster convergence. Thus we expect the following hierarchy (where smaller means convergence in fewer steps)

$$\text{binary64} < \text{posit32} \leq \text{binary32} \leq \text{tensorfloat32} \leq (1) \leq \text{posit16} \leq \text{binary16} \leq (2) < \text{posit8}, \quad (8.16)$$

where **bfloat16** could be either at position (1) or (2), depending on what is more important; precision or number range.

In Figure 8 where the datatype is naively replaced by the simulated datatype, it can be concluded that only datatypes with large enough number ranges converged. **binary64**, **binary32** and **posit32** converged each after **status=27** steps with one reset step. The less precise **tensorfloat32** took **status=27** (+3) and the even less precise **bfloat16** needed **status=32** (+5) steps. Such a hierarchical result was expected since they have the same exponent range and thus approximately the same number range, but differ only in precision (see Table 1). Notice that the less precise the datatype, the more reset steps are needed. This happens because the precision limit of the simulated datatype is reached faster, if the datatype has less precision.

The round-off accumulation error of **posit32** is slightly better than the one of **binary32**, although defeated by 8 orders of magnitude of **binary64** because of its much more precision. It is notable to remark that the round-off accumulation does not increase substantially from step to step, what would be expected from a recursive calculation. The reason for the small difference between **binary32** and **posit32** could be that the involved real numbers are closer to representable numbers in **posit32** than in **binary32**. Posits have a larger number density around 1 compared to floats of the same bit-length, and therefore more precision in that regime (see Figure 5 for the example of **binary16** versus **posit16**). Posits also have more numbers, because they have no NaNs. Round-off accumulation is specially dependent on the precision of the datatype, which makes sense; the lower the precision, the higher the round-off accumulation. The difference in *A*-orthogonality is negligible for **posit32** compared to **binary32**, but again clearly surpassed by **binary64**.

binary16 did not converge (**status=-1**) after the maximal number of **nmx=256** steps. Its footprint is absent in plot 1d, because it consisted only of NaNs and infinities, causing $\alpha_i = 0$ and $\beta_i = 1$. This implied that $\vec{r}_i = \vec{r}_{i+1}$ and $\vec{p}_{i+1} = \vec{r}_{i+1}$ and therefore $\vec{x}_{i+1} = \vec{x}_i$ and the algorithm stalled. This explains the residues not changing in plots 1a and 1b. The reason for the first infinity was an overflow when calculating the norm of \vec{b} in the very first iteration. This suggests that the limited number range of **binary16** might not be enough (at least for a naive replacement), comparing to **bfloat16** with the same bit-length, but larger number range that was able to converge, although very slowly.

The behaviour of **posit8** is very similar to **binary16**, but without the overflow, because posit do not overflow by definition. Instead the biggest representable number is returned or in case of an underflow the smallest representable number is returned [12]. The algorithm stalled at a value of the norm of the recursive residual of $\|\vec{r}_i\| = 8$. The biggest 8-bit posit number with exponent bits $es = 0$ is $2^6 = 64$, so the norm squared cannot be bigger than 64 and the norm itself cannot be bigger than $8 = \sqrt{64}$ (see plot 1b). This happened in the first step, whereas the actual residual in **binary64** was $\sim 10^3$. The amounts $|\alpha_i| \ll 1$ in iterative steps are therefore very small causing $\vec{x}_{i+1} \approx \vec{x}_i$. Significant changes in \vec{x}_i will not happen and convergence is unlikely. Also notice that **posit8** had 256 reset steps, which means that after every step there was a reset step. The steps were caused by the very high precision limit of **posit8**. The value of `PRECISION LIMIT` is `100*MACHINE EPSILON`, which has a value of 3.125 for **posit8**.

The story of **posit16** is very similar, just that the maximal representable value with $es = 1$ is 268,435,456 and the square root of this is 16,384 which is reached after 8 steps (see plot 1b). The actual residual in the 8-th step was $\sim 10^7$, the algorithm diverged and then stalled. Iterative steps are therefore mostly too small and convergence is unlikely.

We observe that number range is more important than precision, when naively replacing the datatype, but the higher the precision, the faster the convergence and the less reset steps needed.

In Figure 9 the replacement utilised the possibility to use **quires** for the posit runs. Therefore, the numbers for the float datatypes are exactly equal to the ones in Figure 8, because floats have no such feature. They are not discussed again.

Comparing plots 1c and 2c and looking at **posit32**, one can see that the round-off accumulation in the residual due to its recursive calculation is slightly better than without using the **quire**. This makes sense, because **quires** introduce deferred rounding. This is exploited specially in the calculation of norms and matrix-vector products. It also results in a somewhat better maintaining of A-orthogonality for the direction vectors.

However, the data points of **posit16** bear little resemblance to its previous or later runs. It comes much closer to the target residual tolerance than in the last simulation, but it is still not reached. The tolerance is within the number range of **posit16**, even so it did not converge. The reason for this is that the smallest representable number in **posit16** is 2^{-28} . The **quire** for **posit16** has the same number range, despite the 128 bits in length. Every norm squared of a non-zero vector must be larger to equal to this number, because posit do not underflow. Therefore the norm is always larger or equal to $\sqrt{2^{-28}} = 2^{-14} \approx 6.1 \cdot 10^{-5}$. The tolerance of $7.8 \cdot 10^{-5}$ - even though larger than that number - is perhaps still too close. Comparing the **lightpink** values, that are **posit16** as well, but the relative residual **res** is set to 10^{-5} instead (the tolerance being one order of magnitude larger), they converged after only **status=24** steps. This suggests that the reason for the strange behaviour lies in the relative residual that was chosen too close to the lowest number above zero of the number regime.

Using the same arguments and analysis, **posit8** had no chance to give a meaningful result.

In Figure 10, a smarter replacement was done. All variables that have a collective role suffer from overflow. For example the norm of a vector $\vec{v} \in \mathbb{R}^n$ is

$$\|\vec{v}\| = \sqrt{\sum_{i=1}^n v_i^2}.$$

The number below the square root may be much bigger before squaring than after. If we calculate the norm in **posit8**, the result will be $\|\vec{v}\| \leq 8$. More importantly, when using a datatype that overflows such as **binary16**, the value after squaring might be perfectly fine, but the value under the square root could be outside the range of representable numbers, $\sqrt{\infty} = \infty$ and $\sqrt{0} = 0$. This is cured if the collective variable is of a datatype with larger number range than the underlying datatype that is summed over. In Figure 10 all collective variables were of type **binary64**.

The data of **binary64** exhibits no significant alterations. Again comparing **binary32** and **posit32** with their previous data points, we see that the round-off accumulation of **binary32** is a little better and **posit32** is approximately the same as with the **quire**, suggesting that when using posits utilising the **quire** is probably sufficient.

Looking at **tensorfloat32**, it has the same exponent range as **binary32**, but less precision and it has the same number of mantissa bits as **binary16**, but at a higher exponent range. Compared to

binary16, both datatypes have the same amount of numbers to be distributed in their respective number range. It is expected to perform worse or equal to **binary32**, but better or equal to **binary16** and **bfloat16**. Therefore it's expected to converge in $27 \leq \text{status} \leq 28$ steps, see equation (8.16). This is indeed the case with $\text{status}=27$ steps. We see that the larger number range compared to **binary16** has little to do with speed of convergence. This is because the number regime is within the **binary16** regime, except for collective variables. This explains as well why **tensorfloat32** performed precisely as in the naive replacement, Figure 8, but the round-off accumulation is better because of the more precise collective variables.

The **bfloat16** with even less precision but comparable number range of **tensorfloat32** converged in $\text{status}=28$ steps as well, but needed two more reset steps, tightening the previous conclusion about speed of convergence.

The most interesting data points are the ones of **binary16** and **posit16** that both were able to converge in $\text{status}=28$ and $\text{status}=34$ steps respectively. They performed quite similar, even though it would be expected that **posit16** would perform a slightly better because of the bigger number range and bigger number density in relevant number regimes (see Figure 5). In plot 3c the increase of round-off accumulation can be observed for **binary16** and **posit16** in steps where the real residue changes (where the algorithm makes progress, see for example: steps 1 to 10). Notice that, when the real residue stalls and the recursive residue still (wrongly) decreases, the round-off accumulation will saturate until the order of magnitude of the two numbers becomes too large such that their difference is dominated by the larger number. This can be seen in the data points of **posit16** in plot 3a. It suggests that the precision limit was chosen too low for the datatype. Notice that the precision limit is defined to be 100 times the **MACHINE_EPSILON** of the datatype. The **MACHINE_EPSILON** for the posit datatypes is quite misleading, because it gives us (by definition) the precision of numbers around 1. This is the regime where posits are most precise, their precision falling off very rapidly when leaving it. Thus for **posit16** in the regime 10^{-1} the **MACHINE_EPSILON** is correct (seen at iteration 14), whereas in the regime 10^{-3} it is chosen too small and we can see a staircase-shape around the reset steps at iterations 28 and 35. Such a stalling of the real residue should be avoided at any cost, because the algorithm stalls as well in that case. The **MACHINE_EPSILON** is defined to be the difference between 1 and the lowest number above 1. For floats this definition makes more sense, because their precision does not fall off that fast, but for posits which are most precise around 1 this gives a too precise value, not reflecting the real precision of posits in their whole number range correctly. Instead, the machine epsilon should be a function of the number regime, increasing when going far away from 1. This is the reason for the staircase-shaped curve of **posit16** in plot 3a. The phenomenon is even more prominent for **posit16** in plot 4a of Figure 11. The **posit32** does not have this problem, because its **MACHINE_EPSILON** is sufficient for the number regime used in the algorithm. When demanding lower relative residuals, staircase-shapes should be expected for **posit32** as well.

Comparing **binary16** with **bfloat16** and **tensorfloat32**, we see again that exponent range is less relevant than precision. Precision determines the amount of reset steps.

Figure 11 shows all the simulated datatypes using a collective datatype of **binary64** just as in Figure 10, but with a relative residual of 10^{-12} instead. This might be a more realistic scenario. The last row resembles the predicted hierarchy (8.16) particularly well. Notice that 10^{-12} is outside the representable number range of **binary16**, **posit16** and **posit8**. This means that these datatypes have no chance to reach the target tolerance, therefore we expected them not to converge. This is indeed the case. We also see that **binary16** and **posit16** both are not able to go below 10^{-5} , meaning the tolerance in the third row was chosen very close to the minimum possible, but still converging tolerance (see also discussion of **posit16** in Figure 9). Both datatypes make no further significant progress after step 45. It can also be seen that even the recursive residue stalls or increases - an indicator that the datatype has reached its limits.

The comparison between **binary32** and **posit32** is again of insight. Their difference is subtle. We see that both needed the same amount of steps. Round-off accumulation and A -orthogonality are again slightly better, making **posit32** the overall better 32-bit datatype for the problem. The reason for this goes down to the higher precision of posits in the relevant number regime. Looking at the **lightgreen** values, that are **posit32** as well, but utilising the **quire** instead of **binary64** as collective variable, we observe the same amount of steps to convergence, but round-off accumulation is slightly worse. It might be an unfair comparison, because **binary64** as collective variable has more precision, surpassing even the deferred rounding employed by the 512-bit **quire** for **posit32**. In plot 4d the **posit32** with **quire** will not go below some fixed value. The reason for this is the lowest **posit32**

value with exponent bits `es=2` is 8^{-30} and the norm of a `posit32`-vector with at least one non-zero component must be bigger or equal to the square root of this; $1.15 \cdot 10^{-18}$. This suggests that when choosing `res` to be smaller than 10^{-18} , we expect `posit32` not to converge anymore in analogy to `posit16` in the second row.

Since `binary16` was able to converge in Figure 10, this suggests that the number regime is within `binary16` giving `posit32` more precision in that regime over `binary32`.

Finally, compare the 3 datatypes with the same exponent range, but different precisions; `binary32`, `tensorfloat32` and `bfloat16`. The less precision, the slower the convergence. The price to go from 23 to 10 mantissa bits results in 1 more conjugate gradient step as well as 4 more reset steps. When going further down to 7 mantissa bits again 1 more regular step and 4 more reset steps were needed to finally bring `bfloat16` to convergence after `status=56` regular conjugate gradient plus 10 reset steps. Bearing in mind that it uses only 16 bits, this is a remarkable result. It performed way better than its 16-bit competitors.

We also see in plot 4a that all datatypes start to converge by the same speed (all slopes are equal). The actual residual of the datatype with the lowest precision, namely `bfloat16` with 7 mantissa bits, resets first, followed by `binary16` and `tensorfloat32` which have both 10 mantissa bits. The next one is `posit16`, because it has more precision than `binary16` in the relevant regime, followed by `binary32` with 23 mantissa bits and later by `posit32`, where the same argument as before holds. The curve of `binary64` would also reset at some point, but that is outside the scale.

Specially plot 4a suggests that we can start to calculate in a datatype with 16 bits of length until we fall below a constant, to be determined value (that depends on the datatype), then continuing the calculation in a datatype with 32 bit-length until that number regime is exhausted as well, again switching to a 64 bit datatype to finish the calculation.

8.3.2 8^4 lattice

In order to make sure that the previous analysis is consistent and the physics involved were relevant, the same data was extracted from a 8^4 lattice and some of the plots were remade from the new data, see Figure 12. Only the datatypes `binary64`, `binary32` and `binary16` were simulated. In principle, the data tells the same story. The main difference to figures 10 and 11 is that more steps were needed to converge, because the Dirac matrix is much larger than before, although only 0.04% of all components were non-zero, compared to 2% in the 4^4 lattice of the previous analysis. In plots 2a to 2f, where the relative residue was chosen to be 10^{-12} , we again see the saturation of `binary16` marking the lower limit of the datatype. After every reset step, a jump in round-off accumulation can be seen, because the residual in the reset step is calculated in higher precision. It is interesting that the round-off accumulation in the final steps of `binary16` come very close to those of `binary32` (see plot 1c). A reason for this could be the clustering of reset steps just before convergence, giving very accurate results with little round-off, even for less precise datatype. We also see that the speed of convergence does not significantly depend on the precision of the datatype, only the amount of reset step does, thus the less steep slope of `binary16`. When the lower limit of the datatype is reached, the slope becomes zero and the residual shows no striking reduction anymore. This is where the datatype should be switched to one with a larger number range.

8.3.3 Conclusion

The decision between floats and posits is not trivial. It highly depends on how fast the machine can perform **FLOPS** and **POPS**. For example division in floating point arithmetic is very expensive (it may exceed 24 CPU cycles, many compiler optimisations evade them), whereas in posit arithmetic it is said to be cheap, because obtaining the inverse of a number is easy.

Another example could be that comparisons between floats are more expensive than for posits. Two posits are equal if their bit representations are equal. Comparing two floats is much more expensive, mainly because of the many **NaNs** and since 0 and -0 are equal but not bit-identical.

On the other hand, there is currently no hardware available, that has dedicated posit units and posits are not studied as intensive as floats. Floats are widespread, well understood and implemented in common hardware.

If one decides to replace `binary32` with posits, the most elegant solution would be to naively replace the datatype and utilise `quires` in collective operations. To use `binary64` collective variables is not recommended, because this would introduce many type conversions between the floating point

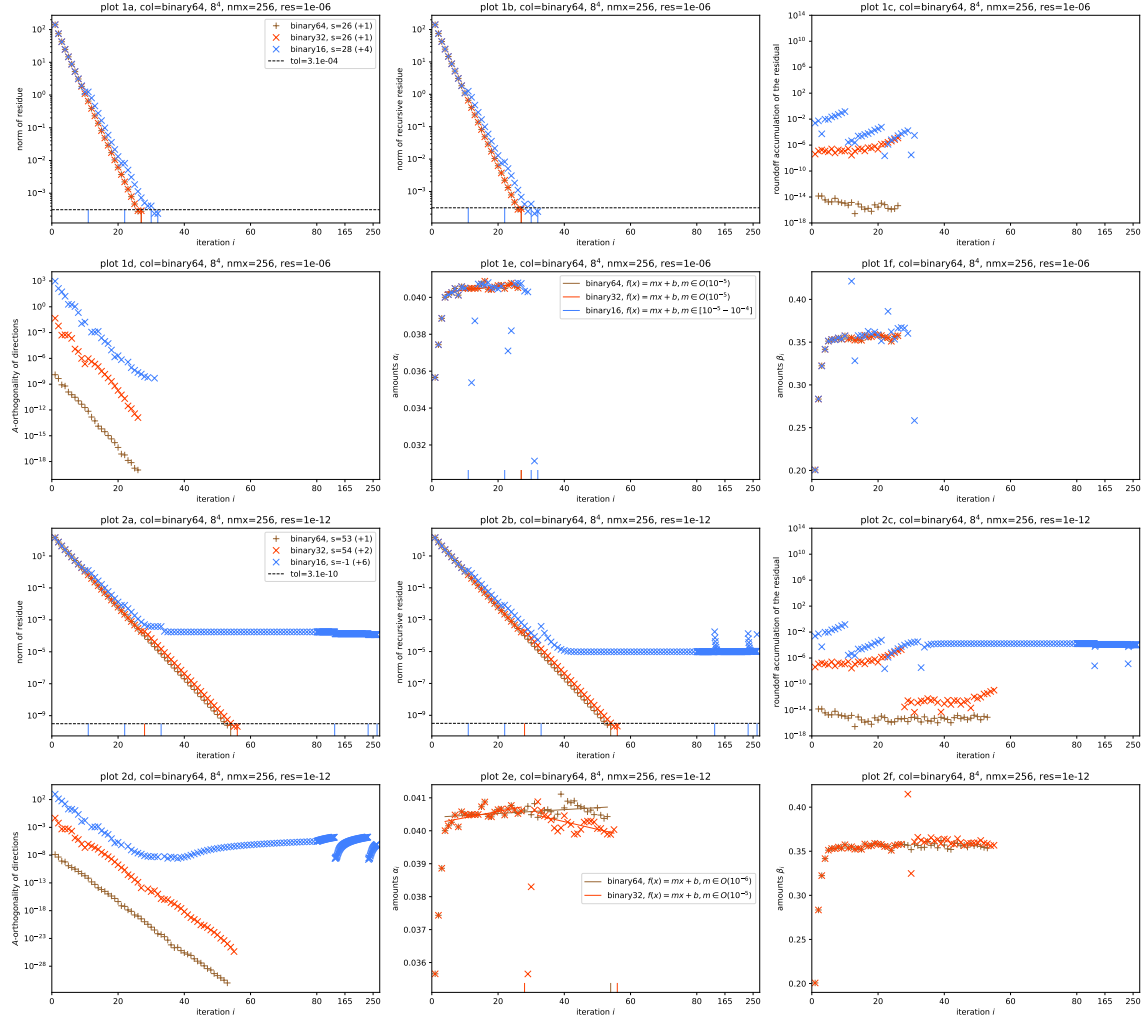


Figure 12: In analogy to figures 10 and 11. This time an 8^4 lattice was used and only the floating point datatypes that are available in hardware nowadays were simulated. The *first and second row* use **binary64** as collective variable and 10^{-6} was the desired relative residual. The *third and fourth row* have the exact same setup, but with a relative residual of 10^{-12} instead.

and the posit format which is assumed to be expensive. The drawback of this method is that `posit16` may only converge if the relative residue is chosen high enough (see plot *2a* in Figure 9).

If the decision goes for floats, which might be the more realistic scenario, then the most elegant solution would be to use collective variables in `binary64`. Type conversions between different IEEE floating point types are not considered to be expensive. The `tensorflow32` compared to `binary32` and `bfloat16` answers the question how important precision is in the calculation. All of them have the same number of exponent bits and therefore approximately the same number range, but very different precisions. We see that all of them were able to converge in any experiment, but with `binary64` as collective variable, the results were closest to each other (see Figure 10 plot *3a*). The only real difference was in the amount of reset steps. If the datatype is lower in bit-length, the memory-boundedness suggests that the calculation performs faster, but the trade-off is the amount of (computationally expensive) reset steps that increases with less precision. However, the datatype for collective operations should be precise and should have a large number range. Since the amount of variables needed in that datatype does not scale with the lattice size, it is perfectly right to use a datatype with large bit-length. Comparing the convergence of `bfloat16` in the naive case (Figure 8 plot *1a*) with the case `binary64` collective variables (Figure 10 plot *3a*), it can be seen that the algorithm converged 21 steps faster, only because the collective datatype was chosen to be `binary64`. On the other hand, comparing the performance of `binary16` in the two plots, we see that the number range of the collective datatype brought `binary16` from no convergence to convergence within `status=35` steps - only marginally slower than `binary32`. These arguments make `binary64` the best choice for variables with a collective role.

Proposal 8.1: Mixed Precision

The above analysis suggests that the calculation of the solution can be (at least partly) conducted in an even less precise datatype than `binary32`. One could for example choose 3 datatypes with different precision. The algorithm can be started using the least precise one. If the tolerance hits a certain value at the boundaries of the datatype, the algorithm switches to the next higher one. The calculation is continued in that datatype until the tolerance reaches the limits of the new datatype. Again the datatype is switched to the next higher one^a. This calculation in mixed precision is not dependent on the algorithm itself and can therefore be applied to every iterative solver algorithm. Algorithm 1 shows an example implementation of such a mixed precision calculation. The array d consists of all available datatypes participating in the calculation in ascending order, meaning the least precise datatype comes first. The function `solve()` performs the underlying algorithm (for example conjugate gradient) in the datatype given by its arguments. It expects at least a starting vector \vec{x}_0 and a tolerance and returns the status ^b, the calculated solution and the

residual up to the given tolerance.

Algorithm 1: Pseudo-code for an iterative algorithm in mixed precision.

```

input: desired norm of relative residual  $rn$ 
input: array of datatypes in  $\{d\}_{k=0}^N$ 
input: iterative algorithm  $solve()$ 
1  $\vec{x}_0, \vec{r}_0, \dots \leftarrow$  initial guess,  $\dots$ ;
2  $\vec{x}, \vec{r} \leftarrow \vec{x}_0, \vec{r}_0$ ;
3  $status \leftarrow 0$ ;
4 for  $k \leftarrow 0, 1$  to  $N$  do
5   convert all variables to datatype  $d[k]$ ;
6    $tol \leftarrow \frac{1}{\|\vec{r}_0\|} \max(rn, \text{MACHINE\_EPSILON of } d[k])$ ;
7    $substatus, \vec{x}, \vec{r}, \dots \leftarrow solve(tol, \vec{x}, \dots)$ ;
8   if  $substatus > 0$  then
9      $status \leftarrow status + substatus$ ;
10  if  $\|\vec{r}\| < rn$  then
11    return  $status, \vec{x}$ ; // success
12 end
13  $status \leftarrow -3$ ;
14 return  $status, \vec{x}_0$ ; // the algorithm failed

```

^aOne obvious choice could be $d = \{\text{binary64}, \text{binary32}, \text{binary16}\}$. When the algorithm is started in **binary16** and a tolerance of $\approx 10^{-4}$ is reached, the algorithm continues in **binary32**, the limit of which is at a tolerance of $\approx 10^{-35}$. A continuing calculation would then be conducted in **binary64**.

^bSee section 8.2

Proposal 8.2: Approximating the amounts α_i

Looking at plot 4e of Figure 11, where the amounts α_i are plotted for every iteration, we see that after every reset step the amounts need 2 – 3 steps to reach a value that is not changing very much for future iterations. This becomes apparent when looking at the fitting lines. The values of the α_i are in the range 10^{-1} and the slopes m of the fitting lines are in the range 10^{-4} - 10^{-5} , suggesting that the value of α_i is not changing from iteration to iteration when only looking at 2 – 3 significant decimal digits.

A possibility to reduce computational cost in each iteration could be to approximate the values of future α_i to be constant. The less precise the datatype, the larger the change in α_i . The large error in α_i of **bfloat16** in all plots suggests that the algorithm is not sensible to errors in α_i . Therefore, it can be expected that the results should not change significantly with a approximated value of α_i .

- Advantage: The residuals can be calculated using $\vec{b} - A\vec{x}$, not recursively. This implies less round-off accumulation.
- Advantage: Only one matrix-vector multiplication per iteration.
- Disadvantage: Since the α_i are just approximated, the number of needed iterations may increase.
- Disadvantage: The Dirac operator D must be given in the form of $A = D^\dagger D$ as *one* operator, else the algorithm still consists of 2 matrix-vector multiplications per iteration. Also, $D^\dagger D$ is less sparse than D .

The results of simulations with approximated values for the α_i can be observed in plot series 13 and 14. The value was approximated based on previous values. The first 5 steps were skipped (thus the algorithm performed natively). In step number 5, the last 3 values of α_i were averaged. In the following steps the constant value calculated in step 5 was reused. After every reset step, the value of α_i had to be recalculated using the above procedure. Therefore a datatype such as **bfloat16** that has reset steps after approximately every 7th

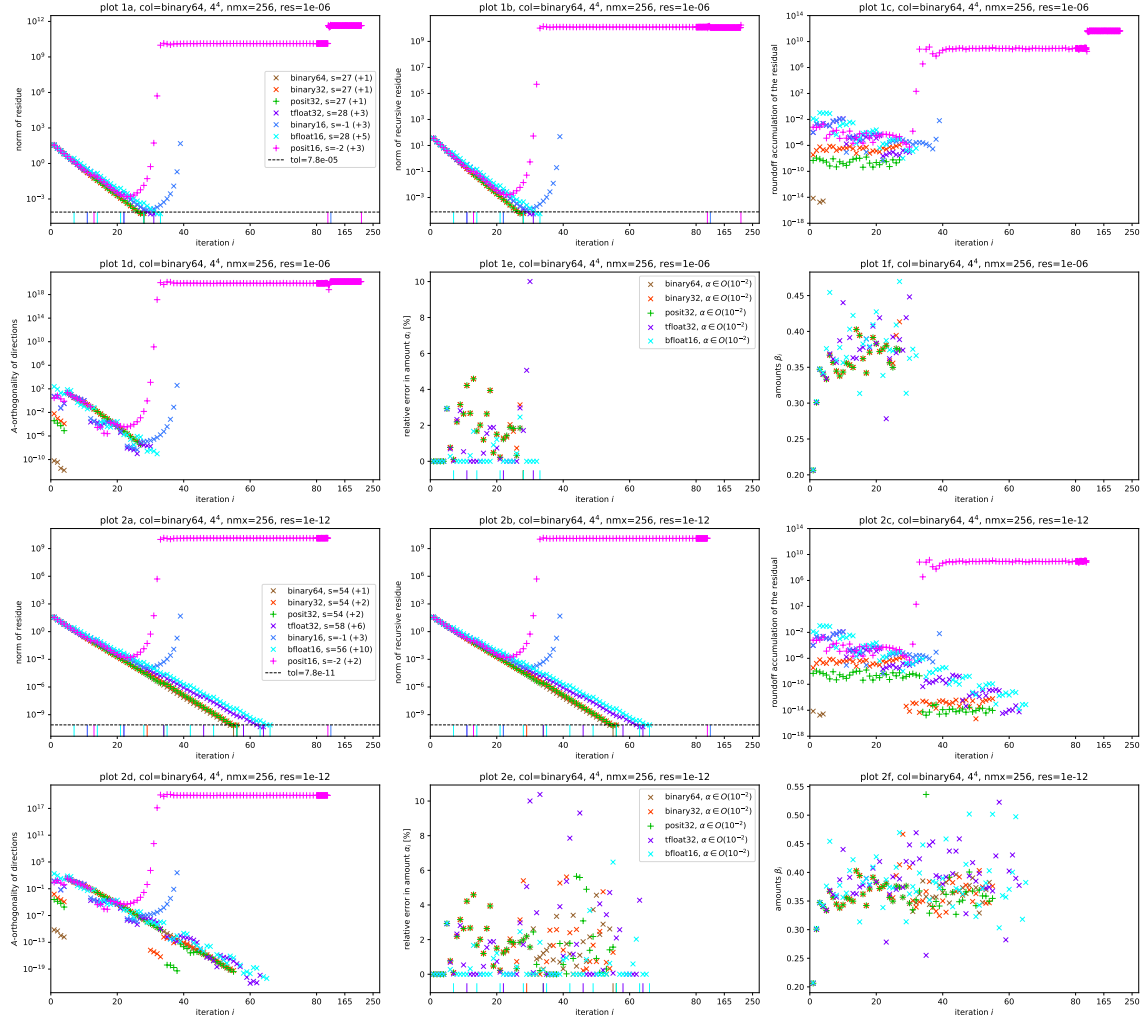


Figure 13: Plots *1a* to *1f* contain the convergence analysis of a conjugate gradient run with a 4^4 lattice, relative residual 10^{-6} and approximated values of α_i . In plots *2a* to *2f* the residual was chosen to be 10^{-12} . Plots *1e* and *2e* contain the relative error in the approximated α_i compared to the real α_i .

regular step, will benefit in only 2 steps per reset step. This is very little difference to native runs compared to datatypes with high precision.

The calculation became more sensible to the number range of the datatype. This can be seen in all plots when looking at `binary16` that was not able to converge anymore, although by a very small amount. `tensorfloat32` on the other hand performed very similar to the regular rounds, it was expected that it needs slightly more iterations. When going with this strategy, it is therefore advisable to perform more regular cg-steps when coming closer to the boundaries of the datatype. One possible solution would be to choose a higher machine epsilon close to the boundaries, forcing the algorithm to perform more reset steps, in turn causing more regular cg-steps and recalculations of α_i .

Notice that with larger lattice sizes, the approximation of the amounts has less error (see plots *1e* and *2e* in figures 13 and 14) and the algorithm is thus more stable.

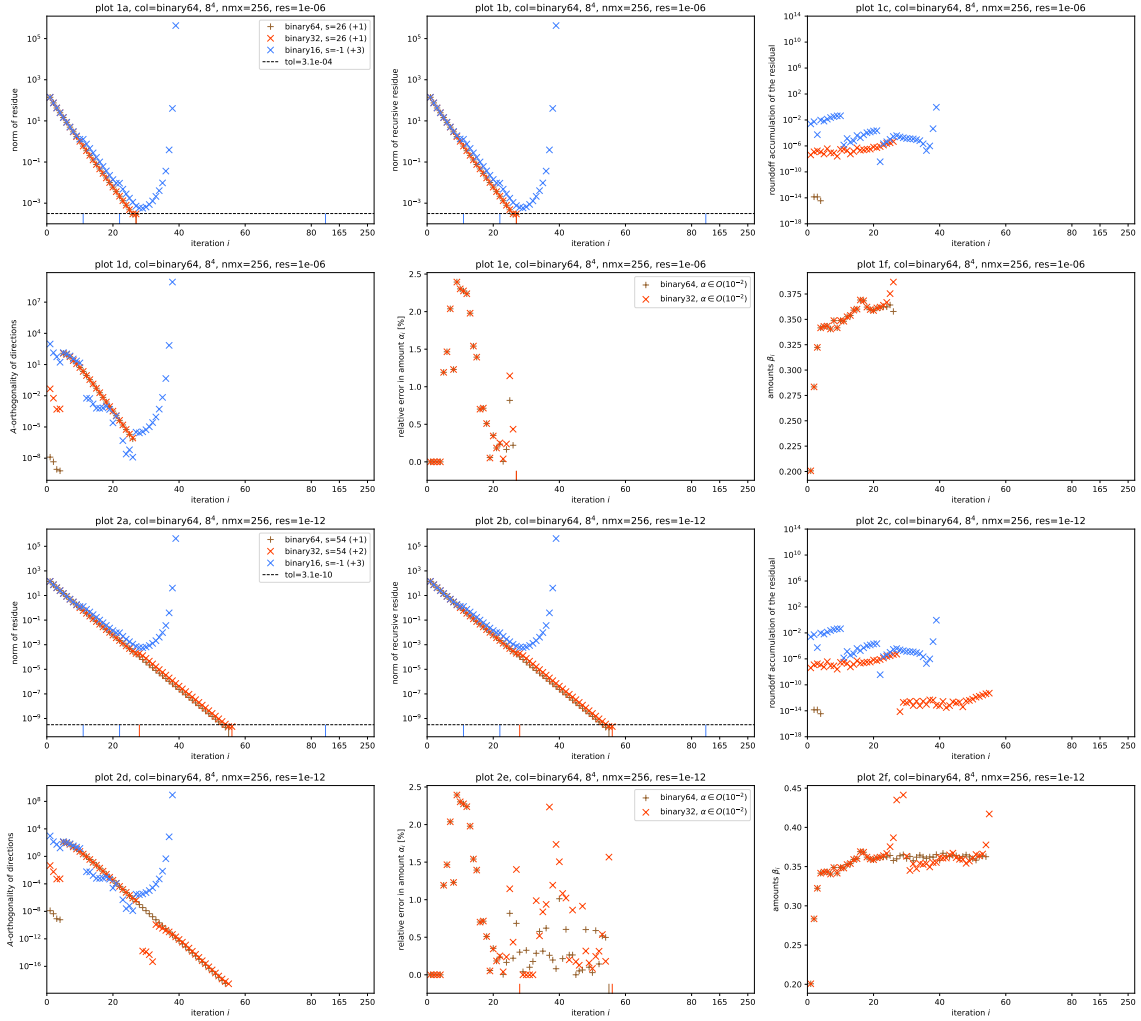


Figure 14: The same setup as figure 13, but with a 8^4 lattice.

9 SAP preconditioned GCR algorithm

The next solver appearing in openQxD is called **SAP_GCR**. It makes use of a multiplicative **Schwarz Alternating Procedure (SAP)** as preconditioner for a flexible **Generalized Conjugate Residual (GCR)** run.

TODO: motivation: parallel processing, chiral regime (spontaneous breaking of chiral symmetry), simulation containing sea-quarks limited to small lattices and large quark masses.

9.1 Even-Odd Preconditioning

Preconditioning in general, when employed in lattice QCD, is expected to have significant impact on the number of iterations of a solver. One way of preconditioning $D\psi = \eta$ on a lattice is

$$LDR\psi' = L\eta,$$

with $\psi = R^{-1}\psi'$ and L, R chosen wisely such that LDR is well conditioned. If $L = \mathbb{I}$, it is called **right preconditioning**, if $R = \mathbb{I}$ it is called **left preconditioning**. If the Dirac-matrix involves only nearest-neighbour interactions it is possible to split the lattice into even and odd sites^{22 23}. If the sites are ordered such that the even sites come first²⁴,

$$D = \begin{pmatrix} D_{ee} & D_{eo} \\ D_{oe} & D_{oo} \end{pmatrix}, \quad \psi = \begin{pmatrix} \psi_e \\ \psi_o \end{pmatrix}$$

D_{ee} (D_{oo}) consists of the interactions of the even (odd) sites among themselves, whereas D_{eo} and D_{oe} consider the interactions of even with odd sites. ψ_e and ψ_o contain the values for even and odd lattice sites of the spinor.

Using specific forms of L and R , D can be brought in a block-diagonal form, namely

$$L = \begin{pmatrix} 1 & -D_{eo}D_{oo}^{-1}D_{oe} \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad R = \begin{pmatrix} 1 & 0 \\ -D_{oo}^{-1}D_{oe} & 1 \end{pmatrix}.$$

After a bit of algebra,

$$LDR = \begin{pmatrix} \hat{D} & 0 \\ 0 & D_{oo} \end{pmatrix}, \quad \text{with} \quad \hat{D} = D_{ee} - D_{eo}D_{oo}^{-1}D_{oe}.$$

This specific preconditioning reduces the amount of iterative steps needed by a factor of 2 approximately, because D_{oo} and \hat{D} are matrices of half the dimension of D . The inversion of D_{oo} is simple, because with only nearest-neighbour-interactions the odd sites do not interact among themselves, only with even sites. Thus D_{oo} exhibits block-diagonal form (all blocks are 6×6 , why?). Using

$$D\psi = \eta \implies \begin{pmatrix} D_{ee} & D_{eo} \\ D_{oe} & D_{oo} \end{pmatrix} \begin{pmatrix} \psi_e \\ \psi_o \end{pmatrix} = \begin{pmatrix} D_{ee}\psi_e + D_{eo}\psi_o \\ D_{oe}\psi_e + D_{oo}\psi_o \end{pmatrix} = \begin{pmatrix} \eta_e \\ \eta_o \end{pmatrix}$$

we can write the preconditioned form, where only the reduced system with even lattice sites has to be solved to determine ψ_e

$$\begin{aligned} \hat{D}\psi_e &= D_{ee}\psi_e - D_{eo}D_{oo}^{-1}D_{oe}\psi_e \\ &= (\eta_e - D_{eo}\psi_o) - D_{eo}D_{oo}^{-1}(\eta_o - D_{oo}\psi_o) \end{aligned}$$

²²It is therefore very similar to a domain decomposition method, see later.

²³Even lattice points are the ones where the sum of the global Cartesian coordinates $(x_0 + x_1 + x_2 + x_3)$ in units of the lattice spacing a is even.

²⁴This is indeed the case in openQxD (see `main/README.global`) in [5].

$$= \eta_e - D_{eo} D_{oo}^{-1} \eta_o,$$

because ψ_o follows from the solution ψ_e via

$$\psi_o = D_{oo}^{-1}(\eta_o - D_{oe}\psi_e).$$

9.2 Schwarz Alternating Procedure

Domain decomposition is a way to partition the large system into (possibly many) smaller sub-problems with regularly updated boundary conditions coming from solutions of neighbouring sub-problems. They fit very well into the notion of parallel processing, because the sub-problem can be chosen to be contained in one single rank. The full lattice is split into sub-lattices called **local lattice**. Each rank has its own local lattice, the size of which is determined at compilation time. The full lattice consists of the ensemble of all local lattices arranged in a grid. It is therefore advisable to choose the size of decomposed sub-domains as a divisor of the local lattice size such that one or more blocks fit into one rank. These sub-problems can then be solved using an iterative solving method.

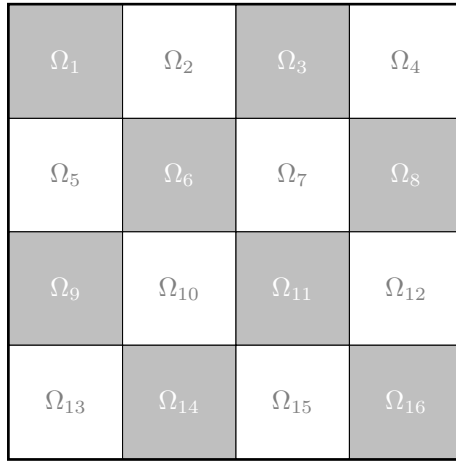


Figure 15: $d = 2$ dimensional example of a decomposition of a lattice Ω into domains named Ω_i .

The idea behind **SAP** is to loop through all blocks Ω_i and solve the smaller sub-problem using boundary conditions given from the most recent global solution (see figure 15). If the original problem only includes nearest-neighbour interactions, the solution of a block Ω_i depends only on that block and its exterior boundary points, which are the adjacent points on the neighbouring blocks with opposite color. For example, the solution of the sub-problem involving Ω_6 , depends only on the solutions of Ω_2 , Ω_5 , Ω_7 and Ω_{10} ²⁵. Therefore all gray (white) sub-problems can be solved simultaneously, with the most recent boundary conditions obtained from the white (gray) domains. Solving all gray, followed by all white sub-problems is called a **Schwarz cycle** and is considered one iteration in the **SAP**. Each sub-problem can be solved with a desired solver separately, again applying some preconditioning²⁶.

²⁵It depends on all other sub-problems as well, but only indirectly.

²⁶Using even-odd preconditioning is perfectly fine with D replaced by the restricted Dirac operator D_i acting only on the points in Ω_i .

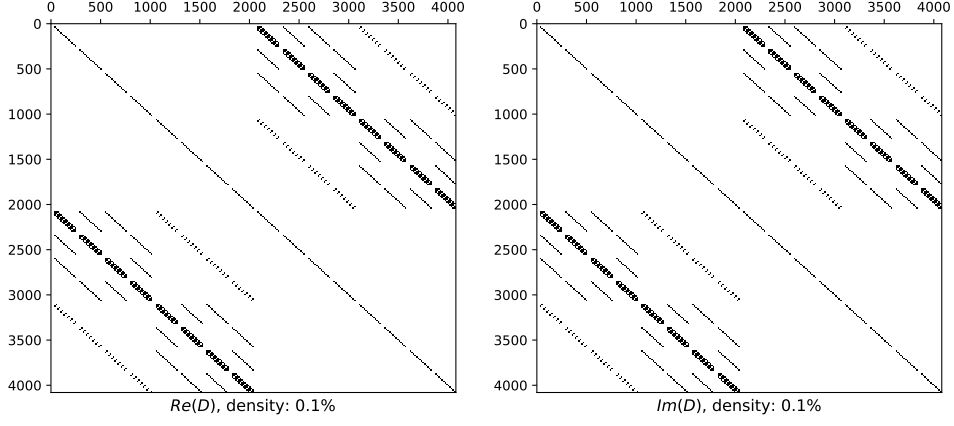


Figure 16: An example plot of a Dirac-matrix of an 8^4 -lattice with SF-boundary conditions. Every pixel consists of 192×192 real numbers. If the average over that numbers is non-zero the pixel is drawn black, else the pixel is drawn white. The density states the percentage of non-zero values over all entries.

Whereas the division into domains on the lattice is straightforward, the representation of the Dirac-operator as a sparse matrix and its decomposition is not. Looking at an actual example of a Dirac-operator as matrix (see figure 16), one observes a lot of structure. While on the diagonal we find the operators restricted to the black and white blocks, the first and the third quadrant describe the operators restricted to the interior and exterior boundaries of the blocks. The operator restricted to the exterior boundaries of the union of all black (white) blocks is denoted by $D_{\partial b}$ ($D_{\partial w}$). The decomposition into $2n$ domains (n gray and n white blocks) can be translated as seen in figure 17. Notice that the restricted operators D_i are well-conditioned, because they exhibit block diagonal form.

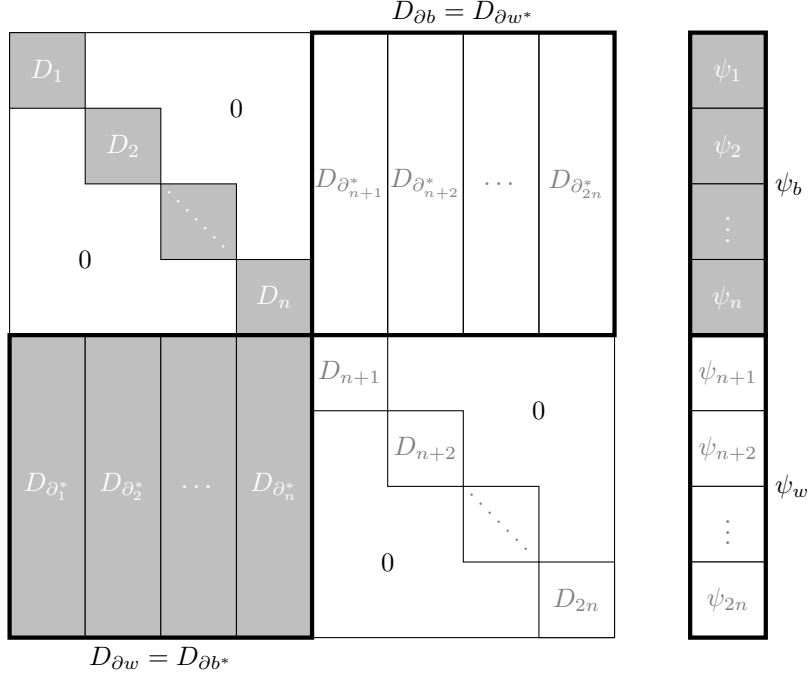


Figure 17: Schematic of the Dirac-operator in term of a large sparse matrix. If the components of the black blocks are arranged such that they appear first, then the decomposition from figure 15 can be translated into a matrix with blocks as in the picture. D_i describes the Dirac-operator restricted to block i and $D_{\partial b}$ ($D_{\partial w}$) is the Dirac-operator restricted to the external boundaries of the black (white) blocks. The color external boundary operators can be decomposed into external boundary operators of the i -th block; $D_{\partial_i^*}$. The right side describes a vector decomposed into the same $2n$ domains ψ_1, \dots, ψ_{2n} . The upper half corresponds to the black blocks and the lower half to the white blocks.

9.3 SAP as a Preconditioner

The multiplicative **Schwarz Alternating Procedure** is such a domain decomposition method coming from the theory of partial differential equations. It can be applied in the form of a right preconditioner M^{-1} making the preconditioned system

$$M^{-1}A\vec{x} = M^{-1}\vec{b} \quad (9.1)$$

to be solved in very few steps, if M^{-1} is a good approximation for A^{-1} . The preconditioning matrix M^{-1} although is never explicitly available during the calculation, such as it is the case in even-odd preconditioning which can also be applied in advance. In order to solve the preconditioned equation (9.1) using an iterative Krylov subspace method, the algorithm must be able to apply M^{-1} and $M^{-1}A$ to an arbitrary vector \vec{v} . If it is possible to implement such operations on multiple ranks in an efficient way and if the preconditioner makes $M^{-1}A$ well conditioned²⁷, we reached the goal. Obviously an application of M^{-1} should be possible without involving A^{-1} . The actions of operators M^{-1} and $M^{-1}A$ on a vector \vec{v} are assembled using a multiplicative **Schwarz Alternating Procedure**, where the blocks are treated by some fixed number of **Minimal Residual (MR)** steps²⁸. The blocks need not to be solved to a certain precision, because the procedure is only used as a preconditioner approximating the solution. This is a motivation for proposal 9.1.

In openQxD the **SAP-GCR** solver is implemented as follows: The large problem is solved using a flexible **GCR** solver, that in each of its **nm**x steps uses a different preconditioner. The preconditioner is given by **ncy** steps of the **Schwarz Alternating Procedure** applied to the current solution vector. Each **SAP** cycle involves approximately solving all gray followed by all white blocks on the whole

²⁷Would it be dingy to expect such a thing from a preconditioner?

²⁸Determined by the value of **nmr** in the solver section of the input file.

lattice each with `nmr` steps of the **MR** method using even-odd preconditioning (`isolv=1`) or not (`isolv=0`).

Proposal 9.1: MR in reduced precision

Since **MR** is memory-bound, it can be conducted in mixed or reduced precision.

Proposal 9.2: Performing the MR steps on the GPU

The preconditioning procedure involves `nmr` **Minimal Residual (MR)** steps to be taken on each block in each Schwarz cycle to approximate a solution to the block problem. Since blocks of the same color are independent on each other and the Dirac operator acting only on a specific block involves no communication whatsoever, we can conclude that the procedure of solving a sub-problem is a problem *local* to the block and self-contained in the sense that it can be solved independently and without MPI communication among ranks. This could be a very handy starting point when going towards GPU-utilisation. Once the source vector and the restricted Dirac operator are transferred to the GPU (both stay constant during the solving process), the problem can be solved on the GPU without involving any communication with other ranks or GPUs. This can also be beneficial, because of the following argument: The local lattice of one single rank, can be subdivided into multiple blocks as well (imagine figure 15 being the local lattice). The actual implementation solves the gray (white) blocks in a local lattice sequentially^a. Since all the gray (white) problems within the local lattice can be solved simultaneously, the code does not exploit the full concurrency potential of the procedure. Solving the sub-problems on the GPU, one could launch **MR** solvers on all gray blocks simultaneously followed by all white blocks. Keeping in mind proposal 9.1, the **MR** solver can be called in mixed or even reduced precision.

For a specific implementation of the GPU-solver, one option is to encode the restricted Dirac-operator in one of the sparse matrix formats (for example **CSR**) and use already existing libraries (for example [2] for **CUDA**) for an application to a spinor. Comparing the implementation of the Dirac-operator in QUDA (see ref. [6]), it is advisable to not rely on such generic libraries, because they ignore further symmetries and structure of the operator. The problem lies mostly in the memory-boundedness of the procedure.

^aBy iterating over the blocks, see `sap()` at line 717ff in `modules/sap/sap.c` in [5].

9.4 Generalised Conjugate Residual algorithm

TODO: Why GCR? it allows inexact preconditioning without compromising the correctness of the solution.

We wish to solve (8.1) if A is not Hermitian. Comparing to the conjugate gradient algorithm, we minimise the residual \vec{r} of the solution \vec{x} , using the *quadratic form*,

$$\begin{aligned} f(\vec{x}) &= \frac{1}{2} (\vec{b} - A\vec{x})^\dagger (\vec{b} - A\vec{x}) + c \\ &= \frac{1}{2} \|\vec{b} - A\vec{x}\|^2 + c \\ &= \frac{1}{2} \|\vec{r}\|^2 + c. \end{aligned}$$

where $c \in \mathbb{C}$. When taking the derivative of this function with respect to \vec{x} , we find that

$$f'(\vec{x}) = A^\dagger A\vec{x} - A^\dagger \vec{b}.$$

Lemma 9.1 (Uniqueness of the solution). *The solution \vec{x} in equation (8.1) is unique and the global minimum of $f(\vec{x})$, if A is non-singular.*

Proof. Let us rewrite $f(\vec{p})$ at an arbitrary point $\vec{p} \in \mathbb{C}$ in terms of the solution vector \vec{x} ,

$$\begin{aligned}
f(\vec{p}) &= \frac{1}{2} (\vec{b} - A\vec{p})^\dagger (\vec{b} - A\vec{p}) + c + f(\vec{x}) - f(\vec{x}) \\
&= f(\vec{x}) + \frac{1}{2} \vec{p}^\dagger (A^\dagger A) \vec{p} - \frac{1}{2} (A\vec{p})^\dagger \vec{b} - \frac{1}{2} \vec{b}^\dagger (A\vec{p}) + \frac{1}{2} \vec{b}^\dagger \vec{b} \\
&= f(\vec{x}) + \frac{1}{2} (\vec{p} - \vec{x})^\dagger (A^\dagger A) (\vec{p} - \vec{x}) + \frac{1}{2} (A\vec{p})^\dagger (\textcolor{brown}{A}\vec{x}) + \frac{1}{2} (\textcolor{brown}{A}\vec{x})^\dagger (A\vec{p}) - \frac{1}{2} (\textcolor{brown}{A}\vec{x})^\dagger (\textcolor{brown}{A}\vec{x}) \\
&\quad - \frac{1}{2} (A\vec{p})^\dagger \vec{b} - \frac{1}{2} \vec{b}^\dagger (A\vec{p}) + \frac{1}{2} \vec{b}^\dagger \vec{b} \\
&= f(\vec{x}) + \frac{1}{2} (\vec{p} - \vec{x})^\dagger (A^\dagger A) (\vec{p} - \vec{x})
\end{aligned}$$

where to obtain the last line, $\textcolor{brown}{A}\vec{x} = \vec{b}$ as used, thus the term simplified.

In the new form of $f(\vec{p})$, one can directly see that, \vec{x} must minimise the function:

$$\begin{aligned}
f(\vec{p}) &= f(\vec{x}) + \frac{1}{2} (\vec{p} - \vec{x})^\dagger (A^\dagger A) (\vec{p} - \vec{x}) \\
&= f(\vec{x}) + \frac{1}{2} \underbrace{\|A(\vec{p} - \vec{x})\|^2}_{> 0 \text{ for } \vec{p} \neq \vec{x}}.
\end{aligned} \tag{9.2}$$

Therefore \vec{x} is the global unique minimum if A is non-singular. \square

Remark. Notice the similarity of the above equation (9.2) to the analogue of the conjugate gradient algorithm (8.2). The only difference is the substitution of $A \mapsto A^\dagger A$. It is therefore advisable in the derivation of an algorithm to require the directions \vec{p}_i to be $A^\dagger A$ -orthogonal instead of A -orthogonal.

In the same manner as in the derivation of the method of conjugate gradient, we impose a iterative **step equation** to be

$$\vec{x}_{i+1} = \vec{x}_i + \alpha_i \vec{p}_i,$$

again with **directions** \vec{p}_i and **amounts** α_i that have to be determined. The recursively calculated **residual** has again the same formula

$$\vec{r}_{i+1} = \vec{r}_i - \alpha_i A\vec{p}_i.$$

Imposing $A^\dagger A$ -orthogonality instead of regular A -orthogonality between error \vec{e}_{i+1} and direction \vec{p}_i ,

$$\begin{aligned}
0 &\stackrel{!}{=} \vec{e}_{i+1}^\dagger (A^\dagger A) \vec{p}_i \\
&= (\vec{e}_i + \alpha_i \vec{p}_i)^\dagger A^\dagger A \vec{p}_i
\end{aligned}$$

gives an expression for the amounts α_i . Notice the above equation is equivalent to imposing A -orthogonality $0 = \vec{r}_{i+1}^\dagger A\vec{p}_i$. However, we find (compare equation (8.9))

$$\alpha_i = \frac{\vec{r}_i^\dagger (A\vec{p}_i)}{\vec{p}_i^\dagger (A^\dagger A) \vec{p}_i} = \frac{\vec{r}_i^\dagger (A\vec{p}_i)}{\|A\vec{p}_i\|^2}.$$

The **GCR** algorithm does store all previous direction \vec{p}_i as well as $A\vec{p}_i$ in contrast to conjugate gradient. Thus the derivation changes slightly. Let's continue with the determination of the

directions using **Gram-Schmidt orthogonalization** by imposing $A^\dagger A$ -orthogonality instead of A -orthogonality and without imposing all previous β_{ij} to be zero (see definition 8.4). Likewise, we set $\vec{u}_i = \vec{r}_i$ and find

$$\begin{aligned}\vec{p}_0 &= \vec{r}_0 \\ \vec{p}_{i+1} &= \vec{r}_{i+1} + \sum_{j=0}^i \beta_{ij} \vec{p}_j,\end{aligned}$$

with

$$\beta_{ij} = -\frac{\vec{r}_{i+1}^\dagger A^\dagger A \vec{p}_j}{\vec{p}_j^\dagger A^\dagger A \vec{p}_j} = -\frac{(A \vec{r}_{i+1})^\dagger (A \vec{p}_j)}{\|A \vec{p}_j\|^2}.$$

Using the above equations, we find the final form of the **Generalised Conjugate Residuals Method**.

Definition 9.1 (Generalised Conjugate Residuals Method). *The iteration step equation of the Generalised Conjugate Residuals Method is defined as*

$$\vec{x}_{i+1} = \vec{x}_i + \alpha_i \vec{p}_i, \quad (9.3)$$

with

$$\vec{r}_{i+1} = \vec{r}_i - \alpha_i A \vec{p}_i, \quad \alpha_i = \frac{\vec{r}_i^\dagger (A \vec{p}_i)}{\|A \vec{p}_i\|^2}, \quad (9.4)$$

$$\vec{p}_{i+1} = \vec{r}_{i+1} + \sum_{j=0}^i \beta_{ij} \vec{p}_j, \quad \beta_{ij} = -\frac{(A \vec{r}_{i+1})^\dagger (A \vec{p}_j)}{\|A \vec{p}_j\|^2}, \quad (9.5)$$

and initial starting vectors

$$\begin{aligned}\vec{x}_0 &= \text{arbitrary starting point}, \\ \vec{p}_0 &= \vec{r}_0 = \vec{b} - A \vec{x}_0.\end{aligned}$$

There are some remarks to note about the method of **GCR**.

Remark. After calculating \vec{r}_{i+1} and $A \vec{r}_{i+1}$, we can recursively determine $A \vec{p}_{i+1}$ via

$$A \vec{p}_{i+1} = A \vec{r}_{i+1} + \sum_{j=0}^i \beta_{ij} A \vec{p}_j. \quad (9.6)$$

This limits the number of matrix-vector products to one per iteration.

Remark. All previous \vec{p}_i and $A \vec{p}_i$ need to be stored in memory in order to construct the next \vec{p}_{i+1} and $A \vec{p}_{i+1}$.

Remark. Comparing to the conjugate gradient algorithm, we imposed $A^\dagger A$ -orthogonality of the directions \vec{p}_i instead of A -orthogonality as well as A -orthogonality of \vec{r}_{i+1} and \vec{p}_i instead of regular orthogonality. A vanishing of all previous β_{ij} on the other hand was not imposed, leading to the sum in the step equation of \vec{p}_{i+1} .

9.5 GCR in openQxD

The actual implementation of the **GCR** algorithm in openQxD is quite different²⁹, but actually equivalent to definition 9.1 (see lemma 9.2). Ref. [15] explains the implementation of the algorithm in detail. The main **GCR**-loop looks as in Algorithm 2 (see Figure 3 in [15])

²⁹Actually called GMRES recursive algorithm (GMRESR) [22], see `fgcr()` in `modules/linsolv/fgcr.c` lines 212ff in [5].

Algorithm 2: Pseudo-code for the GCR recursion.

```

1  $\rho_0 = \eta$  ;
2 for  $k \leftarrow 0, 1, 2$  to  $n_{kv}$  do
3    $\phi_k = M_{sap}\rho_k$  ;
4    $\chi_k = D\phi_k$  ;
5   for  $l \leftarrow 0$  to  $k-1$  do
6      $a_{lk} = (\chi_l, \chi_k)$  ;
7      $\chi_k = \chi_k - a_{lk}\chi_l$  ;
8   end
9    $b_k = \|\chi_k\|$  ;
10   $\chi_k = \frac{\chi_k}{b_k}$  ;
11   $c_k = (\chi_k, \rho_k)$  ;
12   $\rho_{k+1} = \rho_k - c_k\chi_k$  ;
13 end

```

In algorithm 2, M_{sap} is the **SAP** preconditioner, that might depend on the iteration number k as well, making the algorithm flexible. D is the Dirac-operator and ρ_k the residual in the k -th step. The algorithm does not include an update of the solution vector ψ_{k+1} , instead this is done after n_{kv} iterations all at once,

$$\psi_{k+1} = \sum_{l=0}^k \alpha'_l \phi_k. \quad (9.7)$$

Lemma 9.2. *The iterative algorithm from definition 9.1 is equivalent to algorithm 2 when setting the preconditioning operator $M_{sap} = \mathbb{I}$, the Dirac-matrix $D = A$, the source vector $\eta = \vec{b}$ and the solution vectors $\psi_k = \vec{x}_k$.*

Proof. Noticing that the residual $\rho_k = \vec{r}_k$ from line 12 in algorithm 2 and in definition 9.1 must be identical, we find that χ_k must be proportional to $A\vec{p}_k$. Before the normalisation in line 10, we have $\chi_k = A\vec{p}_k$. The $b_k = \|\chi_k\|$ are set before normalisation of χ_k , therefore $b_k = \|\chi_k\| = \|A\vec{p}_k\|$. Using this we find $a_{lk} = (\chi_l, D\rho_k)$ and since $l < k$ the χ_l are normalised, thus $\chi_l = b_l A\vec{p}_l$ after line 10. Thus $a_{lk} = (A\vec{p}_l, D\rho_k)/b_l = -\beta_{k-1,l}\|A\vec{p}_l\|$. Finally, the c_k are defined after normalisation of the χ_k , therefore they evaluate to $c_k = (\chi_k, \rho_k) = (A\vec{p}_k, \vec{r}_k)/b_k = \alpha_k\|A\vec{p}_k\|$. Using these substitutions we find the same formulas as in definition 9.1, except for the step equation.

The main difference between the step equations (9.3) and (9.7) is that in the former the solution \vec{x}_{i+1} is spanned by the direction vectors \vec{p}_i , whereas in the latter it is spanned by the residuals $\rho_i = \vec{r}_i$. This is not a problem since both sets of vectors span the same space, but the amounts α'_i in equation (9.7) differ heavily from the amounts α_i in equation (9.4).

To determine the amounts α'_i in terms of α_i and β_{ij} , we notice equation (9.6),

$$A\vec{p}_i = A\vec{r}_i + \sum_{j=0}^{i-1} \beta_{i-1,j} A\vec{p}_j \iff b_i \chi_i = D\rho_i - \sum_{j=0}^{i-1} a_{ji} \chi_j \quad (9.8)$$

and the fact that

$$\rho_{k+1} = \eta - \sum_{l=0}^k c_l \chi_l. \quad (9.9)$$

But also

$$\rho_{k+1} = \eta - D\psi_{k+1}$$

$$\begin{aligned}
&= \eta - \sum_{l=0}^k \alpha'_l D \rho_k \\
&= \eta - \sum_{l=0}^k \alpha'_l \left[b_k \chi_k + \sum_{j=0}^{k-1} a_{jk} \chi_j \right], \tag{9.10}
\end{aligned}$$

where in the last step equation (9.8) was inserted. The $\chi_i \propto A \vec{p}_i$ are linearly independent, thus the coefficients from (9.10) can be compared to (9.9), giving for $m = 0, 1, \dots, k$

$$\begin{aligned}
\alpha'_m &= \frac{1}{b_m} \left[c_m + \sum_{l=m+1}^k \alpha'_l a_{ml} \right] \\
&= \alpha_m - \sum_{l=m+1}^k \alpha'_l \beta_{l-1,m}.
\end{aligned}$$

□

Proposal 9.3: GCR in mixed precision

In the current version of openQxD [5], the outer **GCR** solver is performed in pure **binary64**. A mixed precision variant would need the preconditioning M_{sap} to be done in mixed precision as well. Algorithm 1 would directly apply with *solve()* replaced by **fgcr()** with the difference that **fgcr()** has to accept D , M_{sap} , \vec{x}_0 and \vec{b} in the desired precision.

9.6 Simulating SAP_GCR

The complete **SAP_GCR** kernel was implemented using Python in the exact same way as the **fgcr()** function from the source code³⁰. The Dirac operator **Dop_double()** was extracted in the same way as for the **cgne()** kernel previously (see section 8.3) using the same configuration. The python implementation contains a floating point datatype for the reduction variables separately (**rdtype**). It also accepts a "large" datatype (**ldtype**) by which the restart steps are calculated in and a "small" datatype (**sdtype**) in which the regular and the **MR** steps are performed in. The result is obtained in terms of the "large" datatype. There are various configuration settings to choose from (see table 4).

setting	meaning	comment
res	desired relative residual	
nmx	maximal number of GCR steps	
nkx	number of generated Krylov vectors until restarting the algorithm	
ncy	number of SAP-cycles to perform in each iteration	
nmx	number of MR-steps to perform on each block in each SAP-cycle	
bs	block size	
ldtype	"large" datatype	can be binary64 or binary32
rdtype	reduction datatype	
sdtype	"small" datatype	

Table 4: Settings for SAP_GCR and their meanings.

The possible datatypes for **ldtype**, **rdtype** and **sdtype** are **binary64** and **binary32**. Unfortunately there was no possibility to use **binary16**, **bfloat16** or **tensorflow32**, even though modern GPUs

³⁰See line 212ff in `modules/linsolv/fgcr.c` in [5].

such as the one tested on support these datatypes, because the tensor-cores are not able to accelerate sparse matrix-vector products.

The following plot series should give an estimate on how much speed improvement can be expected for a GPU-implementation of the solver algorithm. The results should also hint on how to optimally choose the (many) parameters for the solver. It has to be kept in mind that the transfer of the Dirac-operator (full, boundary or blocked) is not part of the time measurements; it is assumed that the operators already reside on the correct places (CPU memory or GPU memory), only spinors are transferred back and forth. Figures 18 - 21 contain the measurements. Every data-point represents the average of at least 20 runs of the `SAP_GCR` kernel in the given configuration and Dirac-operator. The y-axis denotes the duration in seconds and the x-axis shows the configuration (`ncy`, `nmr`) as well as the block size (`bs`) increasing in computational effort per GCR-step from left to right.

Two configurations are non-standard; $(n_{cy}, n_{mr}) = (0, 0)$ and "adaptive". The former indicates no preconditioning (thus a pure `GCR` run) and in the latter configuration the parameters n_{cy} , n_{mr} where chosen automatically in every iteration anew, the block size was chosen to be the largest possible. The adaptive choice was done as follows: If - after a Schwarz-cycle - the norm of the residual is not lower than the previous residual norm, the preconditioning phase is exited. Thus, at least one Schwarz-cycle is performed. A similar strategy is applied to determine the number of MR-steps: if - after at least 4 MR-steps - the norm of the blocked residual is larger than 0.9 times the previous residual norm, the MR-solver of that block exits. So, every block is treated differently in every cycle. A maximum of 20 Schwarz-cycles and 20 MR-steps on each block would be performed if the above exit conditions will never kick in. Therefore the adaptive version tries to find the optimal configuration in every step and every block separately.

The colors denote the datatype setup (`ldtype`, `rdtype`, `sdtype`) and the marker symbols indicate whether the calculation was performed purely on the CPU (circles; \circ , \circ , \circ), purely on the GPU (crosses; \times , \times , \times) or a hybrid variant, where only the MR-steps are calculated on the GPU and the remainder on the CPU, see proposal 9.2 (diamonds; \diamond , \diamond , \diamond). All combinations of the above configurations are present in the plots. The different plots show results from different matrices. 2 matrices where extracted directly from a run of openQxD (`sf_no_cstar_8x8x8x8` and `sf_no_cstar_16x16x16x16`), whereas 2 further matrices where taken from [8]. The matrix `conf6_0-8x8-2` taken from [8] has a parameter $0 \leq k \leq k_c$. The closer k is to its critical value k_c the worse the matrix is conditioned. The used values for k and k_c are given in the title of the plots. The relative residual was chosen to be 10^{-6} and the number of GCR-steps until restart `nkx=7`.

9.6.1 Discussion of figures 18 - 21

An apparent trend visible in all plots is that the pure-GPU variants are most efficient when the block size is large (the number of blocks is small). The pure-CPU variants behave different - the block size has less influence on the run-time. In general, the pure-GPU variants are faster than the pure-CPU ones. This comes from the fact that they take advantage of concurrency. The hybrid-variants are as expected in-between them.

As further general observation, the power of the `SAP_GCR` solver is only fully unfolded if the condition number of the operator is large; the plots of `sf_no_cstar_8x8x8x8` and `sf_no_cstar_16x16x16x16` show no significant performance improvement of the SAP-preconditioning compared to a pure GCR solver without preconditioning ³¹, whereas the runs on `conf6_0-8x8-2` do.

An analysis of the different datatypes shows that the general trend is the lower the involved datatypes are in bit-length, the faster the solution is obtained, which makes sense in memory-bound problems. The setups with `binary32` in reduction variables and as "small" datatype appear to be the most efficient. On the CPU they display on average a speedup compared to the case where only `binary64` was used of $S = 1.246$. However, of the given datatype setups one should choose the one where the datatype of reduction variables (`rdtype`) is set to `binary64` preventing over- or underflows. This was already discussed in section 8.3.

Looking at the fist plot (`sf_no_cstar_8x8x8x8`), as expected the preconditioning gives no significant improvement; on the CPU only 4 setups were faster than the one without preconditioning, on the GPU even none of the preconditioning setups beats the trivial case. Of the CPU ones that were faster than the trivial case, all had the same configuration $(n_{cy}, n_{mr}) = (1, 4)$, but different block

³¹The runs with configuration $(n_{cy}, n_{mr}) = (0, 0)$

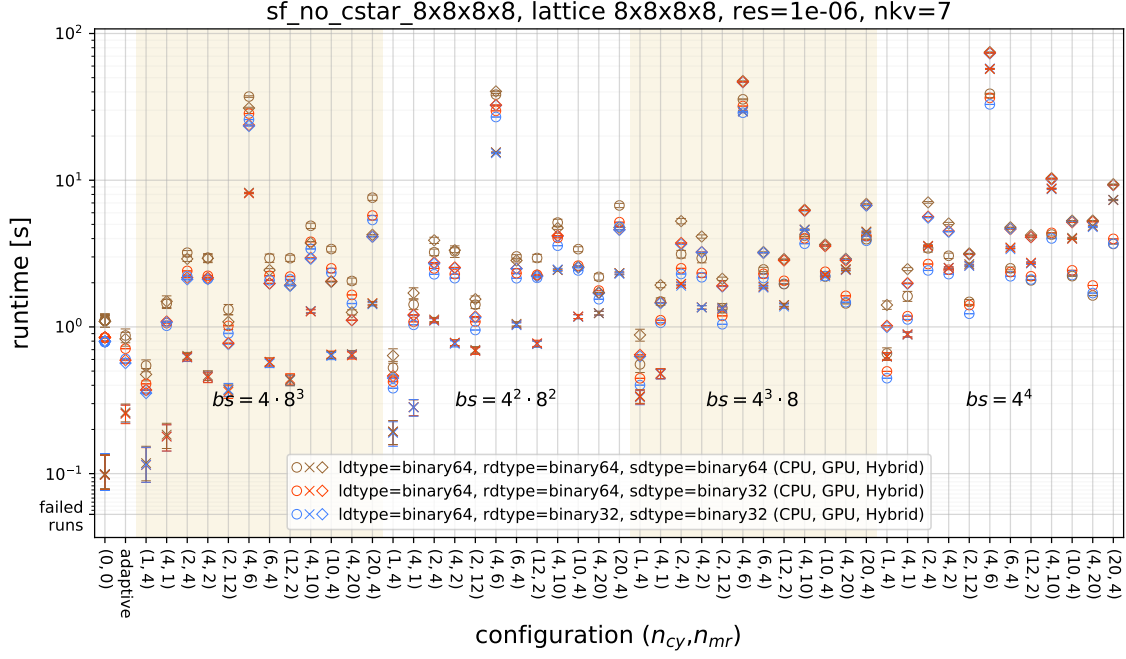


Figure 18: Time measurements for the SAP_GCR kernel on different matrices and configurations. The measurements were conducted on an AMD EPYC 7742 CPU @ 2.25GHz with 512 GB memory and an NVIDIA A100 SXM4 GPU with 40 GB memory.

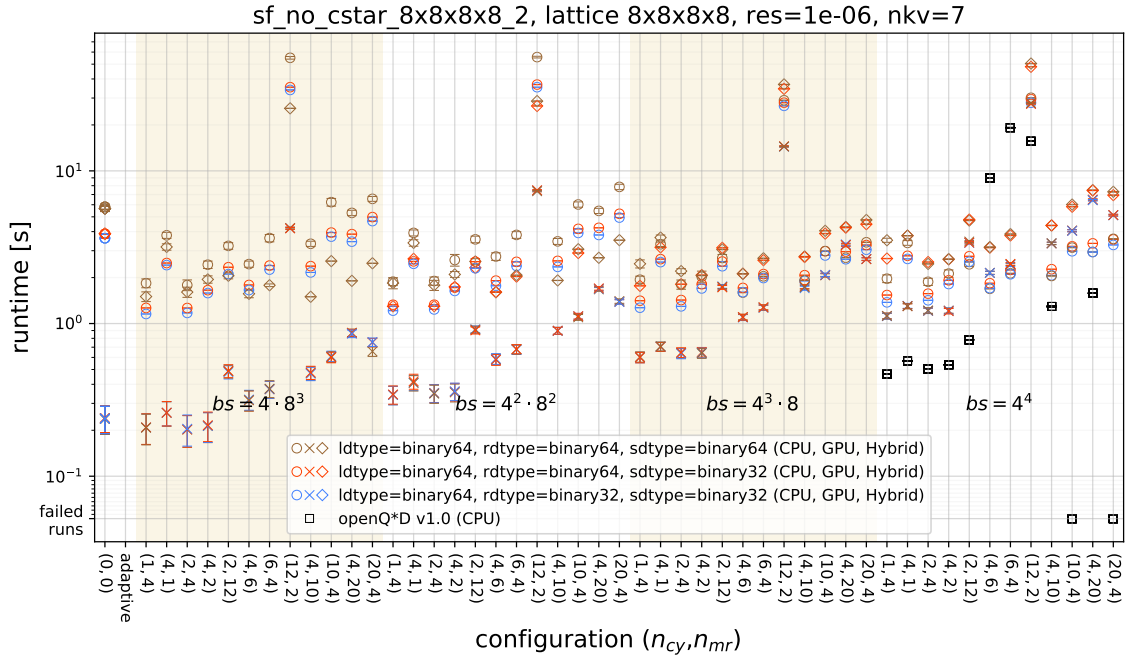


Figure 19: Time measurements for the SAP_GCR kernel on different matrices and configurations. The measurements were conducted on an AMD EPYC 7742 CPU @ 2.25GHz with 512 GB memory and an NVIDIA A100 SXM4 GPU with 40 GB memory.

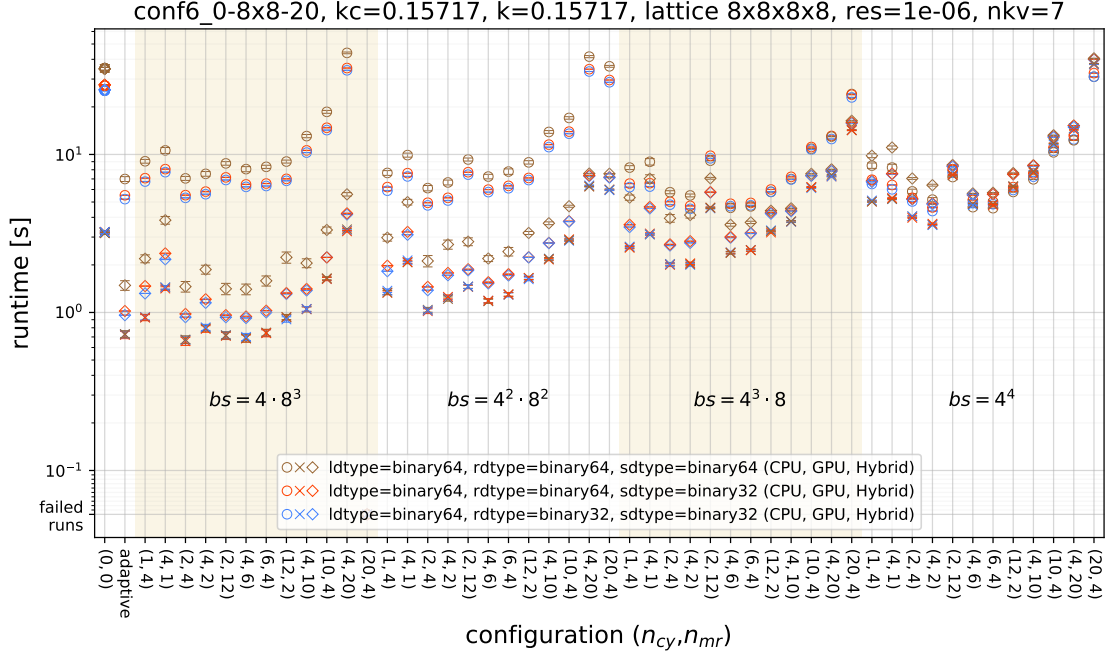


Figure 20: Time measurements for the SAP_GCR kernel on different matrices and configurations. The measurements were conducted on an AMD EPYC 7742 CPU @ 2.25GHz with 512 GB memory and an NVIDIA A100 SXM4 GPU with 40 GB memory.

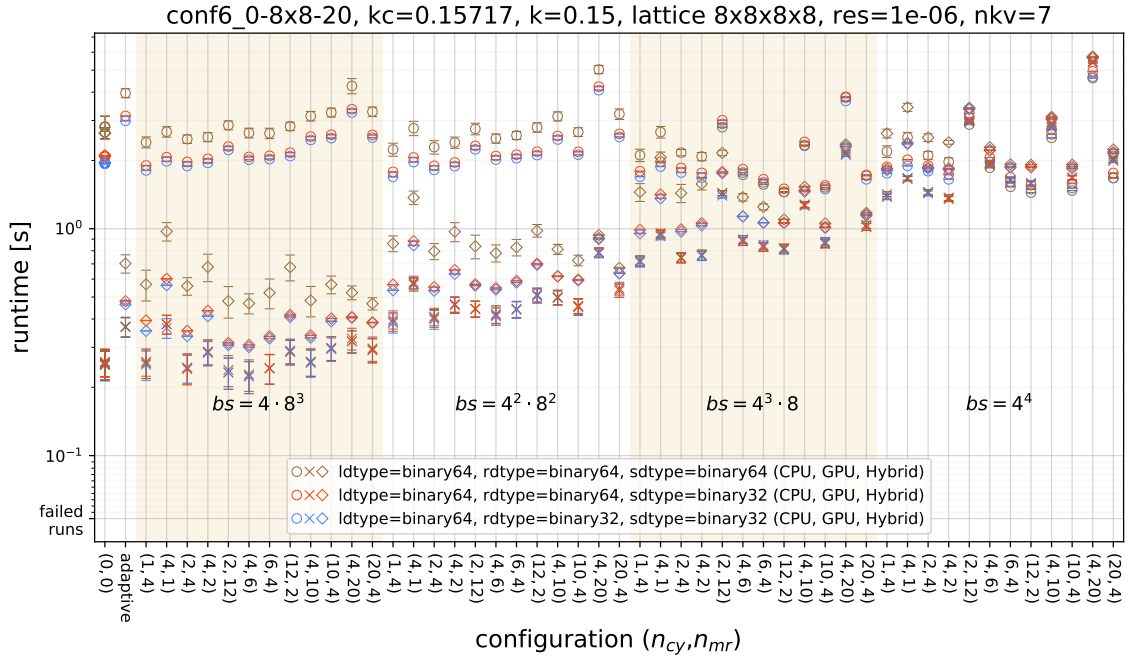


Figure 21: Time measurements for the SAP_GCR kernel on different matrices and configurations. The measurements were conducted on an AMD EPYC 7742 CPU @ 2.25GHz with 512 GB memory and an NVIDIA A100 SXM4 GPU with 40 GB memory.

sizes. This shows that if the operator is well-conditioned, too much preconditioning worsens the performance. $(n_{cy}, n_{mr}) = (1, 4)$ is the configuration with the least amount of preconditioning. The CPU run-time shows a strong dependence on the configuration; there are even certain configurations (for example $(n_{cy}, n_{mr}) = (4, 6)$) that are more than 40 times slower than the no preconditioning case. A wrong choice of configuration parameters can thus lead to a significant performance decrease, whatever exactly "wrong" means in this context. The adaptive choice of parameters should here come to the rescue.

The operator `sf_no_cstar_16x16x16x16` has a very similar behaviour as `sf_no_cstar_8x8x8x8`, because it is well-conditioned as well. The very same configurations give speedup compared to the case without preconditioning, this time both the CPU and GPU variants improved. The claim that the pure-GPU variant slows down with smaller block sizes is even more visible in this plot, since there are 8 block sizes to work on. Looking at the behaviour within a certain constant block size, the algorithm seems to be very sensitive to changes in the configuration. As an example block size $bs = 8^2 \cdot 16^2$, one would expect the run-time to either increase or decrease with respect to the amount of preconditioning. The actual results show a more complex dependence with an exceptional case at $(n_{cy}, n_{mr}) = (2, 8)$. That exceptional case can be seen with all block sizes and on the CPU, GPU as well as in the hybrid case. The plot of `sf_no_cstar_8x8x8x8` features such an exceptional case as well at $(n_{cy}, n_{mr}) = (4, 6)$. However, with both matrices any preconditioning makes the run-time worse, they might not be very representative.

Continuing to the matrix `conf6_0-8x8-2` where the condition number depends on how close the $k = 0.15$ parameter is at its critical value $k_c = 0.15717$. This is the regime where the preconditioning shows benefits. For the pure-CPU cases, we see no strong dependence on the amount of preconditioning, but on the block size. Small block sizes seem to be beneficial, whereas the pure-GPU variant prefers large block sizes. Similar to the above analysis, the good condition number again give not much speedup gain, when comparing the preconditioned cases with the trivial case.

Using the same matrix as above, but at the critical point $k = k_c$, we are in the regime where the **SAP_GCR** algorithm shows its true potential; nearly all cases performed better than the trivial case without any preconditioning. The pure-GPU cases behave as usual - large block sizes are better. This time even the pure-CPU shows a dependence on the block size. Although, it's a weak dependence, but with smaller block sizes the CPU seems to perform slightly better (on contrary to the GPU). The hybrid cases - as usual in-between - are closer to the pure-GPU ones, because despite being hybrid most of the work is done on the GPU. the pattern within a certain block size is repeating and the best amount of preconditioning seems to be at $(n_{cy}, n_{mr}) = (4, 6)$.

9.6.2 Conclusion

In general two different patterns can be observed in the above plot series. Either the matrix is bad conditioned and one can see that most of the preconditioned runs perform better than no preconditioning. Or the matrix is good conditioned and the preconditioned cases perform worse. In both cases the pattern within a block size repeats in other blocks sizes, but shifted. We see two types of patterns. For bad conditioned systems the pattern shows that the ideal case is at some point in the middle. On the other hand for well-conditioned system, the pattern seems random, but a increase of run-time with more preconditioning (higher values of n_{cy} and n_{mr}) can be seen. Since the systems are already well conditioned, too much preconditioning can worsen the run-time.

In openQxD, every rank has its own local lattice to process. The solver algorithms are implemented in a parallel manner, such that the solver is called on all ranks simultaneously. The simulation problem was chosen to be such that the task of one rank was compared on one single core of the CPU and on one single GPU (with all its parallelizability). It is therefore not surprising that most of the time, the GPU solved the problem much faster than the CPU. Keeping this in mind, it is natural to associate a larger local lattice to the available GPUs in the system (maybe 4-8 times larger, depending on the available memory) and let them participate to the solution of the full problem just as they were additional ranks (compare proposal 13.1).

10 Deflated SAP preconditioned GCR algorithm

The low modes of the Dirac operator condensate TODO.

Small quark masses corresponding to real physics are believed to be the cause for the spontaneously breaking of chiral symmetry in lattice QCD [1]. Numerical lattice QCD has the problem that with large lattice volumes and small quark masses simulation techniques become inefficient in the *chiral regime* (where chiral symmetry is spontaneously broken), because the Dirac operators gets more and more ill-conditioned. Thus, the presence of low eigenvalues is a source of difficulty [10]. According to the Bank-Casher relation [1], this is because the number of eigenvalues of D below a fixed value grows with $O(V)$, where V is the total 4D lattice volume. On the other hand, the computational effort scales even worse with $O(V^2)$ [16]. This behaviour goes under the name of *V^2 -problem*.

A solving algorithm that has a flat scaling in with respect to the quark masses can therefore lead to large speedups specially in that regime. By deflating the Dirac operator, it is possible to separate eigenmodes with very small eigenvalues from the others. Thus the space needs to be split in low and high modes without actually calculating the modes, else the problem would be solved already.

10.1 Deflation

Theorem 10.1 (Deflation). *Let A be a linear, invertible operator acting on a vector space Λ , $\vec{b} \in \Lambda$ a arbitrary vector and P_L a projector³² acting on Λ . Also, define the linear operator P_R such that $P_L A = A P_R$ ³³. Consider*

$$\vec{x}^* := P_R \vec{x}_1^* + (1 - P_R) \vec{x}_2^*, \quad (10.1)$$

with \vec{x}_1^* and \vec{x}_2^* being solutions to the "smaller" (projected) systems

$$P_L A \vec{x}_1 = P_L \vec{b} \quad \text{and} \quad (1 - P_L) A \vec{x}_2 = (1 - P_L) \vec{b}$$

respectively. Then

- 1) P_R is a projector,
- 2) \vec{x}^* is the solution to $A \vec{x} = \vec{b}$.

Proof. Using that $P_L^2 = P_L$ is a projector and the defining relation $P_L A = A P_R$,

$$\begin{aligned} P_R^2 &= (A^{-1} P_L A)^2 \\ &= A^{-1} P_L^2 A \\ &= A^{-1} P_L A \\ &= P_R. \end{aligned}$$

By direct calculation,

$$\begin{aligned} A \vec{x}^* &= A P_R \vec{x}_1^* + A (1 - P_R) \vec{x}_2^* \\ &= P_L A \vec{x}_1^* + (1 - P_L) A \vec{x}_2^* \\ &= P_L \vec{b} + (1 - P_L) \vec{b} \\ &= \vec{b}. \end{aligned}$$

□

Remark. Therefore, if we find clever projectors P_L and P_R without involving A^{-1} , we can solve $A \vec{x} = \vec{b}$ by solving the 2 smaller systems of equations and then projecting the solutions using P_R .

³² P_L does not have to be orthogonal or Hermitian.

³³Such a linear operator P_R always exists - just set $P_R := A^{-1} P_L A$, since A is invertible.

Remark. Notice that $P_L A$ as well as $(1 - P_L)A$ are not invertible, therefore there are infinitely many solutions \vec{x}_1^* and \vec{x}_2^* ³⁴. Nonetheless the solution vector \vec{x}^* is still unique after the projection in equation (10.1), because P_R is a projector.

Remark. Comparing deflation to left preconditioning, the difference is that in deflation P_L is a projector and $P_L A$ has condition number infinite whereas in case of preconditioning P_L is invertible (a good approximation of A^{-1}) and the condition number of $P_L A$ is expected to be smaller than of A .

Corollary 10.2. Let A and \vec{b} be as in theorem 10.1. Furthermore let $\{\vec{\omega}_i\}_{i=1}^N$ be an orthonormal basis of a linear subspace $\Omega \subset \Lambda$, called the **deflation subspace** and let the restriction of A to Ω , $\tilde{A} := A|_{\Omega}$ called the **little operator**, be invertible. Define the action of P_L on an arbitrary vector $\vec{x} \in \Lambda$ as

$$P_L \vec{x} := \vec{x} - \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle$$

and let \vec{x}_1^* be one of the (infinite) solutions to the **deflated system** $\hat{A} \vec{x}_1 = P_L \vec{b}$, where $\hat{A} := P_L A$ is called the **deflated operator**. Consider

$$\vec{x}^* := P_R \vec{x}_1^* + \sum_{i,j=1}^N \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{b} \rangle, \quad (10.2)$$

with P_R satisfying $P_L A = A P_R$. Then \vec{x}^* is the unique solution to $A \vec{x} = \vec{b}$.

Proof. Lets first show that $P_L^2 = P_L$ is a projector,

$$\begin{aligned} P_L^2 \vec{x} &= P_L \left(\vec{x} - \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle \right) \\ &= \vec{x} - 2 \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle + \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \sum_{k,l=1}^N \langle \vec{\omega}_j, A \vec{\omega}_k \rangle (\tilde{A}^{-1})_{kl} \langle \vec{\omega}_l, \vec{x} \rangle \\ &= \vec{x} - 2 \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle + \sum_{i,j,l=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_l, \vec{x} \rangle \underbrace{\sum_{k=1}^N \langle \vec{\omega}_j, A \vec{\omega}_k \rangle (\tilde{A}^{-1})_{kl}}_{\substack{= \tilde{A}_{jk} \\ = \delta_{jl}}} \\ &= \vec{x} - 2 \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle + \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle \\ &= \vec{x} - \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle \\ &= P_L \vec{x}. \end{aligned}$$

Now lets show that the second term in equation (10.2) is equal to $(1 - P_R) \vec{x}_2^*$ where \vec{x}_2^* solves the projected system $(1 - P_L) A \vec{x}_2 = (1 - P_L) \vec{b}$.

$$\begin{aligned} (1 - P_R) \vec{x}_2^* &= A^{-1} (1 - P_L) A \vec{x}_2^* \\ &= A^{-1} (1 - P_L) \vec{b} \end{aligned}$$

³⁴Let P be a linear projector (not the identity-operator) and A an invertible linear operator. The system of interest is $PA\vec{x} = P\vec{b}$. There exists at least one solution to this, namely the unique solution to $A\vec{x} = \vec{b}$. Since PA is not invertible, the only two possibilities are zero or infinite solutions and it can't be zero solutions.

$$\begin{aligned}
&= A^{-1} \sum_{i,j=1}^N A\vec{\omega}_i(\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{b} \rangle \\
&= \sum_{i,j=1}^N \vec{\omega}_i(\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{b} \rangle
\end{aligned}$$

which corresponds to the second term of \vec{x}^* in equation (10.2). Therefore by application of theorem 10.1, \vec{x}^* is the unique solution to $A\vec{x} = \vec{b}$. \square

Remark. From P_L in corollary 10.2, the action of P_R on an arbitrary vector \vec{x} can be determined using the defining relation of P_R as

$$\begin{aligned}
P_R \vec{x} &= A^{-1} P_L A \vec{x} \\
&= \vec{x} - \sum_{i,j=1}^N \vec{\omega}_i(\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, A \vec{x} \rangle.
\end{aligned}$$

Remark. An application of P_L to an arbitrary vector \vec{x} involves solving the **little equation** $\tilde{A}\vec{\beta} = \vec{\alpha}$ on Ω for a given $\vec{\alpha} \in \Omega$. To see this, let's look at the k -th component of $P_L \vec{x}$

$$(P_L \vec{x})_k := x_k - \sum_{i,j=1}^N (A\vec{\omega}_i)_k (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle.$$

Define the vector

$$\vec{\alpha}_{\vec{x}} := \begin{pmatrix} \langle \vec{\omega}_1, \vec{x} \rangle \\ \langle \vec{\omega}_2, \vec{x} \rangle \\ \vdots \\ \langle \vec{\omega}_N, \vec{x} \rangle \end{pmatrix}.$$

Then

$$(P_L \vec{x})_k = x_k - \sum_{i=1}^N (A\vec{\omega}_i)_k (\tilde{A}^{-1} \vec{\alpha}_{\vec{x}})_i.$$

By similar analysis, an application of P_R has the same cost with one additional application of A . For an efficient implementation of P_L and P_R , the $2N$ vectors $\{A\vec{\omega}_i\}_{i=1}^N$ and $\{\vec{\omega}_i\}_{i=1}^N$ have to be kept in system memory.

Remark. Assuming that the condition number of A is high and the **spectrum** of A , $\sigma(A)$, is separable in a way such that

$$\sigma(A) = \sigma_l(A) \cup \sigma_h(A) \quad \text{with} \quad \max_{\lambda \in \sigma_l(A)} |\lambda| \ll \min_{\lambda \in \sigma_h(A)} |\lambda|. \quad (10.3)$$

The subscripts stand for "low" and "high", corresponding to the low and high modes of the operator A . So, the property in equation (10.3) states that the bulk of the low and high eigenvalues are somehow clustered in two regions. Consider the linear sub-spaces $\Omega_l, \Omega_h \subset \Lambda$ such that the low and high eigenvectors corresponding to the low and high eigenvalues of A are contained in Ω_l and Ω_h respectively. Then the condition number of A restricted to the low (high) modes is much smaller than the condition number of A . Therefore, if we are able to find an orthonormal basis $\{\vec{\omega}_i\}_{i=0}^N$ of the subspace Ω_l containing the bulk of the low eigenmodes of A , we can apply deflation from corollary 10.2 to solve the little equation that has a significantly smaller condition number than A . Then solve the deflated system and using this solution construct a solution of the full system.

Lemma 10.3. Let A , $\{\vec{\omega}_i\}_{i=1}^N$, Ω , P_L , P_R be as in corollary 10.2 and assume that the spectrum of A is separable (10.3). Define the deflation subspace to be the subspace corresponding to the low eigenmodes, $\Omega := \Omega_l$. Then $\kappa(\hat{A}) \ll \kappa(A)$

Proof. Lets define the orthogonal projector P^\perp to Ω^\perp , the orthogonal complement of the deflation subspace of Ω ,

$$P^\perp \vec{x} := \vec{x} - \sum_{i=1}^N \langle \vec{\omega}_i, \vec{x} \rangle \vec{\omega}_i.$$

The deflated operator $\hat{A} := P_L A$ acts on the orthogonal complement,

$$\begin{aligned} \hat{A} P^\perp \vec{x} &= P_L A P^\perp \vec{x} \\ &= P_L A \vec{x} - \sum_{k=1}^N P_L A \vec{\omega}_k \langle \vec{\omega}_k, \vec{x} \rangle \\ &= A \vec{x} - \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, A \vec{x} \rangle - \sum_{k=1}^N A \vec{\omega}_k \langle \vec{\omega}_k, \vec{x} \rangle + \sum_{i,k=1}^N A \vec{\omega}_i \underbrace{\sum_{j=1}^N (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, A \vec{\omega}_k \rangle}_{\substack{= \tilde{A}_{jk} \\ \delta_{ik}}} \langle \vec{\omega}_k, \vec{x} \rangle \\ &= A \vec{x} - \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, A \vec{x} \rangle - \sum_{k=1}^N A \vec{\omega}_k \langle \vec{\omega}_k, \vec{x} \rangle + \sum_{k=1}^N A \vec{\omega}_k \langle \vec{\omega}_k, \vec{x} \rangle \\ &= A \vec{x} - \sum_{i,j=1}^N A \vec{\omega}_i (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, A \vec{x} \rangle \\ &= P_L A \vec{x} \\ &= \hat{A} \vec{x}. \end{aligned}$$

Define the *minimal and maximal eigenvalue* of A ,

$$\lambda_{\min}(A) := \min_{\lambda \in \sigma(A)} |\lambda| \quad \text{and} \quad \lambda_{\max}(A) := \max_{\lambda \in \sigma(A)} |\lambda|.$$

The condition number of \hat{A} can now be upper bounded,

$$\kappa(\hat{A}) = \frac{|\lambda_{\max}(\hat{A})|}{|\lambda_{\min}(\hat{A})|} \ll \frac{|\lambda_{\max}(A)|}{|\lambda_{\min}(\hat{A})|} \leq \frac{|\lambda_{\max}(A)|}{|\lambda_{\min}(A)|} = \kappa(A),$$

where property (10.3) as used in the first inequality. □

Remark. Lemma 10.3 tells us that the deflated system is significantly better conditioned than the full system and is therefore solved in fewer iterations.

Lemma 10.4. P_L as defined in corollary 10.2 is a projection to the orthogonal complement of Ω , i.e. $\langle \vec{\omega}_k, P_L \vec{x} \rangle = 0$.

Proof. Let \vec{x} be an arbitrary vector, and $k \in \{1, \dots, N\}$, then

$$\langle \vec{\omega}_k, P_L \vec{x} \rangle = \langle \vec{\omega}_k, \vec{x} \rangle - \sum_{i,j=1}^N \langle \vec{\omega}_k, A \vec{\omega}_i \rangle (\tilde{A}^{-1})_{ij} \langle \vec{\omega}_j, \vec{x} \rangle$$

$$\begin{aligned}
&= \langle \vec{\omega}_k, \vec{x} \rangle - \sum_{j=1}^N \langle \vec{\omega}_j, \vec{x} \rangle \sum_{i=1}^N \tilde{A}_{ki} (\tilde{A}^{-1})_{ij} \\
&= \langle \vec{\omega}_k, \vec{x} \rangle - \sum_{j=1}^N \langle \vec{\omega}_j, \vec{x} \rangle \delta_{kj} \\
&= 0.
\end{aligned}$$

□

10.2 Choosing the deflation subspace

Definition 10.1 (Local coherence). *Let $\varepsilon > 0$, $M, N \in \mathbb{N}$ with $M \ll N$ and Λ be the full lattice. Let $B := \{\vec{\rho}_i\}_{i=1}^N$ and $C := \{\vec{\omega}_j\}_{j=1}^M$ with $\vec{\rho}_i, \vec{\omega}_j \in \Lambda$ be two sets of fields. Furthermore define the projector P_Γ to the subspace $\Gamma := \text{span}(C)$ via its action onto an arbitrary vector $\vec{x} \in \Lambda$ as*

$$P_\Gamma \vec{x} := \sum_{j=1}^M \langle \vec{\omega}_j, \vec{x} \rangle \vec{\omega}_j.$$

The fields B are **locally coherent** with the fields C up to ε , if $\forall i \in \{1, \dots, N\}$

$$\|(1 - P_\Gamma) \vec{\rho}_i\| \leq \varepsilon.$$

Remark. Definition 10.1 can be interpreted such that the fields in B can be well approximated by the (much fewer) fields in C , since the approximation error is smaller or equal to ε .

11 Multi-shift Conjugate Gradient algorithm

TODO

Proposal 11.1: MSCG in mixed precision

TODO: Multishift Conjugate Gradient in mixed precision. Currently only in binary64.

12 Dirac operator

Proposal 12.1: Representation of the Dirac-operator

For the implementation of the Dirac-operator on the GPU, the software library QUDA [6] is a good sample. To improve the performance of their Dirac-operator, the authors of QUDA used a representation of the SU(3)-fields with 8 real numbers, a gauge transformation to make almost all of the gauge fields in temporal direction to the identity-matrix and a change of basis in the γ -matrices, such that one of the four matrices has a very simple form. The most interesting one is probably the SU(3)-representation with only 8 real numbers. In openQxD ^a the struct `su3_double` representing a SU(3)-gauge-field consists of 18 double precision numbers. The C-macro for a matrix-matrix multiplication of 2 such `su3_double` ^b consists of 18·12 FLOPs, 2·18 loads and 18 stores. Using `binary64`, the arithmetic intensity is $I = 0.5$ FLOPs per byte, making the problem memory-bound.

Since the rows and columns of SU(3)-matrices form a orthonormal basis of \mathbb{C}^3 , one representation of such matrices can hold only the first two rows or columns and the third row or column is calculated as the vector product of the former two [9]. This is a representation with 12 real numbers. A matrix-matrix multiplication of 2 such matrices ends up in 270 FLOPs, 2·12 loads and 12 stores. This results in an arithmetic intensity of $I = 0.9375$ FLOPs per byte using `binary64` - still memory bound.

If the representation of a $SU(3)$ -gauge-field would be chosen such that the struct contains only 10 numbers [4], then a matrix $A \in SU(3)$ would be represented as $(a_{ij} \in \mathbb{C})$

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & Na_{31}^* & Na_{21} \\ 0 & -Na_{21}^* & Na_{31} \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ 0 & -Na_{13}^* & -Na_{12}^* \\ \frac{1}{N} & -Na_{11}^*a_{12} & -Na_{11}^*a_{13} \end{pmatrix}$$

$$= \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & -N^2(a_{13}a_{31}^* + a_{11}^*a_{12}a_{21}) & -N^2(a_{12}a_{31}^* + a_{11}^*a_{13}a_{21}) \\ a_{31} & -N^2(a_{13}a_{21}^* + a_{11}^*a_{12}a_{31}) & -N^2(a_{12}a_{21}^* + a_{11}^*a_{13}a_{31}) \end{pmatrix}$$

where $N := (1 - |a_{11}|^2)^{-\frac{1}{2}}$. Notice that this representation has a singularity at $|a_{11}| = 1$. The 5 complex numbers $a_{11}, a_{12}, a_{13}, a_{21}, a_{31}$ are subject of the orthonormality constraints

$$|a_{11}|^2 + |a_{12}|^2 + |a_{13}|^2 = |a_{11}|^2 + |a_{21}|^2 + |a_{31}|^2 = 1, \quad (12.1)$$

leading to the observation that all 12 real numbers are in the set $[-1, 1]$.

Finally, in a minimal representation of 8 real numbers [4], can be obtained using the constraint above (12.1). If we write $a_{ij} = x_{ij} + iy_{ij}$ with $x_{ij}, y_{ij} \in \mathbb{R}$ we can eliminate 2 further numbers. One (of many) choices could be

$$y_{31} = \sqrt{1 - |a_{11}|^2 - |a_{21}|^2 - |x_{31}|^2} \quad (12.2)$$

$$y_{13} = \sqrt{1 - |a_{11}|^2 - |a_{13}|^2 - |x_{13}|^2}. \quad (12.3)$$

Using this, only a_{11}, a_{12}, a_{21} and the real parts of a_{13} and a_{31} need to be stored in memory. For a breakdown of the different arithmetic intensities of the representations, see table 5.

Arithmetic intensities in FLOPs per byte			
Reals	I(binary64)	I(binary32)	I(binary16)
18	0.5	1	2
12	0.9375	1.875	3.75
10	1.3667	2.7333	5.4667
8	1.9115	3.8229	7.6458

Table 5: Arithmetic intensities of $SU(3)$ representations with different requirement for real numbers for a matrix-matrix multiplication. In the calculation of the intensities a FLOP-count of 6 was used for the square root of a floating point number and previous results were reused instead of recalculating.

In table 6 one invocation of `Dw_dble()` was called and the macros from `su3.h` where counted.

One invocation of <code>Dw_dble()</code>					
# calls	macro name	I(18)	I(12)	I(10)	I(8)
61440	<code>_vector_add_assign()</code>	0.0416	0.0416	0.0416	0.0416
24576	<code>_su3_multiply()</code>	0.3	0.625	0.9323	1.0989
24576	<code>_su3_inverse_multiply()</code>	0.3	0.625	0.9323	1.0989
12288	<code>_vector_add()</code>	0.0416	0.0416	0.0416	0.0416
12288	<code>_vector_i_add()</code>	0.0416	0.0416	0.0416	0.0416
12288	<code>_vector_i_add_assign()</code>	0.0416	0.0416	0.0416	0.0416
12288	<code>_vector_sub()</code>	0.0416	0.0416	0.0416	0.0416
12288	<code>_vector_i_sub()</code>	0.0416	0.0416	0.0416	0.0416
12288	<code>_vector_sub_assign()</code>	0.0416	0.0416	0.0416	0.0416
12288	<code>_vector_i_sub_assign()</code>	0.0416	0.0416	0.0416	0.0416

Table 6: Number of C macro calls for one call of `Dw_dble()` on a 8^4 local lattice with 4 ranks.

^aSee line 43 in `include/su3.h` in [5].

^bSee line 490ff in `include/su3.h` in [5].

Definition 12.1 (Hadamard product). *The **Hadamard product** of two vectors \vec{x} and \vec{y} is defined as*

$$\begin{aligned}\odot: \mathbb{R}^n \times \mathbb{R}^n &\rightarrow \mathbb{R}^n \\ (\vec{x}, \vec{y}) &\mapsto \vec{x} \odot \vec{y},\end{aligned}$$

with

$$(\vec{x} \odot \vec{y})_i := (\vec{x})_i (\vec{y})_i,$$

where $(\vec{v})_i$ denotes the i -th component of the vector \vec{v} .

13 GPU Implementation

There are multiple possibilities how to implement GPU-utilisation into the current state of the code. The proposals be divided roughly into 2 categories.

- Use the resources of the GPU to process a larger lattice in the same time as before.
- Use the resources of the GPU to process the same lattice, but faster than before.

Proposal 13.1: GPU Implementation Variant 1

Keeping the results in mind (figures 18 - 21) that the pure GPU implementation of the `SAP_GCR` was by far the fastest, it makes sense to associate a local lattice to each GPU - so to treat a GPU as an additional rank with its own local lattice. Each GPU would then act as another rank and the full lattice can be extended by as many local lattices as there are available GPUs. Some nodes might have GPUs, some not. Each GPU has a **bystander process** running on the CPU. Process contexts associated to a GPU (CPU) are from now on called **GPU ranks (CPU ranks)**. The bystander process will not need a lot of CPU load, since it only hosts control and informational variables and is not involved in any calculations. The communication from a GPU rank to another rank (CPU or GPU) in the system can be done either via the bystander process or via direct injection into the network bypassing the CPU and the bystander process. The former case will not need a change of the MPI communication code among ranks. Since the communication will only take negligible computational effort on the bystander process, it will not *steal* much from the CPU ranks running on the same node. Let me give an example: Let's assume the application runs on one single node with 8 cores and 4 GPUs attached to it. This would involve 8 CPU ranks (separate processes) and 4 GPU ranks with in total 4 (separate) bystander processes. Technically the machine is over-subscribed and the operating system needs to schedule the processes. But since the bystander processes will be in sleeping state most of the time, this will not (or only negligibly) degrade the performance of the 8 CPU ranks.

- Advantage: The application can run on hybrid machines in the sense that some nodes can have GPUs attached to them and others don't.
- Disadvantage: The full system consists of two types of ranks - GPU ranks and CPU ranks. Inevitably, either the GPU or the CPU ranks will be faster on the same local lattice size. This means that either the GPUs or the CPUs will have to wait for the others to finish. Utilisation of resources is not perfect.

Proposal 13.2: GPU Implementation Variant 2

Since the GPUs are fast on the solver algorithms, a pure GPU implementation of all currently implemented solvers ^a can lead to a significant speed up of the program.

- Advantage: Only the small subset of solver algorithm code has to be changed.
- Disadvantage: The Dirac-operator has to be held in the main memory as well as in the GPU memory. Both need to be in sync. This will lead to a lot of intra-node traffic, which itself should not be a problem. The Dirac-operator is stored in terms of its gauge-fields. They will be redundant and thus the full amount of GPU and CPU memory is not optimally utilised.
- Disadvantage: Since the solvers run on the GPUs, the CPUs will be stale during that time (although the waiting time will be smaller than the time they would need to run the solvers themselves). In the HMC-calculations done on the CPUs, the GPUs will be stale. Again the full potential of performance is not utilised.

^aCGNE, MSCG, SAP_GCR and DFL_SAP_GCR

Proposal 13.3: GPU Implementation Variant 3

If the problem is split among CPU ranks and GPU ranks, one will always have the waiting problem in the sense that either the CPU or the GPU ranks will have to wait for the others to complete (see proposal 13.1). In order for the GPU to speed up the process, every rank should receive the same amount of help from the GPUs. Then every rank solves the same problem faster. This is possible if all nodes participating in the calculation share the same specifications, the same number GPUs (not zero) and the same amount of memory. Let's assume this is given. Starting with a hybrid implementation as in proposal 9.2, only the blocked problems are solved on the GPU and the internal color boundary operator as well as the full Dirac-operator are performed on the CPU as usual (see figures 18 - 21, diamonds; \diamond , \diamond , \diamond). This means that *all* blocked problems of all ranks on a single node are solved (or rather *nmr* MR-steps are performed) on the GPUs of that node. During that time the CPUs are stale. To solve this problem, we can use the fact that all blocks of the same color are independent of each other. So, not all blocks need to be transferred to the GPU - some of them can still reside and be processed in the current rank on the CPU. The work here should be divided such that both - CPU and GPU - need approximately the same amount of time for the processing of their blocks. The question on *how* to divide the blocks still remains. It is evident that the GPU might be able to process more blocks than the CPU rank in the same time, but this highly depends on its occupation. A robust solution might implement the division of blocks in an adaptive manner.

- Advantage: This proposal can be a good starting point from where to go further.
- Disadvantage: The GPU is only utilised in one part of one solver algorithm, else the GPU is stale.

14 Algorithm-independent considerations

Proposal 14.1: Choice of starting vectors

During a simulation, a lot of linear systems of equation of the type $A\vec{x} = \vec{b}$ have to be solved. Most solving algorithms have the possibility to choose a initial starting vector \vec{x}_0 from where to start the iterative process (see definition 8.15 for example). Consider the sequence of matrices and source vectors for which the linear system of equations should be solved $\{A_1, \vec{b}_1\}, \{A_2, \vec{b}_2\}, \dots, \{A_n, \vec{b}_n\}$ and the matrix A not changing very much among steps

$(A_i \approx A_j)$. If the difference $\|\vec{b}_i - \vec{b}_j\|$ is small, then the difference of the solution vectors $\|\vec{x}_i - \vec{x}_j\|$ can be expected to be small too. Assuming that the system $A_i \vec{x}_i = \vec{b}_i$ is already solved, the iterative solver for $A_j \vec{x}_j = \vec{b}_j$ can use \vec{x}_i as its initial starting vector and reduce the amount of steps needed to solve the latter system.

15 Summary

TODO

16 Future

17 References

- [1] T. Banks and A. Casher. Chiral symmetry breaking in confining theories. *Nuclear Physics B*, 169(1-2):103–125, 1980.
- [2] N. Bell and M. Garland. Efficient sparse matrix-vector multiplication on cuda. Technical report, Citeseer, 2008.
- [3] I. Buck. Taking the plunge into gpu computing. *GPU Gems*, 2:509–519, 2005.
- [4] B. Bunk and R. Sommer. An 8 parameter representation of $su(3)$ matrices and its application for simulating lattice qcd. *Computer physics communications*, 40(2-3):229–232, 1986.
- [5] I. Campos, P. Fritzsche, M. Hansen, M. Krstić Marinković, A. Patella, A. Ramos, and N. Tantalo. openq*d. <https://gitlab.com/rcstar/openQxD>, 2018. Accessed: 2021-01-06.
- [6] M. A. Clark, R. Babich, K. Barros, R. C. Brower, and C. Rebbi. Solving lattice qcd systems of equations using mixed precision solvers on gpus. *Computer Physics Communications*, 181(9):1517–1528, 2010.
- [7] W. Cody. Towards sensible floating-point arithmetic. Technical report, Argonne National Lab., IL (USA), 1980.
- [8] T. A. Davis and Y. Hu. The university of florida sparse matrix collection. *ACM Transactions on Mathematical Software (TOMS)*, 38(1):1–25, 2011.
- [9] P. De Forcrand, D. Lellouch, and C. Roiesnel. Optimizing a lattice qcd simulation program. *Journal of Computational Physics*, 59(2):324–330, 1985.
- [10] L. Giusti, C. Hoelbling, M. Lüscher, and H. Wittig. Numerical techniques for lattice qcd in the ϵ -regime. *Computer physics communications*, 153(1):31–51, 2003.
- [11] D. Göddeke, R. Strzodka, and S. Turek. *Accelerating double precision FEM simulations with GPUs*. Univ., 2005.
- [12] P. W. Group et al. Posit standard documentation - release 3.2-draft. *Posit Standard Documentation*, 2018.
- [13] J. L. Gustafson and I. T. Yonemoto. Beating floating point at its own game: Posit arithmetic. *Supercomputing Frontiers and Innovations*, 4(2):71–86, 2017.
- [14] R. Krashinsky, O. Giroux, S. Jones, N. Stam, and S. Ramaswamy. Nvidia ampere architecture in-depth. *NVIDIA blog*: <https://devblogs.nvidia.com/nvidia-ampere-architecture-in-depth>, 2020.
- [15] M. Lüscher. Solution of the dirac equation in lattice qcd using a domain decomposition method. *Computer physics communications*, 156(3):209–220, 2004.
- [16] M. Lüscher. Local coherence and deflation of the low quark modes in lattice qcd. *Journal of High Energy Physics*, 2007(07):081, 2007.
- [17] I. of Electrical, E. E. C. S. S. Committee, and D. Stevenson. *IEEE standard for binary floating-point arithmetic*. IEEE, 1985.
- [18] I. of Electrical, E. E. C. S. S. Committee, and D. Stevenson. *IEEE standard for binary floating-point arithmetic*. IEEE, 2008.
- [19] R. G. T. REMOVE. Todo remove github repository: Source code of the implementation. <http://github.com/chaoos/TODO>, 2021. Accessed: 2021-01-01.
- [20] J. R. Shewchuk et al. An introduction to the conjugate gradient method without the agonizing pain, 1994.
- [21] P. Van Nieuwenhuizen and A. Waldron. A continuous wick rotation for spinor fields and supersymmetry in euclidean space. *arXiv preprint hep-th/9611043*, 1996.

- [22] C. Vuik. New insights in gmres-like methods with variable preconditioners. *Journal of computational and applied mathematics*, 61(2):189–204, 1995.
- [23] S. Wang and P. Kanwar. Bfloat16: the secret to high performance on cloud tpus. *Google Cloud Blog*, 2019.
- [24] C.-N. Yang and R. L. Mills. Conservation of isotopic spin and isotopic gauge invariance. *Physical review*, 96(1):191, 1954.

Appendices

A Proofs

Theorem .1 (Plaquette). *Let a function $F(x)$ be defined as*

$$F_x(\epsilon) := e^{-\epsilon f(x+\epsilon a_1)} e^{-\epsilon g(x+\epsilon a_2)} e^{\epsilon f(x+\epsilon a_3)} e^{\epsilon g(x+\epsilon a_4)},$$

where x, a_1, a_2, a_3, a_4 are 4-vectors and $f(x), g(x)$ are non commuting functions of x . Then the derivatives of $F_x(\epsilon)$ are

$$\left. \frac{\partial F_x}{\partial \epsilon} \right|_{\epsilon=0} = 0 \quad (.1)$$

$$\left. \frac{\partial^2 F_x}{\partial \epsilon^2} \right|_{\epsilon=0} = 2[f(x), g(x)] + 2[(a_3 - a_1)^\mu f'_\mu(x) + (a_4 - a_2)^\mu g'_\mu(x)], \quad (.2)$$

where a short-hand notation for the derivatives was used; $f'_\mu(x) := \left. \frac{\partial f}{\partial x^\mu} \right|_x$.

Proof. Let's introduce some shorthand notation:

$$F_x(\epsilon) := \underbrace{e^{-\epsilon f_1(x+\epsilon a_1)}}_{=:F_1(\epsilon)} \underbrace{e^{-\epsilon f_2(x+\epsilon a_2)}}_{=:F_2(\epsilon)} \underbrace{e^{\epsilon f_3(x+\epsilon a_3)}}_{=:F_3(\epsilon)} \underbrace{e^{\epsilon f_4(x+\epsilon a_4)}}_{=:F_4(\epsilon)}, \quad (.3)$$

such that for $j \in \{1, 2, 3, 4\}$ we have $F_j(\epsilon) = \exp((-1)^{\lceil j/2 \rceil} \epsilon f_j(x + \epsilon a_j))$ with $f(x) = f_1(x) = f_3(x)$ and $g(x) = f_2(x) = f_4(x)$.

Starting with equation (.1),

$$\begin{aligned} \left. \frac{\partial F_x}{\partial \epsilon} \right|_{\epsilon=0} &= \left[F_1(\epsilon)(-f(x + \epsilon a_1) - \epsilon f'_\mu(x + \epsilon a_1))F_2(\epsilon)F_3(\epsilon)F_4(\epsilon) \right. \\ &\quad + F_1(\epsilon)F_2(\epsilon)(-g(x + \epsilon a_2) - \epsilon g'_\mu(x + \epsilon a_2))F_3(\epsilon)F_4(\epsilon) \\ &\quad + F_1(\epsilon)F_2(\epsilon)F_3(\epsilon)(f(x + \epsilon a_3) + \epsilon f'_\mu(x + \epsilon a_3))F_4(\epsilon) \\ &\quad \left. + F_1(\epsilon)F_2(\epsilon)F_3(\epsilon)F_4(\epsilon)(g(x + \epsilon a_4) + \epsilon g'_\mu(x + \epsilon a_4)) \right] \Big|_{\epsilon=0} \\ &= -f(x) - g(x) + f(x) + g(x) \\ &= 0. \end{aligned}$$

TODO: second derivative.

□

B Code

All code used in this report is open source and can be found in the GitHub repository [19]

C List of Proposals

1.1 Example proposal	2
8.1 Mixed Precision	35
8.2 Approximating the amounts α_i	36
9.1 MR in reduced precision	43
9.2 Performing the MR steps on the GPU	43
9.3 GCR in mixed precision	47

11.1 MSCG in mixed precision	56
12.1 Representation of the Dirac-operator	56
13.1 GPU Implementation Variant 1	58
13.2 GPU Implementation Variant 2	58
13.3 GPU Implementation Variant 3	59
14.1 Choice of starting vectors	59

Acronyms

BLAS	Basic Linear Algebra Subprograms. 26
CG	Conjugate Gradient. 10
CGNE	Conjugate Gradient on the normal equations. 26, 58
CSR	Compressed Sparse Row. 42
CUDA	Compute Unified Device Architecture. 42
DFL_SAP_GCR	Deflated Generalized Conjugate Residual with Schwarz Alternating Preconditioning. 58
FLOPS	Floating Point Operations Per Second. 33
GCR	Generalized Conjugate Residual. 39, 41, 43, 44, 46, 47
MR	Minimal Residual. 41, 42, 46
MSCG	Multishift Conjugate Gradient. 55, 58
NaN	not a number. 12–15, 30, 33
POPS	Posit Operations Per Second. 33
QCD	Quantum chromodynamics. 1
SAP	Schwarz Alternating Procedure. 39–41, 44
SAP_GCR	Generalized Conjugate Residual with Schwarz Alternating Preconditioning. 39, 46, 49, 50, 58
SF	Schrödinger functional. 28, 40

Glossary

bfloat16	Googles Brain float [23] floating point number representation with encoding in length of 16 bits. 12, 14, 17, 28, 30, 32, 33, 35, 36, 46
binary16	IEEE754 2008 [18] conformant floating point number representation with encoding in length of 16 bits. 12, 14, 17, 18, 28–33, 35–37, 46
binary32	IEEE754 2008 [18] conformant floating point number representation with encoding in length of 32 bits. 10, 12, 14, 17, 18, 26–28, 30–33, 35, 36, 46, 49
binary64	IEEE754 2008 [18] conformant floating point number representation with encoding in length of 64 bits. 10, 12, 14, 26–36, 46, 49, 55

fused multiply-add A multiply-add operation $a + bc$ in one shot, where the rounding is deferred. 16

posit16 Posit Standard [12] conformant storage format for real number representation with encoding in length of 16 bits and an exponent size of **es=1**. 16, 17, 28–33, 35

posit32 Posit Standard [12] conformant storage format for real number representation with encoding in length of 32 bits and an exponent size of **es=2**. 16, 17, 28, 30–33

posit64 Posit Standard [12] conformant storage format for real number representation with encoding in length of 64 bits and an exponent size of **es=3**. 16

posit8 Posit Standard [12] conformant storage format for real number representation with encoding in length of 8 bits and an exponent size of **es=0**. 16, 28, 30–32

quire Posit Standard [12] conformant special fixed-size data type that can be thought of as a dedicated register that permits dot products, sums, and other operations to be performed with rounding error deferred to the very end of the calculation [?]. 16, 29, 31–33

rank In **MPI** a process is identified by its rank, which is an integer between $[0, N - 1]$, where N is the size of the MPI process group. 19

sparse matrix A matrix, where most of the entries are 0. 19

tensorfloat32 Nvidias TensorFloat-32 [14] floating point number representation with encoding in length of 32 bits, but only 19 bits are used. 12, 14, 17, 28, 30–33, 35, 37, 46