

# On the suitability of the Brillouin action as a kernel to the overlap procedure

Stephan Dürr<sup>a,b</sup> and Giannis Koutsou<sup>c</sup>

<sup>a</sup>University of Wuppertal, Gaußstraße 20, 42119 Wuppertal, Germany

<sup>b</sup>IAS/JSC, Forschungszentrum Jülich GmbH, 52425 Jülich, Germany

<sup>c</sup>Cyprus Institute, CaSToRC, 20 Kavafi Street, Nicosia 2121, Cyprus

## Abstract

We investigate the Brillouin action in terms of its suitability as a kernel to the overlap procedure, with a view on both heavy and light quark physics. We use the diagonal elements of the Kenney-Laub family of iterations for the sparse matrix sign function, since they grow monotonically and facilitate cascaded preconditioning strategies with different rational approximations to the sign function. We find that the overlap action with the Brillouin kernel is significantly better localized than the version with the Wilson kernel.

## 1 Introduction

One of the key issues in a numerical study of lattice QCD is a suitable choice of the lattice Dirac operator, as this choice has a major impact on the overall cost, in terms of CPU time, of the computation. Whenever processes with non-zero momentum transfer are considered (e.g. in meson and baryon form factors which are relevant for semileptonic decays) the lattice dispersion relation is of interest, i.e. how much the continuum relation  $(aE)^2 - (a\mathbf{p})^2 = (am)^2$  is violated, where we use the lattice spacing  $a$  to build dimensionless quantities.

The Wilson Dirac operator [1, 2] and the Brillouin Dirac operator [3]

$$D_{\text{wil}}(x, y) = \sum_{\mu} \gamma_{\mu} \nabla_{\mu}^{\text{std}}(x, y) - \frac{a}{2} \Delta^{\text{std}}(x, y) + m_0 \delta_{x, y} - \frac{c_{\text{SW}}}{2} \sum_{\mu < \nu} \sigma_{\mu\nu} F_{\mu\nu} \delta_{x, y} \quad (1)$$

$$D_{\text{bri}}(x, y) = \sum_{\mu} \gamma_{\mu} \nabla_{\mu}^{\text{iso}}(x, y) - \frac{a}{2} \Delta^{\text{bri}}(x, y) + m_0 \delta_{x, y} - \frac{c_{\text{SW}}}{2} \sum_{\mu < \nu} \sigma_{\mu\nu} F_{\mu\nu} \delta_{x, y} \quad (2)$$

both show cut-off effects  $\propto a$  which can be reduced to  $\propto a^2$  by proper tuning of the coefficient  $c_{\text{SW}}$  [4–7]. The only difference is the discretization used for the covariant derivative  $\nabla_{\mu}$  and the gauged Laplacian  $\Delta$ ; the former operator uses a 9-point stencil for the Laplacian, while the latter operator uses a 81-point stencil (the Nabla operator always uses a subset of that stencil). The larger stencil allows for an improved dispersion relation (see Refs. [3, 8, 9] and below), but obviously the numerical cost is increased.

Regardless whether  $c_{\text{SW}}$  is zero or tuned to remove the  $O(a)$  on-shell cut-off effects, the operators (1, 2) are subject to limitations concerning the (renormalized) quark mass (which derives from the bare quark mass  $m_0$ ) that can be used in a simulation. For light quarks there

is an algorithmic bound for such non-chiral actions [10,11], and for heavy quark masses cut-off effects tend to proliferate (unless special measures are taken, see e.g. Refs. [9,12,13]).

The algorithmic limitation how light a quark mass may be taken at a given value of the gauge coupling  $\beta = 6/g_0^2$  is absent for chiral actions, i.e. for actions which satisfy the Ginsparg-Wilson relation [14–17]. The overlap construction (here and below this term is meant to include both the “domain-wall” [18–20] and the “overlap” [21,22] emanation of this idea) manages to upgrade a non-chiral into a chiral action. This is a highly practical procedure, though it is somewhat expensive in terms of CPU time (see below).

As a side effect, the overlap construction leads to automatic  $O(a)$  improvement. In other words no tuning of a coefficient like  $c_{\text{SW}}$  in (1, 2) is needed; the requirement of chiral symmetry automatically kills odd powers of  $a$  in on-shell quantities [23]. This is the reason why the overlap action with the Wilson kernel has proven very useful in heavy quark physics, see for instance the charm physics programs by the Kentucky group, JLQCD, and RBC/UKQCD [24–26].

In this paper we wish to explore whether there is any relevant improvement if one replaces, in the overlap construction, the Wilson kernel by the Brillouin kernel. Ideally, such an action would enable one to use a uniform relativistic formulation to simulate all hadronizing quarks ( $d, u, s, c, b$ ) at their physical mass values, on accessible lattices. The first technical question is whether the improved dispersion relation of the Brillouin operator for light quark masses (both at the quark-level [3,9] and for hadronic quantities [8]) would persist after the overlap procedure has been applied. The second question is whether the CPU requirements of the Brillouin-overlap action are roughly comparable to those of the Wilson-overlap formulation or whether they proliferate. The third question is whether there is any notable technical difference between the two overlap formulations, e.g. in terms of operator locality.

The remainder of this paper addresses these questions in due turn, intertwined with a few reminders on the overlap formulation and its technical implementation in a sparse matrix setup to make it self-contained. Sec. 2 presents an investigation of the free-field dispersion relations of both the Wilson and Brillouin kernel, along with their overlap descendents. Sec. 3 summarizes some knowledge about the Kenney-Laub family of iterations for the matrix sign function, since the diagonal members of this family show properties which we consider particularly convenient for the implementation of an overlap-times-vector application. Sec. 4 gives a quick review of the overlap construction and discusses a way of introducing the mass in the overlap operator which avoids any “extra prescription” if the Green’s function is used in the computation of a decay constant or matrix element. Sec. 5 illustrates the eigenvalues of some low-order Kenney-Laub iterates of the Wilson and Brillouin kernels on small lattices where all eigenvalues can be calculated. Sec. 6 presents the spectral flow, i.e. eigenvalues of the shifted hermitean kernels (for both formulations) on some selected gauge backgrounds. Sec. 7 addresses the aforementioned technical issues, such as the operator locality, and reports on a pilot spectroscopy calculation on  $40^3 \times 64$  lattices generated by QCDSF. Sec. 8 is a reminder that the Kenney-Laub family of matrix iterations offers many possibilities for cascaded preconditioning strategies where very-low-order polynomial approximations to the sign function are used to speed-up computations with not-so-low-order approximations. Sec. 9 gives reasons why we feel optimistic about the use of the framework portrayed in this article in future studies of full QCD with exact (i.e. arbitrarily good) chiral symmetry. Sec. 10 contains our summary, and some technical material is arranged in three appendices. A preliminary account of this work appeared in Ref. [27].

## 2 Quark-level dispersion relations

In this section we discuss the free-field dispersion relations of the Wilson and Brillouin operators, as well as those of their overlap descendents.

For the Wilson operator the dispersion relation reads (see App. A for details)

$$2 \cosh(aE) \left[ 4 + am - \sum_i \cos(ap_i) \right] = 1 + \sum_i \sin^2(ap_i) + \left[ 4 + am - \sum_i \cos(ap_i) \right]^2 \quad (3)$$

and an expansion in powers of  $a$  yields [9]

$$\begin{aligned} (aE)^2 - (a\mathbf{p})^2 &= \left[ (am)^2 - (am)^3 + \frac{11}{12}(am)^4 - \frac{5}{6}(am)^5 \right] \\ &+ \left[ -\frac{2}{3}(am)^2 + \frac{7}{6}(am)^3 \right] (a\mathbf{p})^2 \\ &+ \left[ -\frac{2}{3} + \frac{am}{2} \right] \left( \sum_{i < j} a^4 p_i^2 p_j^2 + \sum_i (ap_i)^4 \right) + O(a^6) . \end{aligned} \quad (4)$$

For the Brillouin operator the dispersion relation reads (see App. A for details)

$$\begin{aligned} &\frac{1}{729} \sum_i s_i^2 \prod_{j \neq i} \{c_j + 2\}^2 \{\cosh^2 + 4 \cosh + 4\} + \frac{1}{729} \{1 - \cosh^2\} \prod_i \{c_i + 2\}^2 \\ &+ \frac{1}{64} \prod_i \{c_i + 1\}^2 \{\cosh^2 + 2 \cosh + 1\} - \frac{1}{4} \prod_i \{c_i + 1\} \{\cosh + 1\} [2 + am] + [2 + am]^2 = 0 \end{aligned} \quad (5)$$

with  $s_i = \sin(ap_i)$ ,  $c_i = \cos(ap_i)$ , and an expansion of the physical solution yields [9]

$$\begin{aligned} (aE)^2 - (a\mathbf{p})^2 &= \left[ (am)^2 - (am)^3 + \frac{11}{12}(am)^4 - \frac{5}{6}(am)^5 \right] \\ &+ \left[ 0 + \frac{1}{12}(am)^3 \right] (a\mathbf{p})^2 \\ &+ \left[ 0 + \frac{am}{12} \right] \left( \sum_{i < j} a^4 p_i^2 p_j^2 + \sum_i (ap_i)^4 \right) + O(a^6) . \end{aligned} \quad (6)$$

As was already noted in Ref. [9], a comparison of (4) and (6) shows that the Brillouin construction manages to reduce the amount of isotropy breaking (the term  $\propto a^4$  in the last line vanishes, and the term  $\propto a^5$  receives a factor  $1/6$ ). However, the momentum independent part in the first line, which is an expansion of  $\log^2(1 + am)$ , is unchanged from the Wilson case [9]. This suggests that the Brillouin construction brings an advantage for heavy quark spectroscopy only in case non-zero spatial momenta are involved.

This conclusion is supported by the plots shown in Fig. 1. The Wilson operator at  $am = 0$  shows significant deviations from the continuum dispersion relation and strong isotropy violations (differences between the momentum directions). Furthermore, at the heavy quark mass  $am = 0.75$  strong cut-off effects even at  $a\mathbf{p} = \mathbf{0}$  become visible. The Brillouin operator at  $am = 0$  features a significantly improved dispersion relation with much smaller isotropy violations, but at  $am = 0.75$  the cut-off effects are equally large as in the Wilson case.

For the overlap operator with the Wilson kernel the dispersion relation follows from searching for zeros of  $c^2 + 2cd(\frac{a}{2}\hat{p}^2 - \frac{e}{a})[\bar{p}^2 + (\frac{a}{2}\hat{p}^2 - \frac{e}{a})^2]^{-1/2} + d^2 = 0$  with  $c, d, \hat{p}^2, \bar{p}^2$  given in App. A, and

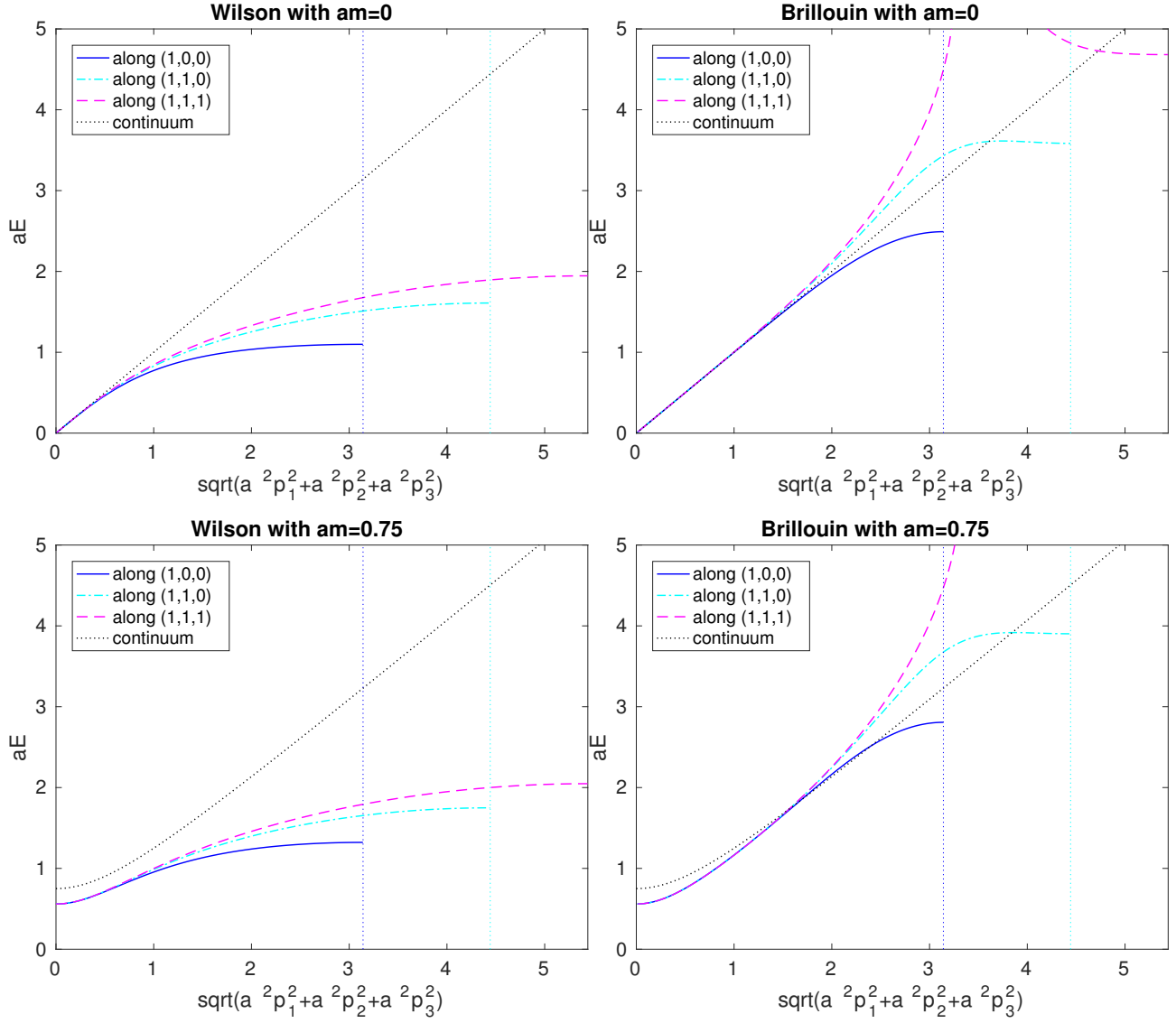


Figure 1: Free field dispersion relations of  $D_W$  (left) and  $D_B$  (right) for the bare quark masses  $am = 0$  (top) and  $am = 0.75$  (bottom). We plot the spatial directions  $(1, 0, 0)$ ,  $(1, 1, 0)$ , and  $(1, 1, 1)$ , where the Brillouin zone ends at  $\pi/a$ ,  $\sqrt{2}\pi/a$ , and  $\sqrt{3}\pi/a$ , respectively.

an expansion in powers of  $a$  yields

$$\begin{aligned}
 (aE)^2 - (a\mathbf{p})^2 &= \left[ (am)^2 - \frac{2\rho^2 - 6\rho + 3}{6\rho^2} (am)^4 \right] \\
 &+ \left[ -\frac{2}{3} (am)^2 + 0 \right] (a\mathbf{p})^2 \\
 &+ \left[ -\frac{2}{3} + 0 \right] \left( \sum_{i < j} a^4 p_i^2 p_j^2 + \sum_i (ap_i)^4 \right) + O(a^6). \quad (7)
 \end{aligned}$$

For the overlap operator with the Brillouin kernel the dispersion relation follows from searching for zeros of  $c^2 + 2cd(\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2} + d^2 = 0$  with  $c, d, \tilde{p}^2, \tilde{p}^2$  given in App. A, and

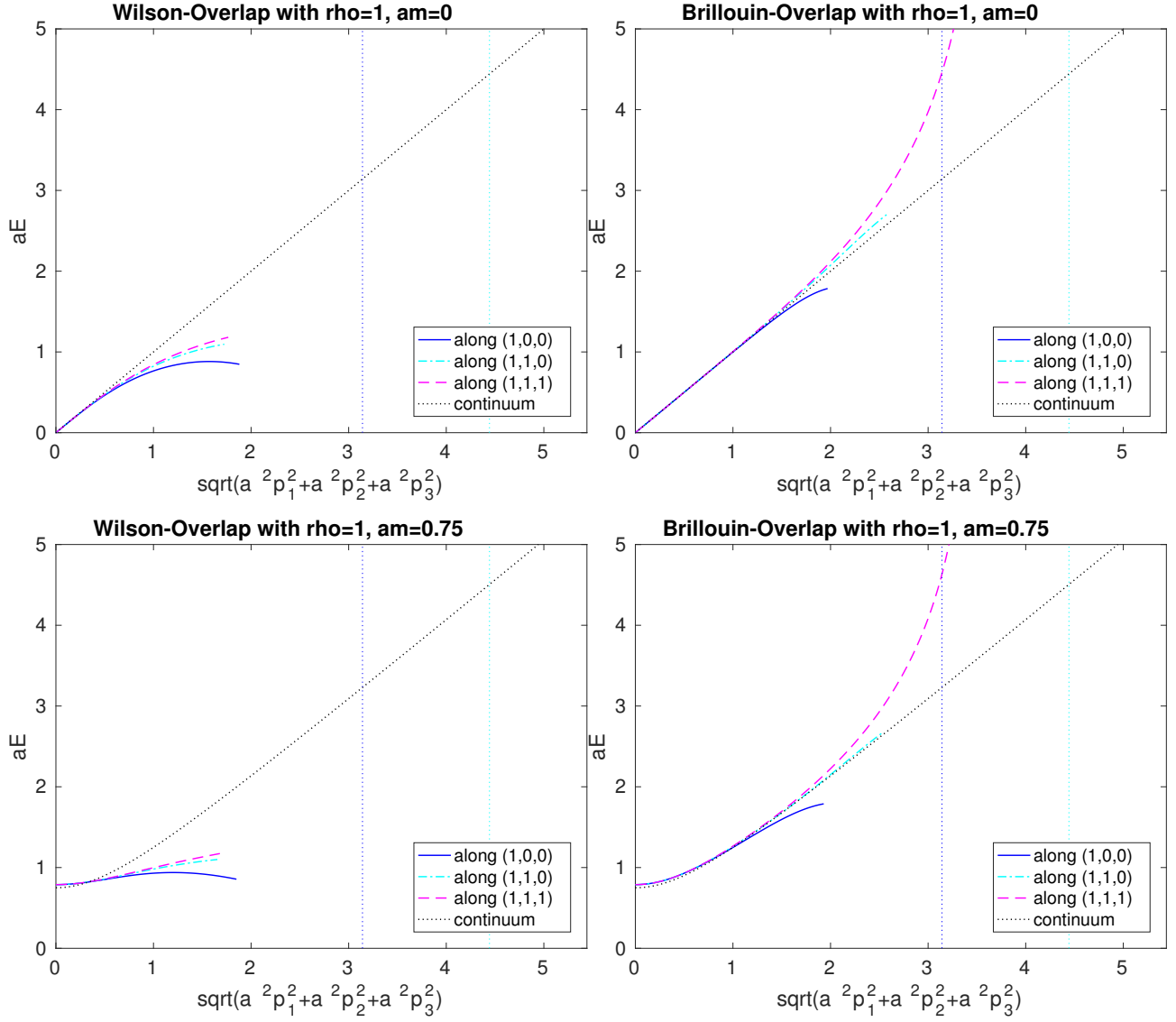


Figure 2: Same as Fig. 1 but for the overlap actions based on the Wilson (left) and Brillouin (right) kernels at  $\rho = 1$ . The overlap mass is again  $am = 0$  (top) and  $am = 0.75$  (bottom).

an expansion in powers of  $a$  yields

$$\begin{aligned}
(aE)^2 - (a\mathbf{p})^2 &= \left[ (am)^2 - \frac{2\rho^2 - 6\rho + 3}{6\rho^2} (am)^4 \right] \\
&+ \left[ 0 + 0 \right] (a\mathbf{p})^2 \\
&+ \left[ 0 + 0 \right] \left( \sum_{i < j} a^4 p_i^2 p_j^2 + \sum_i (ap_i)^4 \right) + O(a^6) .
\end{aligned} \tag{8}$$

Comparing (7) and (8) suggests that the Brillouin overlap inherits the reduced isotropy breaking from its kernel action. The cut-off effects at  $a\mathbf{p} = \mathbf{0}$  are still identical for the two overlap actions, and the removal of odd powers of  $am$  clearly reduces the momentum-dependent cut-off effects in comparison to the non-chiral predecessors. Note that in (7, 8) the coefficient

of  $(am)^4$  can be made zero by choosing  $\rho = (3 - \sqrt{3})/2 \simeq 0.634$ ; in this case the free-field Brillouin overlap dispersion relation is free of cut-off effects through  $O(a^5)$ .

This conclusion is supported by the plots shown in Fig. 2. The  $\rho = 1$  Wilson overlap operator at  $am = 0$  shows similar (or even worse) deviations from the continuum dispersion relation as its predecessor, but at  $a\mathbf{p} = \mathbf{0}$  the cut-off effects for a heavy quark mass are much mitigated. The  $\rho = 1$  Brillouin overlap operator at  $am = 0$  still enjoys a rather good dispersion relation, and at  $a\mathbf{p} = \mathbf{0}$  the cut-off effects are equally small as for the Wilson overlap operator.

Evidently, the nice behavior of the free field dispersion relation of the Brillouin overlap operator at arbitrary quark mass and generic  $\rho$  is a necessary (and not a sufficient) condition for this formulation to be useful in real physics applications. However, given this property, we think it is worth while to investigate the Brillouin overlap action in more detail.

### 3 Kenney-Laub iterates for the matrix sign function

#### 3.1 Definition

Kenney and Laub proposed a family of iterations to compute the matrix sign function (equivalently to compute the unitary factor in the polar decomposition) that have some remarkable properties [28]. The  $(m, n)$  iteration for a matrix  $A$  with no purely imaginary eigenvalue is

$$X_{k+1} = X_k p_{mn}(I - X_k^2) [q_{mn}(I - X_k^2)]^{-1} \equiv f_{mn}(X_k), \quad X_0 = A \quad (9)$$

where  $r_{mn}(t) = p_{mn}(t)/q_{mn}(t)$  is the  $(m, n)$  Padé approximant to  $h(t) = (1 - t)^{-1/2}$ . Here  $I$  is the identity,  $p_{mn}$  is a polynomial of order  $m$  in  $t = 1 - x^2$  (or of order  $2m$  in  $x$ ), and  $q_{mn}$  is a polynomial of order  $n$  in  $t = 1 - x^2$  (or  $2n$  in  $x$ ). To compute the polar decomposition  $A = UP$  with unitary  $U$  and positive semi-definite  $P$  one simply replaces  $X_k^2 \rightarrow X_k^\dagger X_k$ .

In Tab. 1 the first few members  $f_{m,n}$  of this family are listed (which one also finds in the literature [28]) and in Tab. 2 the elements  $f_{n,n}$  with  $n = 4, \dots, 8$  are given. The convergence order [in  $k$ ] of the  $(m, n)$  element is  $m + n + 1$ . Two subsets of this family have special properties. First of all, the elements on the diagonal ( $m = n$ ) and first upper diagonal ( $n - m = 1$ ) are globally convergent, i.e. they work with any argument [28]. Second, the elements in the first column ( $n = 0$ ) do not require an inverse, but they tend to be numerically unstable. In fact, the element ( $m = 1, n = 0$ ) is the Newton-Schulz iteration for the matrix sign function, which

	$n = 0$	$n = 1$	$n = 2$	$n = 3$
$m = 0$	—	$\frac{2x}{1+x^2}$	$\frac{8x}{3+6x^2-x^4}$	$\frac{16x}{5+15x^2-5x^4+x^6}$
$m = 1$	$\frac{x(3-x^2)}{2}$	$\frac{x(3+x^2)}{1+3x^2}$	$\frac{4x(1+x^2)}{1+6x^2+x^4}$	$\frac{8x(3+5x^2)}{5+45x^2+15x^4-x^6}$
$m = 2$	$\frac{x(15-10x^2+3x^4)}{8}$	$\frac{x(15+10x^2-x^4)}{4(1+5x^2)}$	$\frac{x(5+10x^2+x^4)}{1+10x^2+5x^4}$	$\frac{2x(3+10x^2+3x^4)}{1+15x^2+15x^4+x^6}$
$m = 3$	$\frac{x(35-35x^2+21x^4-5x^6)}{16}$	$\frac{x(35+35x^2-7x^4+x^6)}{8(1+7x^2)}$	$\frac{x(35+105x^2+21x^4-x^6)}{2(3+42x^2+35x^4)}$	$\frac{x(7+35x^2+21x^4+x^6)}{1+21x^2+35x^4+7x^6}$

Table 1: Iteration functions  $f_{mn}(x) = x p_{mn}(x^2)/q_{mn}(x^2)$  for  $m, n = 0, \dots, 3$  from the Kenney-Laub family (9) for the matrix sign function. The element  $f_{00}(x) = x$  is not useful.

$m = n = 4$	$\frac{x(9+84x^2+126x^4+36x^6+x^8)}{1+36x^2+126x^4+84x^6+9x^8}$
$m = n = 5$	$\frac{x(11+165x^2+462x^4+330x^6+55x^8+x^{10})}{1+55x^2+330x^4+462x^6+165x^8+11x^{10}}$
$m = n = 6$	$\frac{x(13+286x^2+1287x^4+1716x^6+715x^8+78x^{10}+x^{12})}{1+78x^2+715x^4+1716x^6+1287x^8+286x^{10}+13x^{12}}$
$m = n = 7$	$\frac{x(15+455x^2+3003x^4+6435x^6+5005x^8+1365x^{10}+105x^{12}+x^{14})}{1+105x^2+1365x^4+5005x^6+6435x^8+3003x^{10}+455x^{12}+15x^{14}}$
$m = n = 8$	$\frac{x(17+680x^2+6188x^4+19448x^6+24310x^8+12376x^{10}+2380x^{12}+136x^{14}+x^{16})}{1+136x^2+2380x^4+12376x^6+24310x^8+19448x^{10}+6188x^{12}+680x^{14}+17x^{16}}$

Table 2: Diagonal iteration functions  $f_{nn}(x)$  for  $n = 4, \dots, 8$  from the Kenney-Laub family (9).

derives from the Newton method  $X_{k+1} = \frac{1}{2}(X_k + X_k^{-1})$  through an additional expansion of the inverse. The  $(m = 0, n = 1)$  element generates the inverses of the Newton-Schulz series, and the  $(m = 1, n = 1)$  element is sometimes named after Halley.

It seems to us that the diagonal mappings  $f_{n,n}$  with  $n \geq 1$  are most convenient for practical use. They have the special algebraic property that the polynomial  $q_{nn}(x)$  in the denominator is the *mirror polynomial* of  $p_{nn}(x)$  in the numerator, i.e. the coefficients show up in reverse order (e.g.  $5 + 10x^2 + x^4$  versus  $1 + 10x^2 + 5x^4$  in  $f_{2,2}$ ).

### 3.2 Principal Padé iteration functions

The diagonal  $(m = n)$  and first upper diagonal  $(m = n - 1)$  elements of the family (9) are singled out as the “principal Padé iteration functions”. For these  $m$  and  $n$  one defines

$$g_\ell(x) \equiv g_{m+n+1}(x) \equiv f_{m,n}(x) \quad (10)$$

which means that the index  $\ell$  counts them in Tab.1 in a zig-zag fashion, i.e.  $g_1(x) = x$ ,  $g_2(x) = 2x/(1+x^2)$ ,  $g_3(x) = x(3+x^2)/(1+3x^2)$ ,  $g_4(x) = 4x(1+x^2)/(1+6x^2+x^4)$ , and so on. These functions have a number of important properties [29]:

- (i) The coefficients of the numerator and the denominator follow from the binomial theorem

$$g_\ell(x) = \frac{(1+x)^\ell - (1-x)^\ell}{(1+x)^\ell + (1-x)^\ell} \quad (11)$$

and the numerator/denominator are thus the odd/even parts of  $(1+x)^\ell$ .

- (ii) This implies the following symmetry properties (for  $x > 0$ ) about  $x = 1$

$$g_{2n}(1/x) = g_{2n}(x) \quad [\text{upper diagonal}], \quad g_{2n+1}(1/x) = 1/g_{2n+1}(x) \quad [\text{diagonal}] \quad (12)$$

and ditto (for  $x < 0$ ) about  $x = -1$ , since the overall functions are all odd in  $x$ .

- (iii) All principal Padé iteration functions allow a  $\tanh(\cdot)$  representation, viz.

$$g_\ell(x) = \tanh(\ell \operatorname{artanh}(x)) . \quad (13)$$

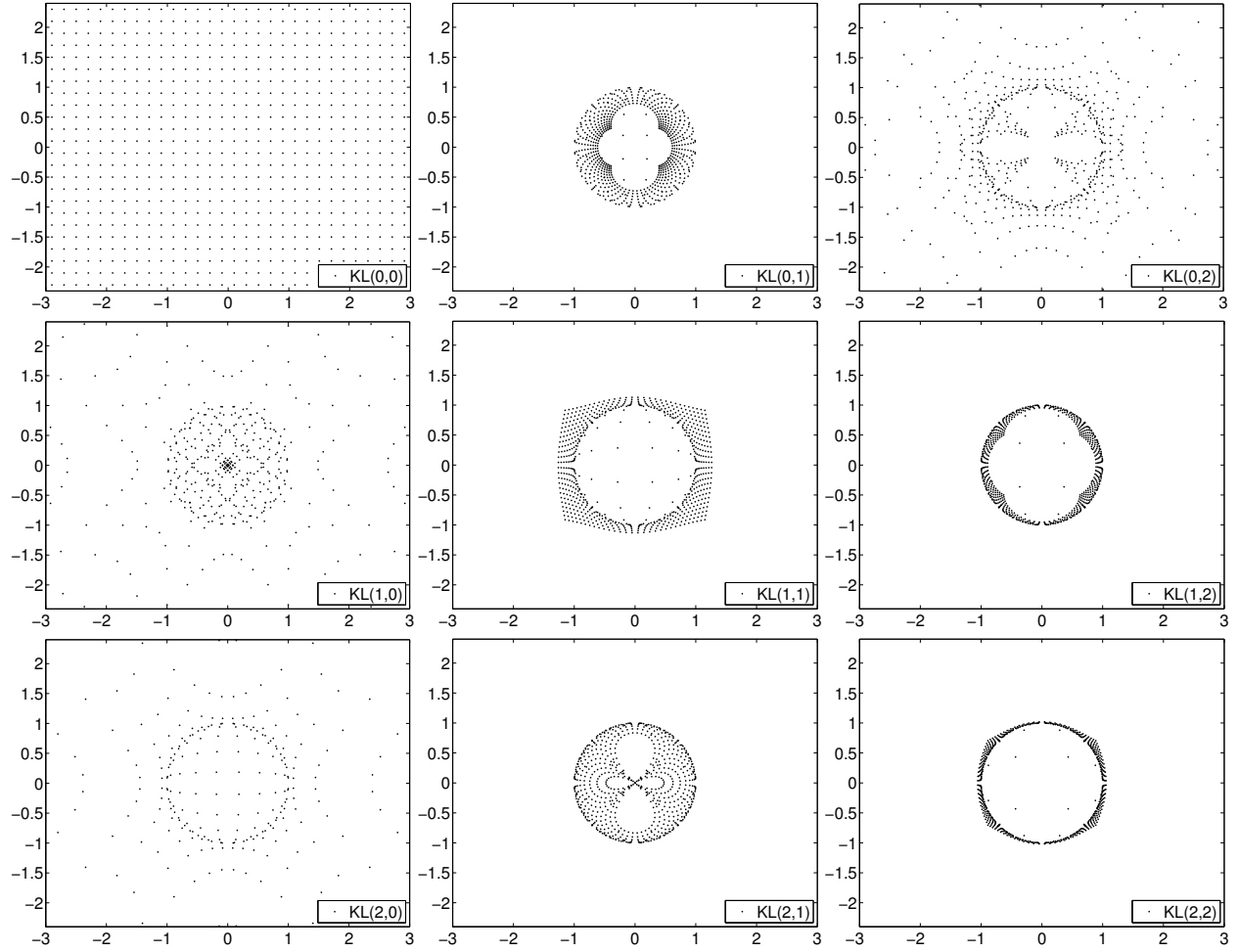


Figure 3: Image of regularly spaced points  $x \in \mathbf{C}$  with  $-3 \leq \text{Re}(x) \leq 3$  and  $-2.4 \leq \text{Im}(x) \leq 2.4$  under one iteration of  $f_{0,0}(\cdot)$  [identity, top left] to  $f_{2,2}(\cdot)$  [bottom right] as defined in Tab. 1.

(iv) Nesting two principal Padé iteration functions yields another one, viz.

$$g_{\ell'}(g_{\ell''}(x)) = g_{\ell'\ell''}(x) . \quad (14)$$

Since the product of two odd numbers is odd, and the product of two even numbers is even, it follows that both the diagonal and the first upper diagonal Kenney-Laub mappings satisfy this “semigroup property” separately. Obviously, this implies that nestings commute for principal Padé iteration functions, while this does not hold in general, i.e.  $f_{mn}(f_{pq}(x)) \neq f_{pq}(f_{mn}(x))$  for arbitrary  $m, n, p, q$ .

In addition,  $g_{\ell}(x)$  admits a partial fraction form which we will discuss in the next subsection.

The effect of the Kenney-Laub mappings for the unitary projection [i.e. for (9) with  $X_k^2 \rightarrow X_k^{\dagger} X_k$ ] may be visualized in the complex plane. The nine panels in Fig. 3 illustrate the effect of one operation of  $f_{m,n}$  with  $0 \leq m, n \leq 2$ , as given in Tab. 1. It seems plausible that the mappings on the diagonal and first upper diagonal are globally convergent, while mappings far away from the diagonal work only for near-unitary arguments  $X^{\dagger} X \simeq I$ .

To visualize the approximations to the sign function that derive from the Kenney-Laub mappings it suffices to restrict the discussion to  $x > 0$ , since each  $f_{m,n}(x)$  is an odd function



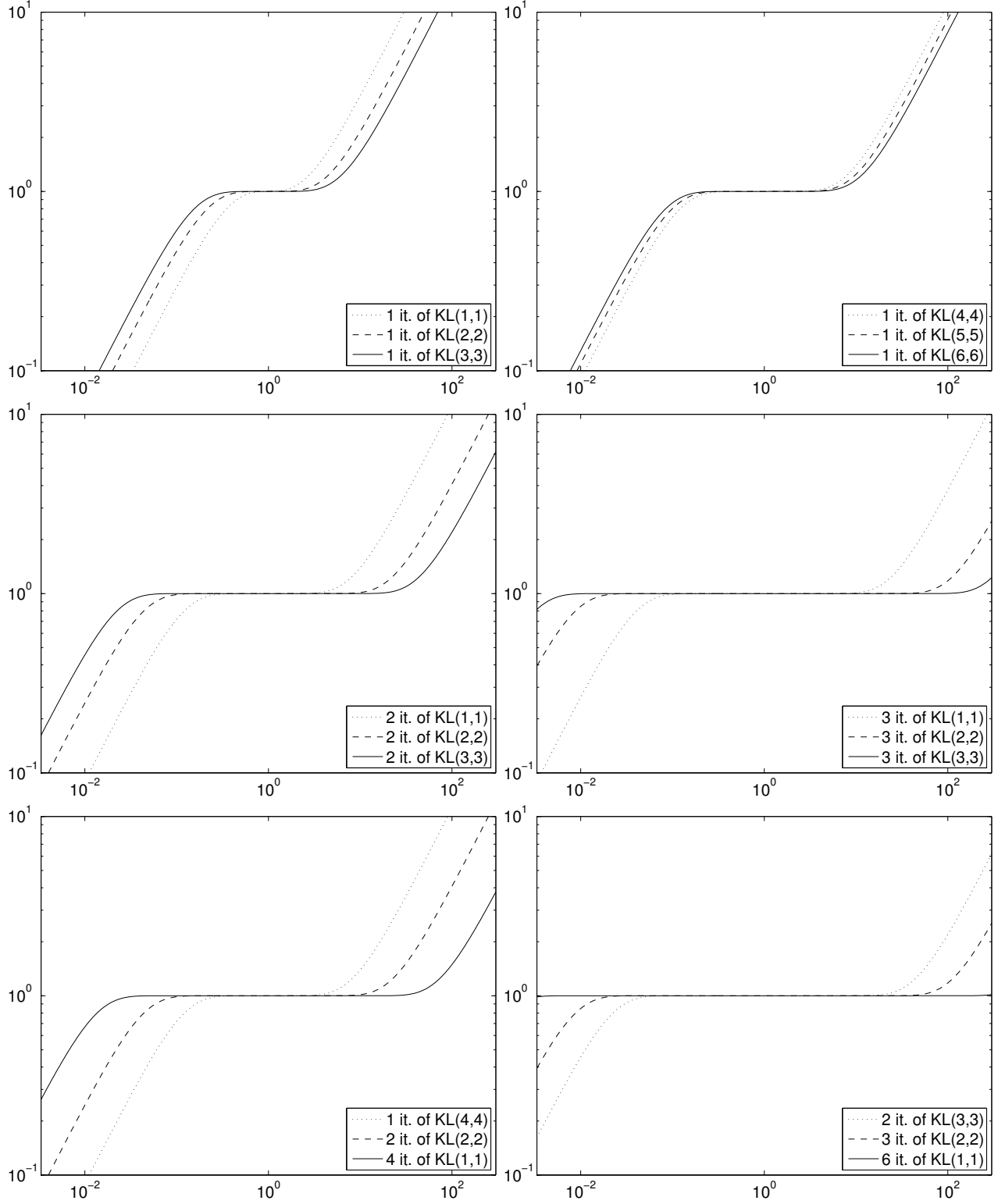


Figure 4: Image of the interval  $x \in [0.003, 300]$  under one or several iterations of the diagonal Kenney-Laub mappings  $x \rightarrow f_{nn}(x) = x p_{nn}(x^2)/q_{nn}(x^2)$  in log-log representation.

of  $x$ . The panels of Fig. 4) illustrate various combinations of diagonal ( $m = n$ ) mappings on the interval  $x \in [0.003, 300]$ . Evidently, these approximations work best for  $x \simeq 1$ , with monotonically decreasing quality for small ( $x \ll 1$ ) and large ( $x \gg 1$ ) arguments. That this decrease in quality is symmetric about 1 is a direct consequence of (12).

In lattice QCD the elements of the first upper diagonal have been used before [30]. The most obvious difference to the diagonal mappings which we advocate is that the former set of functions assumes a maximum/minimum at  $x = \pm 1$ , respectively, while the diagonal functions  $f_{n,n}(x) = g_{2n+1}(x)$  increase without any bound. In a similar vein we emphasize that – unlike optimal rational approximations [31–33] – diagonal Kenney-Laub functions show *no wiggles*; the value  $\pm 1$  at  $x = \pm 1$  is approached monotonically, both from the origin and from  $\pm\infty$ .

### 3.3 Partial fraction and continued fraction representations

In view of numerical applications let us rewrite the diagonal Kenney-Laub mappings in partial fraction form. The first two diagonal mappings can be brought into the form

$$g_3(X) = f_{1,1}(X) = \frac{X}{3} \left( 1 + \frac{8/3}{X^\dagger X + 1/3} \right) \quad (15)$$

$$g_5(X) = f_{2,2}(X) = \frac{X}{5} \left( 1 + \frac{4(1 - 1/\sqrt{5})}{X^\dagger X + 1 - 2/\sqrt{5}} + \frac{4(1 + 1/\sqrt{5})}{X^\dagger X + 1 + 2/\sqrt{5}} \right) \quad (16)$$

while for higher  $n$  the roots of the denominator polynomial can only be given over the field of complex numbers (though they happen to be real). The general formula reads [29]

$$g_{2n+1}(x) = \frac{x}{2n+1} \sum_{i=0}^n \frac{2 - \delta_{i,n}}{\sin^2((2i+1)\pi/(4n+2)) + \cos^2((2i+1)\pi/(4n+2))x^2} \quad (17)$$

and from the explicit form provided in Tab. 6 of the appendix it is easy to see that the smallest shift gets progressively smaller with increasing  $n$ . Moreover, the coefficients in the numerator are all positive, and they grow synchronously with the shift in the denominator. This formula is reminiscent of the one for the first upper diagonal [29, 30]

$$g_{2n}(x) = \frac{x}{n} \sum_{i=0}^{n-1} \frac{1}{\sin^2((2i+1)\pi/(4n)) + \cos^2((2i+1)\pi/(4n))x^2} \quad (18)$$

except that the former expression has a constant contribution ( $i = n$ ), while the latter one has not. The bottom line is that one can use a multi-shift conjugate gradient (CG) solver to evaluate  $f_{n,n}(X)v$  on a given vector  $v$  [34, 35]. In our view it is convenient that the coefficients can be worked out beforehand, i.e. independent of the spectral properties of  $A \equiv X^\dagger X$ .

For  $f_{n,n}$  with  $n \geq 2$  also a continued fraction representation can be given, for instance

$$f_{2,2}(X) = \frac{X}{5} \left( 1 + \frac{8}{X^\dagger X + 7/5 - \frac{16/25}{X^\dagger X + 3/5}} \right) \quad (19)$$

$$f_{3,3}(X) = \frac{X}{7} \left( 1 + \frac{16}{X^\dagger X + 3 - \frac{24/7}{X^\dagger X + 5/3 - \frac{8/63}{X^\dagger X + 1/3}}} \right) \quad (20)$$

$$f_{4,4}(X) = \frac{X}{9} \left( 1 + \frac{80/3}{X^\dagger X + 77/15 - \frac{264/25}{X^\dagger X + 139/45 - \frac{520/891}{X^\dagger X + 103/117 - \frac{96/1859}{X^\dagger X + 3/13}}}} \right) \quad (21)$$

$$f_{5,5}(X) = \frac{X}{11} \left( 1 + \frac{40}{X^\dagger X + 39/5 - \frac{624/25}{X^\dagger X + 73/15 - \frac{160/99}{X^\dagger X + 61/39 - \frac{408/1859}{X^\dagger X + 131/221 - \frac{8/289}{X^\dagger X + 3/17}}}} \right) \quad (22)$$

$$f_{6,6}(X) = \frac{X}{13} \left( 1 + \frac{56}{X^\dagger X + 11 - \frac{352/7}{X^\dagger X + 7 - \frac{272/77}{X^\dagger X + 31/13 - \frac{1064/1859}{X^\dagger X + 227/221 - \frac{616/5491}{X^\dagger X + 53/119 - \frac{16/931}{X^\dagger X + 1/7}}}}} \right) \quad (23)$$

where all coefficients are found to be given by (small-over-small) rational numbers.

## 4 Overlap operator construction

Given any undoubled (or doubled but with one chirality in the physical branch) “kernel” Dirac operator  $D_m^{\text{ke}}$  at a quark mass  $m$ , the massless overlap operator  $D^{\text{ov}}$  is defined as a backshifted version of the (unique) unitary part of the kernel at negative mass  $-\rho/a$  [21, 22]

$$aD^{\text{ov}} \equiv aD_0^{\text{ov}} = \begin{cases} \rho \left[ aD_{-\rho/a}^{\text{ke}} (a^2 D_{-\rho/a}^{\text{ke}\dagger} D_{-\rho/a}^{\text{ke}})^{-1/2} + 1 \right] = \rho \left[ \gamma_5 \text{sign}(\gamma_5 aD_{-\rho/a}^{\text{ke}}) + 1 \right] \\ \rho \left[ (a^2 D_{-\rho/a}^{\text{ke}} D_{-\rho/a}^{\text{ke}\dagger})^{-1/2} aD_{-\rho/a}^{\text{ke}} + 1 \right] = \rho \left[ \text{sign}(aD_{-\rho/a}^{\text{ke}} \gamma_5) \gamma_5 + 1 \right] \end{cases} \quad (24)$$

where  $0 < \rho < 2$  is an arbitrary parameter (its canonical value is 1). The equivalence of the two lines follows from the singular value decomposition  $aD_{-\rho/a}^{\text{ke}} = USV^\dagger$  with unitary  $U, V$  and  $S > 0$ , by means of which  $a^2 D_{-\rho/a}^{\text{ke}\dagger} D_{-\rho/a}^{\text{ke}} = VS^2V^\dagger$  and  $a^2 D_{-\rho/a}^{\text{ke}} D_{-\rho/a}^{\text{ke}\dagger} = US^2U^\dagger$ . This implies  $aD^{\text{ov}} = \rho[USV^\dagger VS^{-1}V^\dagger + 1] = \rho[UV^\dagger + 1]$  and  $aD^{\text{ov}} = \rho[US^{-1}U^\dagger USV^\dagger + 1] = \rho[UV^\dagger + 1]$ , respectively, which completes the proof. Note that the reformulation in terms of the matrix sign function in eqn. (24) holds only if the kernel is  $\gamma_5$ -hermitean, i.e.  $\gamma_5 D^{\text{ke}} \gamma_5 = D^{\text{ke}\dagger}$ .

The massless overlap operator (24) fulfills the Ginsparg-Wilson (GW) relation [14]

$$D\gamma_5 + \gamma_5 D = \frac{a}{\rho} D\gamma_5 D \iff D\gamma_5 \left(1 - \frac{aD}{2\rho}\right) + \left(1 - \frac{aD}{2\rho}\right) \gamma_5 D = 0 \quad (25)$$

and  $D^{\text{ov}}$  is thus said to be “chirally symmetric”, regardless the details of the kernel. In practice there is a choice to be made regarding the type of kernel (e.g. Wilson or Brillouin), how much link smearing one wants to apply, and whether the kernel shall be equipped with a clover term. Whenever  $D^{\text{ke}}$  is  $\gamma_5$ -hermitean, this property extends to  $D^{\text{ov}}$ , and in this case multiplying (25) with  $\gamma_5$  from the left or the right yields (note that  $[D, D^\dagger] = 0$  is implied)

$$D^\dagger + D = \frac{a}{\rho} D^\dagger D = \frac{a}{\rho} D D^\dagger \iff D^\dagger \left(1 - \frac{aD}{2\rho}\right) + \left(1 - \frac{aD}{2\rho}\right)^\dagger D = 0. \quad (26)$$

### 4.1 Kenney-Laub iterates of shifted Dirac kernels

For the sake of clarity let us consider the use of a diagonal Kenney-Laub mapping  $f_{n,n}$  to define, for a given kernel  $aD^{\text{ke}}$ , an approximation to the overlap operator (24). With  $X = aD^{\text{ke}} - \rho$  the

	$f_{1,1}^{(k)}$	$f_{2,2}^{(k)}$	$f_{3,3}^{(k)}$	$f_{4,4}^{(k)}$	$f_{5,5}^{(k)}$	$f_{6,6}^{(k)}$
$k = 1$	$2.7 \cdot 10^{-3}$	$3.4 \cdot 10^{-5}$	$4.2 \cdot 10^{-7}$	$5.2 \cdot 10^{-9}$	$6.3 \cdot 10^{-11}$	$7.9 \cdot 10^{-13}$
$k = 2$	$5.2 \cdot 10^{-9}$	$2.8 \cdot 10^{-24}$	$3.5 \cdot 10^{-47}$	$1.0 \cdot 10^{-77}$	$6.9 \cdot 10^{-116}$	$1.1 \cdot 10^{-161}$
$k = 3$	$3.4 \cdot 10^{-26}$	$1.0 \cdot 10^{-119}$	$1.0 \cdot 10^{-327}$	$4.6 \cdot 10^{-696}$	$1.6 \cdot 10^{-1270}$	$6.8 \cdot 10^{-2097}$

Table 3: Image of the would-be zero-mode  $\lambda = 0.2$  under 1 to 3 iterations of the diagonal Kenney-Laub mappings  $\lambda \rightarrow f_{n,n}(\lambda-1) + 1$  as defined in Tabs. 1, 2. The order of convergence (in  $k$ ) of the columns is 3, 5, 7, 9, 11, 13, respectively.

relation  $Y = X(X^\dagger X + 3)(3X^\dagger X + 1)^{-1}$  defines a Dirac operator  $aD^{\text{it}} = \rho[Y + 1]$  with improved chiral symmetry. After another iteration, which may involve a different mapping, for instance  $Z = Y(Y^\dagger Y Y^\dagger Y + 10Y^\dagger Y + 5)(5Y^\dagger Y Y^\dagger Y + 10Y^\dagger Y + 1)^{-1}$ , the redefinition  $aD^{\text{it}} = \rho[Z + 1]$  yields a Dirac operator with an even smaller violation of the relation (25).

In usual applications one cannot hold any of these matrices in memory. The challenge is thus to implement the forward application  $D^{\text{it}}x$  on a given vector  $x$  in such a form that everything boils down to repeated matrix-vector multiplications of the form  $D^{\text{ke}\dagger} D^{\text{ke}} y$  and  $D^{\text{ke}} z$ .

For the diagonal iterations  $f_{n,n}$  the image of the “would-be zero-mode”  $\lambda = 0.2$  of  $D^{\text{ke}}$  under 1 to 3 iterations is summarized in Tab. 3. For such a mode and in double-precision arithmetics the mappings  $f_{1,1}$  and  $f_{2,2}$  achieve exact chiral symmetry after 3 and 2 iterations, respectively. Note that the nesting formula (14) says  $f_{1,1}^{(3)} = g_3^{(3)} = g_{27}^{(1)} = f_{13,13}^{(1)}$  and  $f_{2,2}^{(2)} = g_5^{(2)} = g_{25}^{(1)} = f_{12,12}^{(1)}$ .

## 4.2 Massive overlap action – traditional version

Let  $\lambda$  be an eigenvalue of the massless overlap operator  $aD^{\text{ov}}$  with parameter  $\rho$ , i.e.  $\lambda = \rho(1 + e^{i\varphi})$  with  $\varphi \in ]0, 2\pi[$ . This circular eigenvalue spectrum is mapped onto the imaginary axis through the stereographic projection  $\lambda \rightarrow \tilde{\lambda} \equiv \lambda/(1 - \lambda/[2\rho]) = 2i\rho/\tan(\varphi/2)$ . The massive overlap operator follows by shifting this line by  $am$  to the right, and inverting the mapping.

The traditional way of doing this is to multiply  $\tilde{\lambda} + am$  with the factor  $(1 - \lambda/[2\rho])$  which then leads to  $\lambda + am(1 - \lambda/[2\rho])$ . In operator language this means that the massive overlap operator  $D_m^{\text{ov}}$  follows by adding a “chirally rotated” scalar term [21, 22]

$$\begin{aligned} D_m^{\text{tra}} &\equiv D^{\text{ov}} + m\left(1 - \frac{a}{2\rho} D^{\text{ov}}\right) = \left(1 - \frac{am}{2\rho}\right) D^{\text{ov}} + m \\ &= \left(\frac{\rho}{a} - \frac{m}{2}\right) \gamma_5 \text{sign}(\gamma_5 a D_{-\rho/a}^{\text{ke}}) + \left(\frac{\rho}{a} + \frac{m}{2}\right) \end{aligned} \quad (27)$$

which yields an operator with a circular eigenvalue spectrum of radius  $\rho - am/2$  around the point  $(\rho + am/2, 0)$  in the complex plane. Obviously this implies the constraint  $am < 2\rho$ . An ad hoc way of removing this constraint would be to replace  $m$  in (27) by  $\tilde{m} \equiv 1/(1/m + a/[2\rho]) = m/(1 + am/[2\rho])$ , so that  $a\tilde{m} \leq 2\rho$  for all  $am$ . With the traditional definition (27) solving the massive Dirac equation  $D_m^{\text{tra}} x = b$  for  $x$  with a given right-hand side  $b$  is equivalent to solving

$$(D^{\text{ov}} + \tilde{m}) x = \tilde{b} \quad \text{with} \quad \tilde{m} = \frac{m}{1 - am/[2\rho]} \quad \text{and} \quad \tilde{b} = \frac{b}{1 - am/[2\rho]} = \frac{\tilde{m}}{m} b \quad (28)$$

for  $x$ , with the massless  $D^{\text{ov}}$  defined in (24).

### 4.3 Massive overlap action – complete version

Alternatively, one might start from the proper inversion of the stereographic mapping, which is  $\tilde{\lambda} \rightarrow \tilde{\lambda}/(1 + \tilde{\lambda}/[2\rho]) = \lambda$ , and by adding the mass to  $\tilde{\lambda}$  one ends up with the massive spectrum

$$\frac{\tilde{\lambda} + am}{1 + \tilde{\lambda}/[2\rho] + am/[2\rho]} = \frac{\lambda/(1 - \lambda/[2\rho]) + am}{1 + \lambda/(1 - \lambda/[2\rho])/[2\rho] + am/[2\rho]} = \frac{\lambda + am(1 - \lambda/[2\rho])}{1 + am(1 - \lambda/[2\rho])/[2\rho]}$$

which does not entail any constraint on  $am$  (with  $am \rightarrow \infty$  the eigenvalue spectrum shrinks to a point at  $2\rho$ ). In operator language this means that the complete definition

$$\begin{aligned} D_m^{\text{com}} &\equiv \frac{D^{\text{ov}} + m\left(1 - \frac{a}{2\rho}D^{\text{ov}}\right)}{1 + \frac{am}{2\rho}\left(1 - \frac{a}{2\rho}D^{\text{ov}}\right)} = \frac{(1 - \frac{am}{2\rho})D^{\text{ov}} + m}{(1 + \frac{am}{2\rho}) - \frac{a^2m}{2\rho^2}D^{\text{ov}}} \\ &= \frac{(\frac{\rho}{a} - \frac{m}{2})\gamma_5 \text{sign}(\gamma_5 a D_{-\rho/a}^{\text{ke}}) + (\frac{\rho}{a} + \frac{m}{2})}{1 - \frac{am}{2\rho}\gamma_5 \text{sign}(\gamma_5 a D_{-\rho/a}^{\text{ke}})} \end{aligned} \quad (29)$$

looks superficially similar to the Moebius kernel that was proposed for the massless case [36]. Note that the fractional notation in (29) is well-defined, since the normality of  $D^{\text{ov}}$  ensures that the numerator and the inverse of the denominator would commute. Hence with the definition (29) solving the massive Dirac equation  $D_m^{\text{com}}x = b$  for  $x$  with a given  $b$  amounts to solving

$$(D^{\text{ov}} + \tilde{m})x = \tilde{b} \quad \text{with} \quad \tilde{m} = \frac{m}{1 - am/[2\rho]} \quad \text{and} \quad \tilde{b} = \frac{1 + am/[2\rho](1 - aD^{\text{ov}}/[2\rho])}{1 - am/[2\rho]} b \quad (30)$$

for the vector  $x$ , with the massless  $D^{\text{ov}}$  defined in (24). A comparison with (28) shows that the procedure is the same, except that the right-hand side  $\tilde{b}$  is now defined in a different manner.

### 4.4 Massive overlap action – proof of equivalence

It is well known that the traditional form (28) of the massive overlap action is to be used in conjunction with a “chiral symmetry ensuring factor”  $(1 - aD/[2\rho])$  to be attached to the external densities  $S$ ,  $P$  and currents  $V_\mu$ ,  $A_\mu$ . This leads to an effective Green’s function (propagator)

$$S_m^{\text{tra}} = \frac{1 - \frac{a}{2\rho}D^{\text{ov}}}{D^{\text{ov}} + m\left(1 - \frac{a}{2\rho}D^{\text{ov}}\right)} = \frac{1}{\frac{D^{\text{ov}}}{1 - \frac{a}{2\rho}D^{\text{ov}}} + m} \quad (31)$$

which, thanks to the “extra prescription”, has the same form as in the continuum [37–40].

In the complete approach, when inverting (29) without any extra factors one arrives at

$$S_m^{\text{com}} = \frac{1 + \frac{am}{2\rho}\left(1 - \frac{a}{2\rho}D^{\text{ov}}\right)}{D^{\text{ov}} + m\left(1 - \frac{a}{2\rho}D^{\text{ov}}\right)} = \frac{\frac{1}{1 - \frac{a}{2\rho}D^{\text{ov}}} + \frac{am}{2\rho}}{\frac{D^{\text{ov}}}{1 - \frac{a}{2\rho}D^{\text{ov}}} + m} \quad (32)$$

which differs from (31) by just a contact term ( $I$  is the identity)

$$S_m^{\text{com}} - S_m^{\text{tra}} = \frac{a}{2\rho}I. \quad (33)$$

In short, we recommend to abandon the definition (27, 28) and to use (29, 30) instead. In this complete form chiral symmetry is genuinely built in, and there is no need for invoking any “extra prescription” if decay constants and other matrix elements are to be determined.

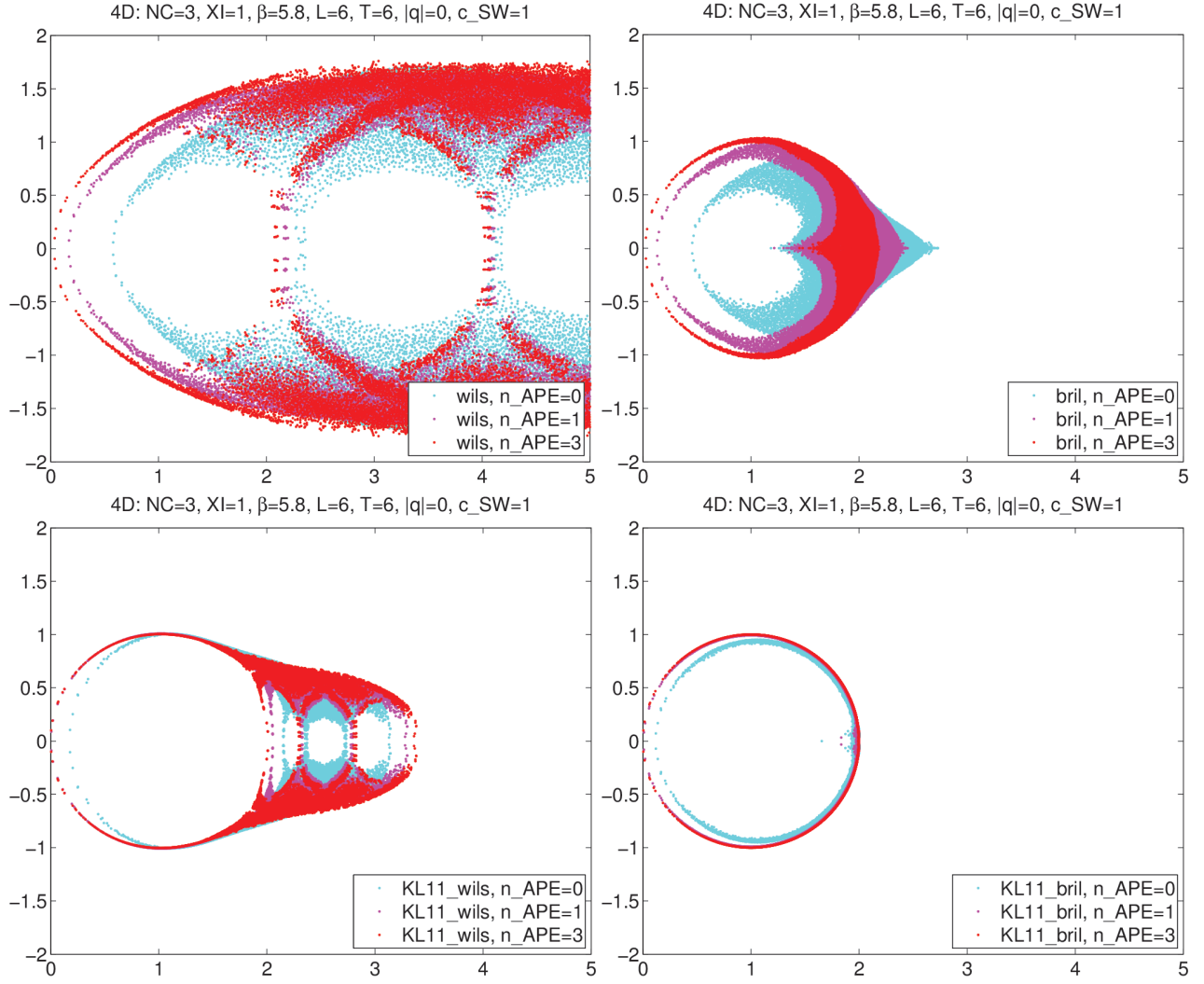


Figure 5: *Eigenvalue spectra of the Wilson (left) and Brillouin (right) operators with 0,1,3 APE smearings before (top) and after (bottom) one iteration of  $f_{1,1}$ . The titles specify the number of colors ( $N_c = 3$ ), the anisotropy coefficient ( $\xi = 1$  means  $a_s = a_t$ ), the box size ( $N_s/a = N_t/a = 6$ ), the absence of topological charge ( $q_{\text{top}} = 0$ ), and tree-level clover improvement ( $c_{\text{SW}} = 1$ ).*

## 5 Eigenvalue spectra with Wilson and Brillouin kernel

To gain an understanding of the difference between an approximate overlap operator with Wilson kernel and the same fixed-order approximant with the Brillouin kernel it is useful to take a look at the eigenvalue spectrum of either kernel on a given background.

The first row of Fig. 5 displays such eigenvalue spectra on a thermalized SU(3) gauge configuration. The Wilson operator has 5 branches with multiplicities 1,4,6,4,1 (from left to right, the last two are cut off), respectively. Only the first (leftmost) branch contributes to continuum physics. The Brillouin operator has only two branches, with multiplicities 1,15, respectively. Again, only the first branch contributes in the continuum, but the advantage is that the unphysical species are more condensed; they all sit near  $a\lambda = 2$  (which proves useful in the overlap projection, see below). The figures show the effect of the link smearing combined with tree-level

(that is  $c_{\text{SW}} = 1$ ) clover improvement. In the Wilson case the horizontal “jitter” in the physical branch gets ameliorated by the smearing; after 3 smearings the segment of the physical branch close to the origin is fairly close to a GW circle. Also in the Brillouin case both the additive mass shift and the “jitter” in the physical branch get reduced by the smearing; after 3 steps the Brillouin eigenvalue spectrum looks similar to that of a “parameterized fixed point action” (which is the practical implementation of the “perfect action”) [15, 16, 41–45].

The second row of Fig. 5 displays the eigenvalues of the Kenney-Laub iterate  $f_{1,1}$  of the two kernels at  $\rho = 1$ . With either kernel the eigenvalue spectrum gets attracted (compared to the first row) towards the unit circle, but the effect is more stringent with the Brillouin kernel. With the Wilson kernel (left) there is a significant left-over from the 15 unphysical branches, now at  $1.9 < \text{Re}(z) < 3.4$ . With the Brillouin kernel (right) the eigenvalue spectrum is essentially a GW circle, at least if the version with smearing is considered. Clearly, further iterations of the Kenney-Laub mapping (tantamount to higher  $n$  in  $f_{n,n}$ ) will bring the eigenvalue spectrum of the resulting operator arbitrarily close to a GW circle, and a higher value of  $n$  is needed with the Wilson kernel to reach a certain level of proximity than with the Brillouin kernel.

## 6 Spectral flow with Wilson and Brillouin kernel

Both kernels considered (Wilson or Brillouin) are  $\gamma_5$ -hermitean, but neither one is normal (i.e.  $[D^{\text{ke}}, D^{\text{ke}\dagger}] \neq 0$  for either  $D^{\text{ke}} = D^{\text{wils}}$  or  $D^{\text{ke}} = D^{\text{bril}}$ ). Accordingly, the spectral properties of  $D^{\text{ke}}$  and  $\gamma_5 D^{\text{ke}}$  cannot be deduced from each other.

Fig. 6 shows the eigenvalues of  $\gamma_5 D^{\text{wils}}$  on one gauge configuration for a scan of  $m_0$  in the range from  $-2$  to  $0$ , with and without link smearing, as well as with and without a clover term. Ideally, one wants to choose all tunable parameters such that the “eye” in the hermitean eigenvalue spectrum (the leftmost “bay” in Fig. 6; the “open sea” to the right of the “straights” is not shown) is wide open. Clearly, link smearing helps a lot in this respect. In comparison, the choice  $c_{\text{SW}} = 0$  versus  $c_{\text{SW}} = 1$  seems less important. Still, what speaks in favor of a clover term in the Wilson kernel is that the “magic” value  $\rho \simeq 0.634$  of Sec. 2 then fares reasonably, while without the clover term the opening of the “eye” is far from optimal at this value of  $\rho$ .

Fig. 7 shows the eigenvalues of  $\gamma_5 D^{\text{bril}}$  on one gauge configuration for a scan of  $m_0$  in the range from  $-2$  to  $0$ , with and without link smearing, as well as with and without a clover term. Again, link smearing is found to have a very beneficial effect on the opening of the “eye”. Also with the Brillouin kernel the choice  $c_{\text{SW}} = 0$  versus  $c_{\text{SW}} = 1$  seems insignificant regarding the maximum width of the “eye”, but it affects the position (i.e. the value of  $-m_0 = \rho$ ) at which the maximum is realized. Interestingly, with  $c_{\text{SW}} = 1$  the “magic” value  $\rho \simeq 0.634$  of Sec. 2 more-or-less coincides with the choice of  $\rho$  which maximizes the opening of the “eye”.

In Fig. 8 we show similar eigenvalues (in fact eigenvalues of  $\gamma_5 D^{\text{ke}} \gamma_5 D^{\text{ke}}$ , i.e. without the sign information) on much larger lattices (the  $40^3 \times 64$  lattices by QCDSF that will be discussed in Sec. 7). In view of the lesson just learned, we restrict ourselves to the version with link smearing. The spacing in  $m_0 = -\rho$  is too wide to allow for individual eigenvalue tracking. Still, it is evident that the main difference between the Wilson and the Brillouin kernel is the *upper* end of the eigenvalue spectrum; with the Brillouin kernel it is at least an order of magnitude lower. What matters for the CPU time spent in large-scale computations is the effective condition number  $\lambda_{\text{max}}/\lambda_n$  after  $n - 1$  low modes are projected away (we show the situation for  $n = 1, 10, 30, 100$ ). It seems on such big lattices the difference between the Wilson

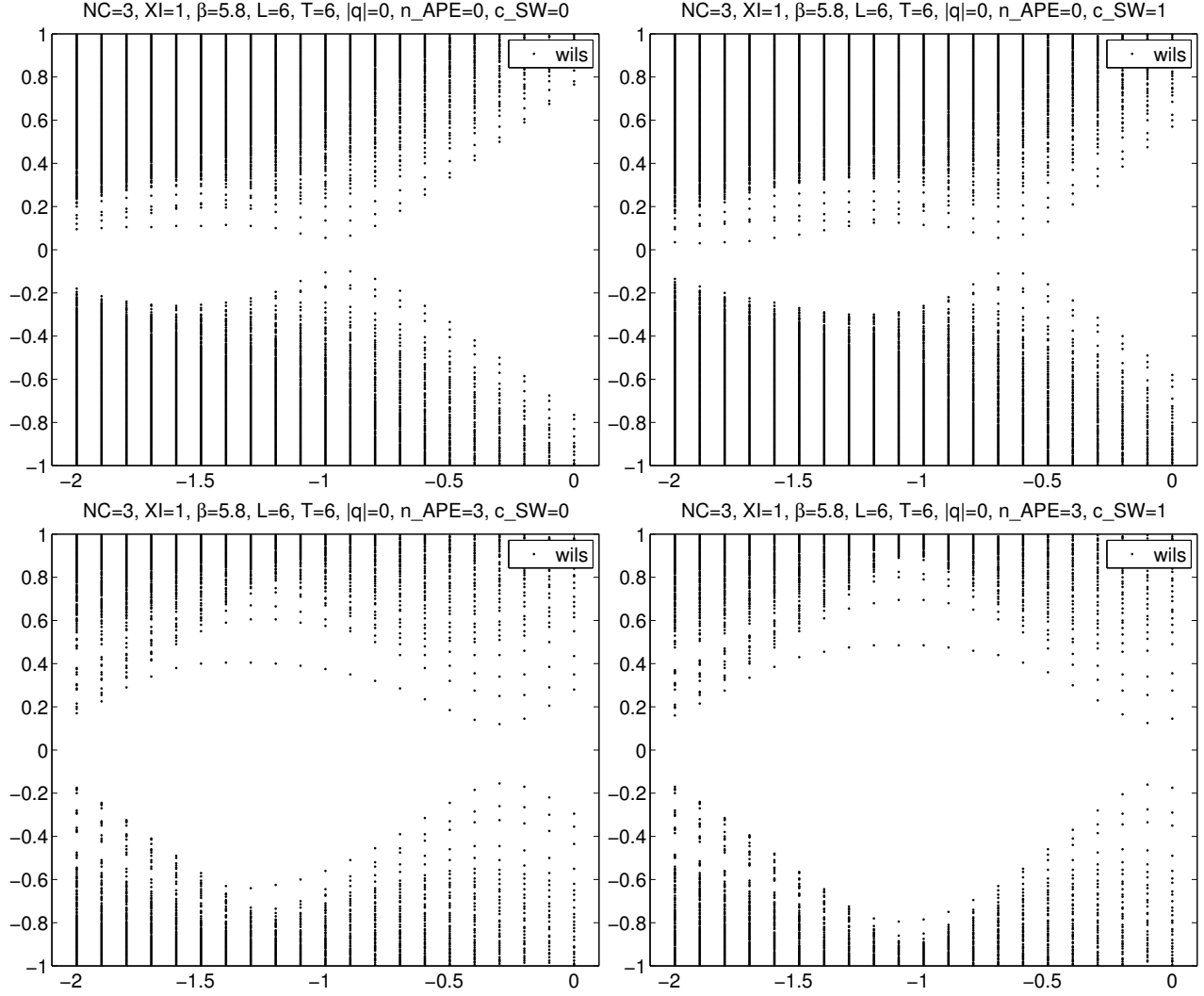


Figure 6: Spectrum of  $\gamma_5 D_{-\rho}^{\text{wils}}$  on one gauge configuration as a function of  $m_0 = -\rho$  with  $c_{\text{SW}} = 0$  (left) or  $c_{\text{SW}} = 1$  (right) and 0 (top) or 3 (bottom) APE smearings. Compare to Fig. 7.

and the Brillouin kernel is less pronounced than it appeared on the small lattices. Still, it is encouraging to see that with the Brillouin kernel the “magic” choice  $\rho \simeq 0.634$  fares well, both for  $c_{\text{SW}} = 0$  and  $c_{\text{SW}} = 1$ .

## 7 Numerical tests with Brillouin and Wilson kernel

The massless overlap operator  $D^{\text{ov}}$  as defined in (24) differs from the kernel  $D^{\text{ke}}$  in several ways: (i)  $D^{\text{ov}}$  is normal, i.e. it commutes with  $D^{\text{ov}\dagger}$ , (ii)  $D^{\text{ov}}$  satisfies the GW relation (25), (iii)  $D^{\text{ov}}$  is not ultralocal but just exponentially localized (with the fall-off pattern being a measure of the quality of the resulting operator). Here we verify these properties numerically on matched quenched lattices (i.e. with a fixed physical box size  $L$ ), using clover improved kernels ( $c_{\text{SW}} = 1$ ) and 1 or 3 steps of  $\alpha = 0.72$  APE smearing. In addition, we explore the inversion cost of the fixed-order Kenney-Laub overlap operator on large  $N_f = 2$  lattices generated by QCDSF.



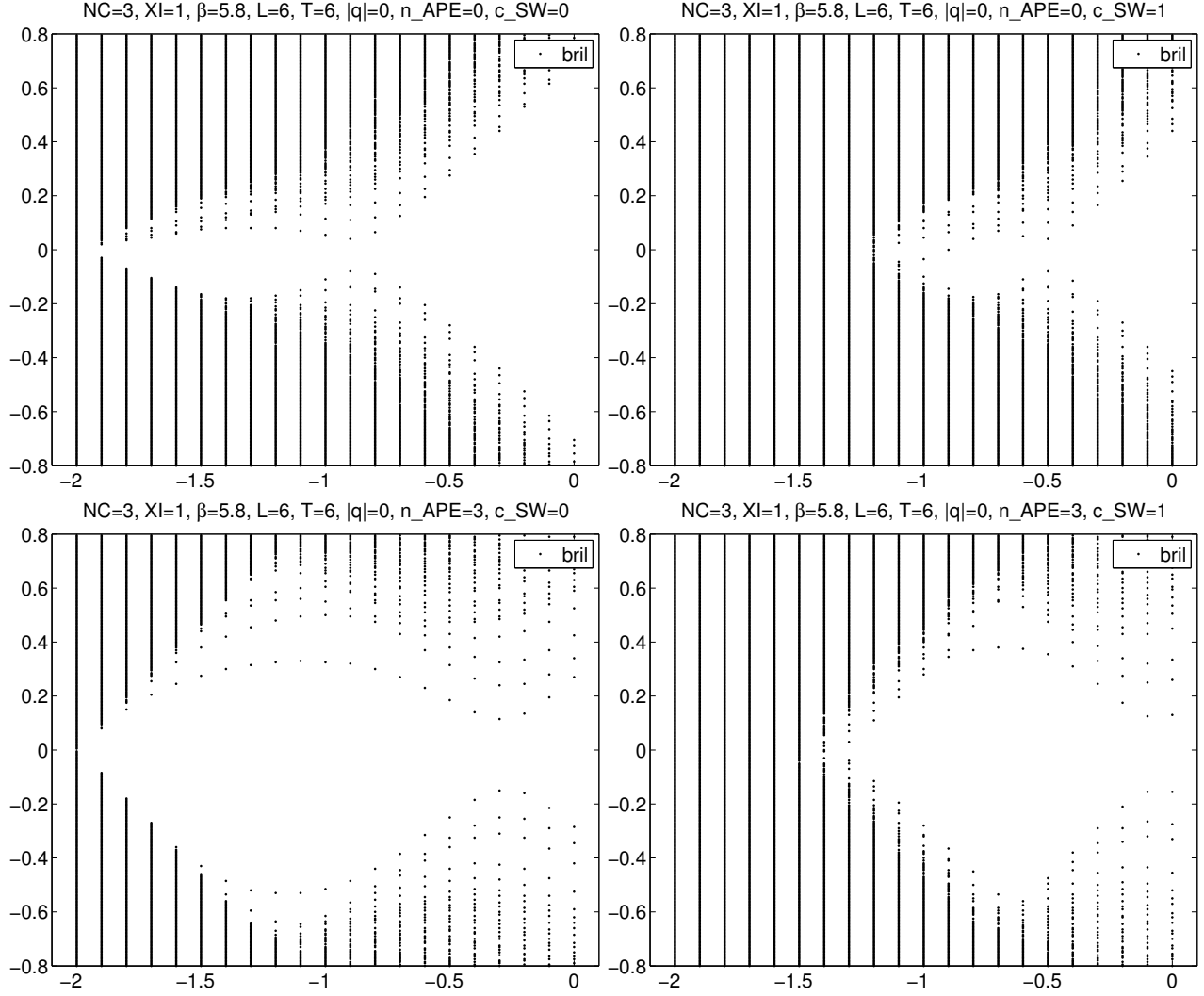


Figure 7: Same as Fig. 6, but now for  $\gamma_5 D_{-p}^{\text{bril}}$ . For the title details see the caption of Fig. 5.

## 7.1 Operator normality

We select the fixed rational approximation to the sign function implied by the Kenney-Laub iterate  $D = D_{-1}^{(1)} + 1$  with  $D_{-1}^{(1)} = f_{1,1}(D_{-1}^{(0)})$  and  $D^{(0)}$  being the Wilson or Brillouin kernel. We measure  $\|(DD^\dagger - D^\dagger D)\eta\|$  for a few dozen normalized Gaussian random vectors  $\eta$  on 40 configs for each  $\beta$  used in Ref. [3], and Fig. 9 shows the result. Both in the Wilson and in the Brillouin case, the operator with 3 steps of link smearing in the kernel exhibits smaller deviations from normality than the one with 1 step of smearing in the kernel. The main lesson to be learned is that both operators with Brillouin kernel have a smaller violation of normality than the two operators with Wilson kernel. Evidently, in order to reach a fixed level of normality violation, e.g.  $\|(DD^\dagger - D^\dagger D)\eta\| < 10^{-12}$ , the order of the rational approximation must be enhanced most drastically for the unsmeared Wilson kernel and least so for the smeared Brillouin kernel.

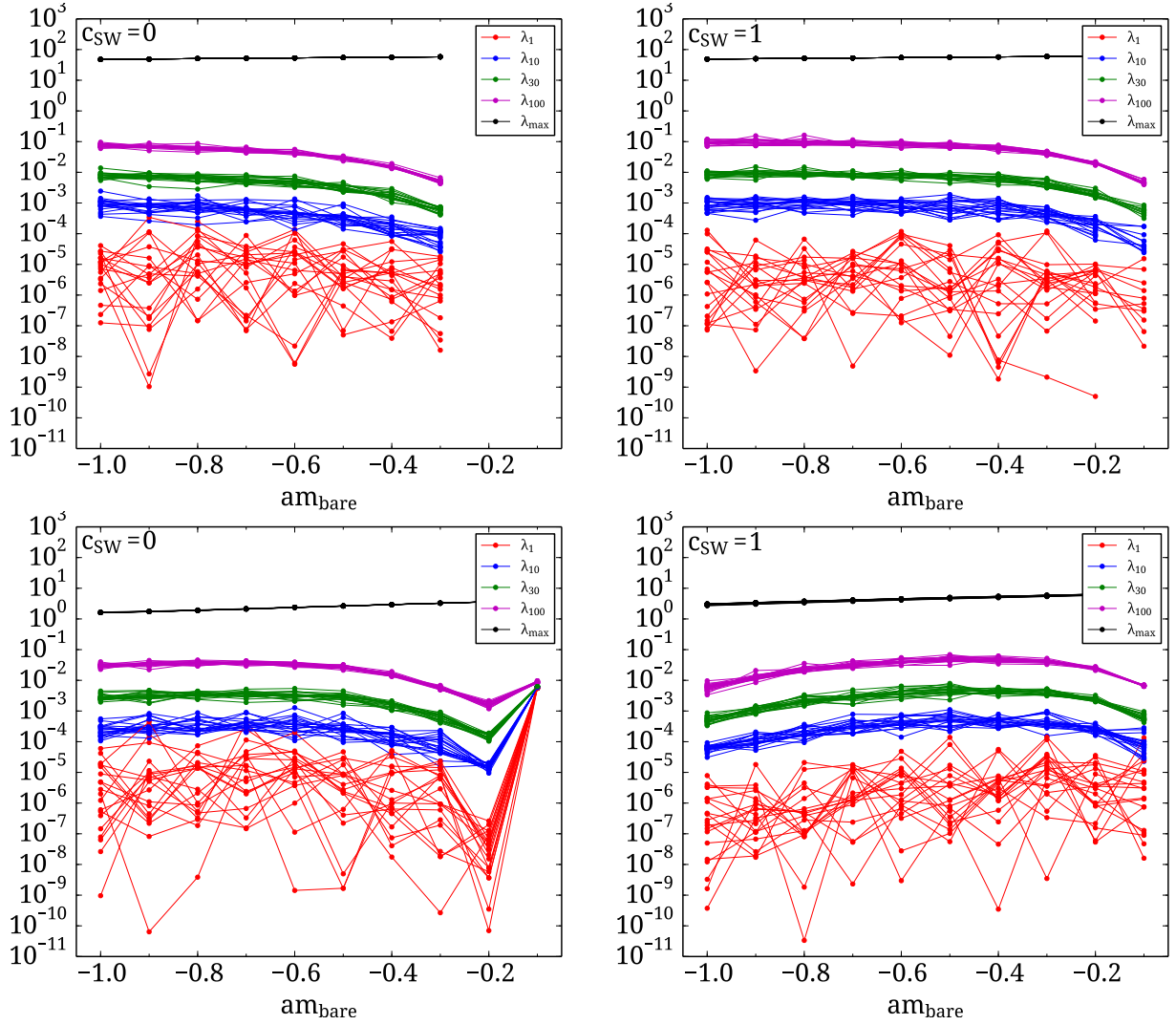


Figure 8: Spectrum of  $H^2_{-\rho} = D^\dagger_{-\rho} D_{-\rho}$  versus  $am_{\text{bare}} = -\rho$  on the QCDSF lattices, using the Wilson (top) and Brillouin (bottom) kernel with  $c_{\text{SW}} = 0$  (left) and  $c_{\text{SW}} = 1$  (right). We show the first, tenth, thirtieth and 100th eigenvalues ( $\lambda_{1,10,30,100}$ ), as well as the largest one ( $\lambda_{\text{max}}$ ).

## 7.2 Ginsparg-Wilson relation

We use the same (low-order) rational approximation to the sign function and the same pure gauge ensembles as in the previous subsection. We measure  $\|(D\gamma_5 + \gamma_5 D - D\gamma_5 D)\eta\|$ , which we will refer to as the “GW defect”, for a few dozen normalized Gaussian random vectors  $\eta$  on 40 configurations of each ensemble, see Fig. 10. The GW defect with the Wilson kernel is several orders of magnitude larger than with the Brillouin kernel. With the Wilson kernel the difference between the two smearing levels is barely visible, while with the Brillouin kernel increasing  $N_{\text{APE}}$  from 1 to 3 significantly reduces the GW defect. Moreover, in the Brillouin case pushing to the continuum (i.e. to smaller  $g_0^2$ ) has a beneficial effect, too, while no such effect is visible with the Wilson kernel. Evidently, in order to achieve a fixed level of GW violation, say  $\|(D\gamma_5 + \gamma_5 D - D\gamma_5 D)\eta\| < 10^{-12}$ , the order of the rational approximation needs to be increased much more drastically for the Wilson kernel than for the Brillouin kernel.

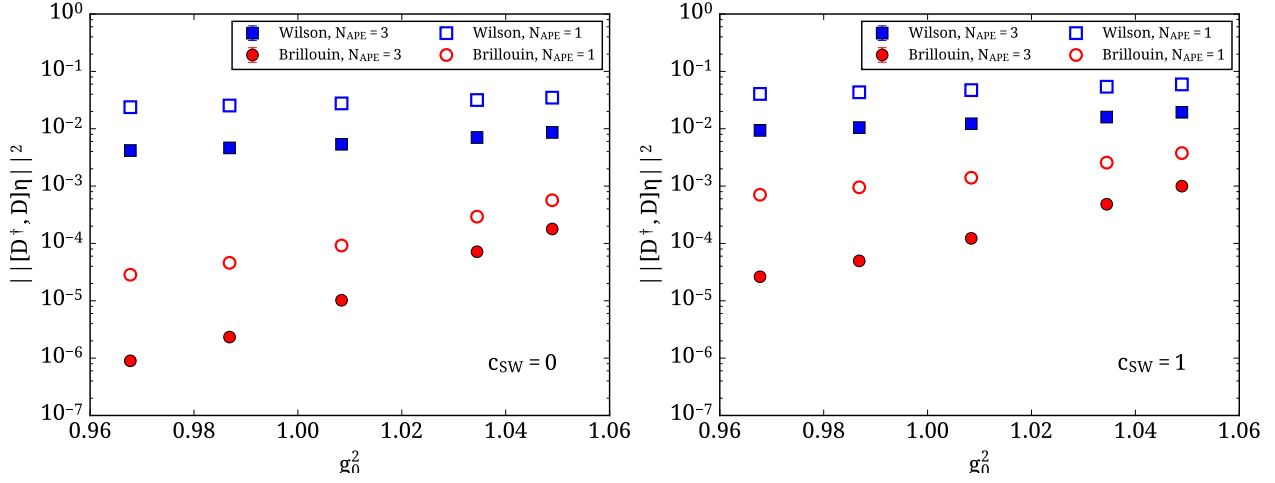


Figure 9: Normality of the iterate  $D_{-1}^{(1)} + 1$  with the Kenney-Laub mapping  $f_{1,1}$ , and  $D_{-1}^{(0)}$  the Wilson or Brillouin kernel with 1 or 3 APE smearings, and  $c_{\text{SW}} = 0$  (left) or  $c_{\text{SW}} = 1$  (right). Results on volume-matched ensembles of 40 quenched lattices each are plotted versus  $6/\beta$ .

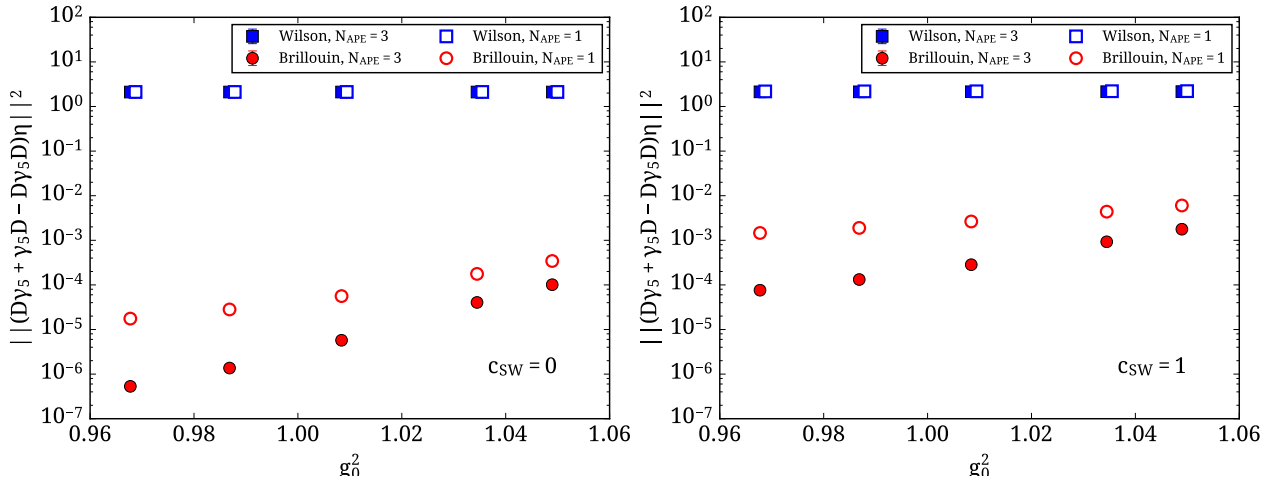


Figure 10: GW defect of the iterate  $D_{-1}^{(1)} + 1$  with the Kenney-Laub function  $f_{1,1}$ , and  $D_{-1}^{(0)}$  the Wilson or Brillouin kernel with 1 or 3 APE smearings, and  $c_{\text{SW}} = 0$  (left) or  $c_{\text{SW}} = 1$  (right). Results on volume-matched ensembles of 40 quenched lattices each are plotted versus  $6/\beta$ .

### 7.3 Exponential operator localization

The locality of the overlap action with the Wilson kernel was first studied in Ref. [46]. Ref. [47] demonstrated that a more extended (but still ultralocal) kernel can significantly improve the coordinate-space locality of the resulting overlap action. In Refs. [48–50] it was shown that even a slight modification through some link-smearing can lead to a considerable improvement. Therefore, one may hope that trading the Wilson kernel for the Brillouin kernel leads to a noticeable improvement of the locality of the overlap operator. Note that all of this holds up to some gauge coupling  $g_0^{\text{max}}$ , since Refs. [51, 52] pointed out that, once the gauge background becomes too rough, eigenmodes of the underlying shifted kernel  $D_{-\rho}^{\text{ke}}$  delocalize, and mix into a band, with the effect that the overlap operator may cease to be exponentially localized.

We measure the fall-off of  $|\zeta(x)|$ , with  $\zeta = D\eta$  and  $\eta$  a normalized Gaussian random vector

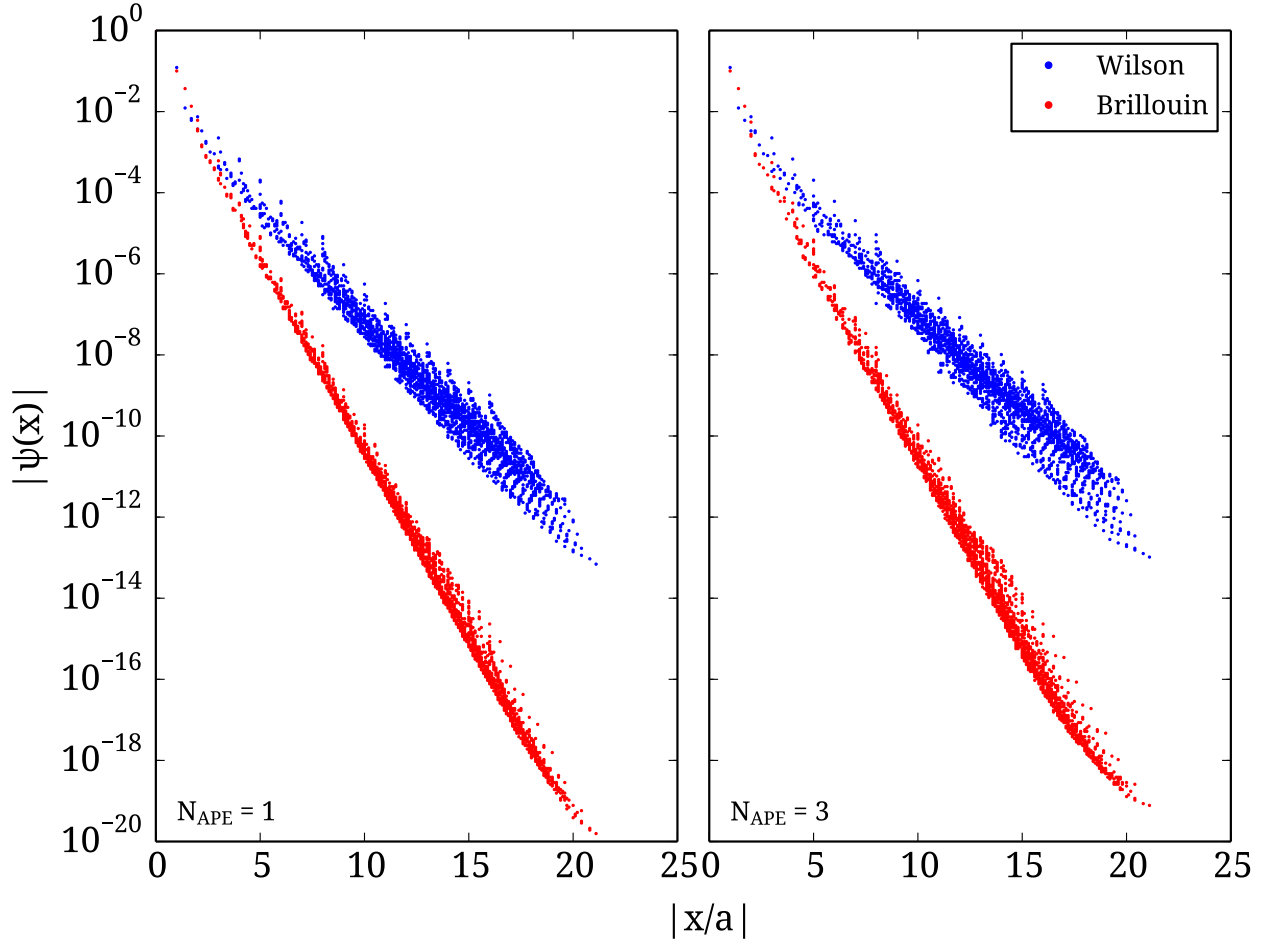


Figure 11: *Localization of the approximate overlap action defined through  $f_{1,1}$  with Wilson or Brillouin kernel, 1 or 3 APE smearings and  $c_{\text{sw}} = 1$  prior to averaging over various directions with a fixed value of  $|x|$ . Quenched lattices at  $\beta = 5.95$  and  $L/a = 16, T/a = 32$  are used.*

(in spinor/color space) with support at the single site 0, for about a dozen  $\eta$  per config on 20 quenched configs per  $\beta$ . Fig. 11 shows the result at the lattice spacing  $a \simeq 0.1$  fm as a function of the Euclidean distance  $|x|$ . The norm falls off exponentially with distance, but there are signs of rotational symmetry breaking – different directions with a common  $|x|$  do not lie on top of each other. Clearly, this rotational symmetry breaking is more pronounced with the Wilson kernel, and a higher smearing level does not help. The other marked difference is that the fall-off rate with the Brillouin kernel is better than with the Wilson kernel. Fig. 12 shows the locality after averaging over directions  $x$  with the same  $|x|$ . The fall-off pattern looks even more exponential than previously, and the better localization (roughly by a factor 2 at fixed  $\rho = 1$ ) of the Brillouin version is found to be virtually independent of the lattice spacing. Note that the first row of Fig. 12 demonstrates that the combination of  $c_{\text{sw}} = 1$  and some link smearing ensures that either overlap action is exponentially localized on lattices with  $a \simeq 0.16$  fm.

How this localization, i.e. the “effective mass”  $\delta$  in  $|\psi(x)| \propto \exp(-\delta|x|)$ , varies as a function of  $\rho$  is shown in Fig. 13. With an unsmeared and unimproved Wilson kernel frequently a value  $\rho \simeq 1.4$  is chosen to optimize locality on coarse lattices [46]. This, however, creates a clash with the free-field behavior where optimum locality is reached for  $\rho \simeq 0.6$  [50]. Our figure shows that

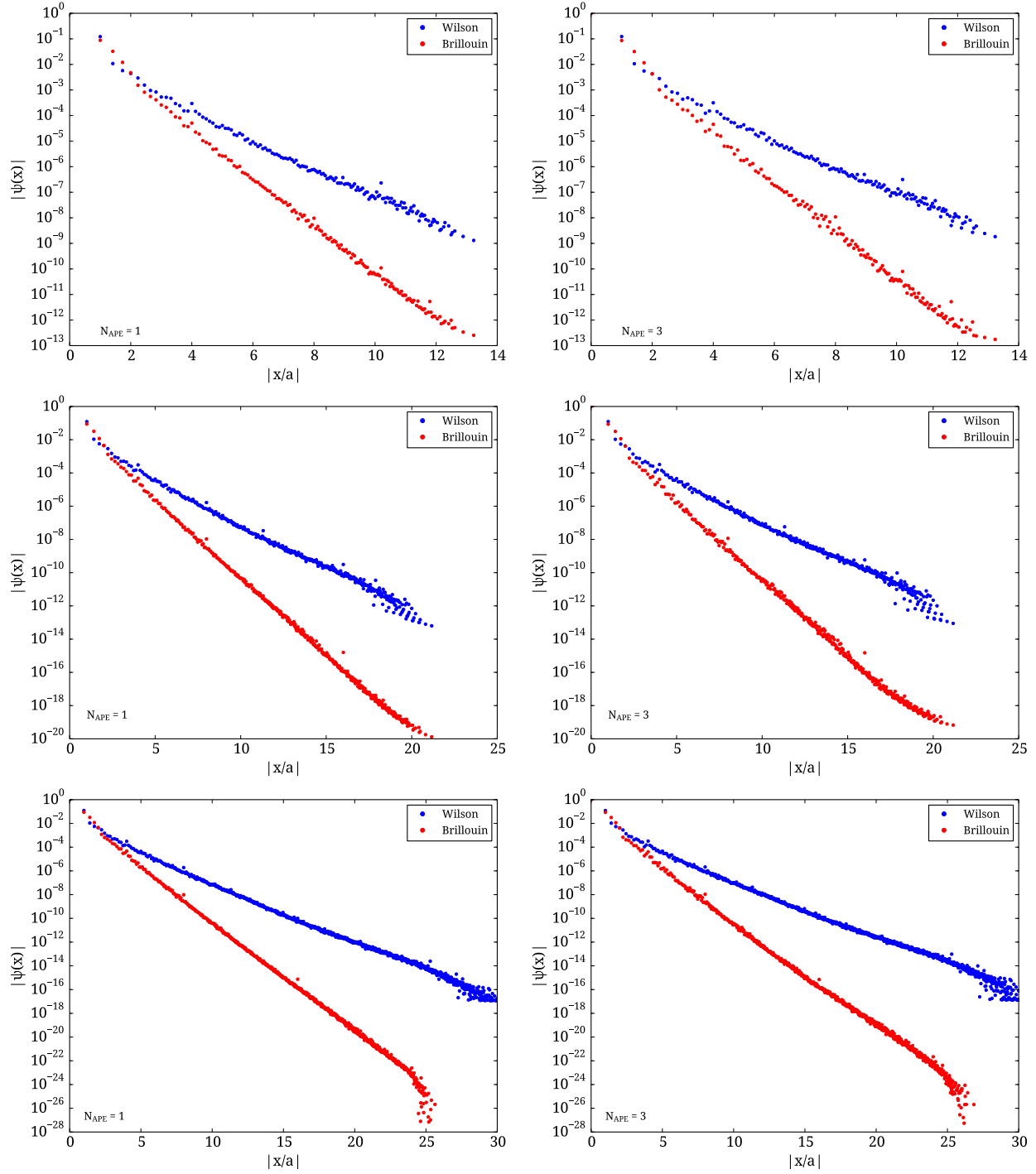


Figure 12: Same as Fig. 11 but after averaging over various directions with a common  $|x|$ . The panel rows feature  $\beta = 5.72, 5.95, 6.20$  and  $L/a = 10, 16, 24$  with  $T = 2L$  (from top to bottom).

even for the Wilson kernel this clash is resolved by some link smearing and putting  $c_{\text{SW}} = 1$ ; then the optimum is assumed at  $\rho \simeq 0.6$ . Similarly, the Brillouin kernel with link smearing and  $c_{\text{SW}} = 1$  has an optimum locality which (for accessible lattice spacings) is at  $\rho \simeq 1$ , but also the “magic” value  $\rho \simeq 0.634$  of Sec. 2 fares quite well.

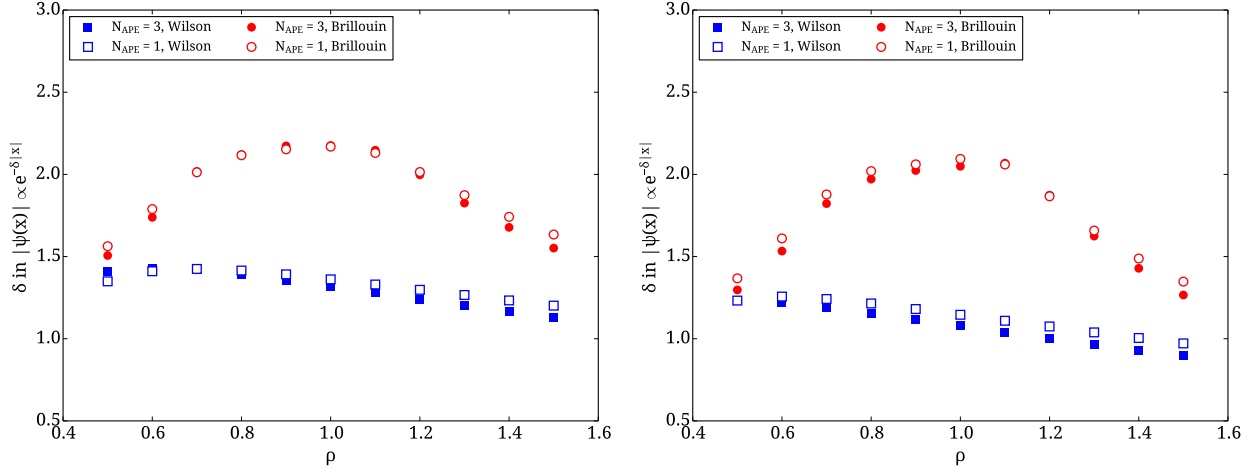


Figure 13: Inverse localization  $\delta$  for the  $f_{1,1}$  overlap actions with Wilson and Brillouin kernel, using 1-3 APE steps and  $c_{\text{SW}} = 1$ , as a function of  $\rho$ . We use the coarsest ( $a^{-1} = 1.236$  GeV,  $10^3 \times 20$  grid, left) and finest ( $a^{-1} = 2.964$  GeV,  $24^3 \times 48$  grid, right) lattices of Ref. [3].

In conclusion we find that the Brillouin kernel diminishes the anisotropy effects and results in an overlap operator that falls off significantly faster than the one with the Wilson kernel. This may turn out to be relevant for QCD studies of bulk thermodynamic properties [53].

## 7.4 Exploration of inversion cost and residual mass

To further assess the suitability of the Brillouin operator as a kernel to the overlap procedure we conduct a pilot study of overlap inversions with a given source vector, as is typical in spectroscopy calculations. The overall setup is standard [30–33]; we use a BiCGstab (“outer”) solver and a Kenney-Laub  $f_{n,n}$  (“inner”) approximation to the matrix sign function.

We use a freely available  $N_f = 2$  ensemble by QCDSF, with geometry  $40^3 \times 64$ , sea pion mass  $M_\pi \simeq 280$  MeV and lattice spacing  $a \simeq 0.0728(05)(19)$  fm deduced from  $a^{-1} = 2.71(2)(7)$  GeV at  $\beta = 5.29$  [54, 55]. Given the results in the previous subsections, we focus on the overlap operator with 3 smearings (at  $\alpha_{\text{APE}} = 0.72$  each) and no ( $c_{\text{SW}} = 0$ ) or tree-level ( $c_{\text{SW}} = 1$ ) clover improvement in the kernel. The shift parameter is pinned to the canonical value  $\rho = 1$  to avoid any tuning overhead; using the “magic” value  $\rho \simeq 0.634$  is not expected to bring any significant change. The lattices are sufficiently long in Euclidean time such that we can identify clear effective mass plateaus for all bare quark masses studied. A selection of such plateaus is shown in Fig. 14. With fixed statistics, the statistical errors grow at small quark masses, but it is always evident that excited states contributions disappear at large  $t/a$ .

We measure  $aM_\pi$  and monitor the number  $N_{\text{iter}}$  of outer iterations (i.e. of BiCGstab) for a selection of  $f_{n,n}$  overlap masses  $am$ . Results with  $c_{\text{SW}} = 0$  and  $c_{\text{SW}} = 1$  in the Brillouin kernel are shown in Fig. 15 and presented in Tab. 4. Since these are approximate overlap fermions, the additive mass shift is non-zero. With the  $c_{\text{SW}} = 0$  kernel, using  $f_{1,1}$  brings more than an order of magnitude reduction, compared to the bare Brillouin action, and using  $f_{4,4}$  makes it consistent with zero within our statistical precision. On the other hand, with the  $c_{\text{SW}} = 1$  kernel, the bare Brillouin action has a comparatively small mass shift, using  $f_{1,1}$  makes it consistent with zero within  $\sim 2\sigma$ , while using  $f_{4,4}$  makes it consistent with zero within  $\sim 1\sigma$ .

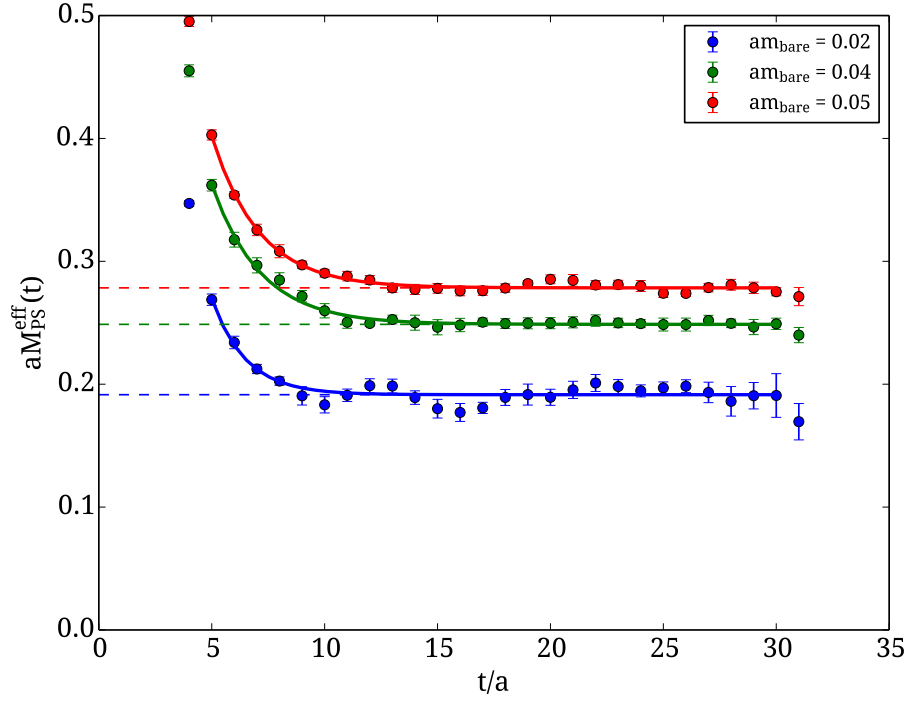


Figure 14: *Effective masses of pions built from two degenerate  $f_{1,1}$  overlap fermions based on the Brillouin kernel with 3 APE steps and  $c_{\text{SW}} = 1$ . The bare quark masses  $am = 0.02, 0.04, 0.05$  correspond to pion masses of 520 MeV, 670 MeV, 750 MeV, respectively.*

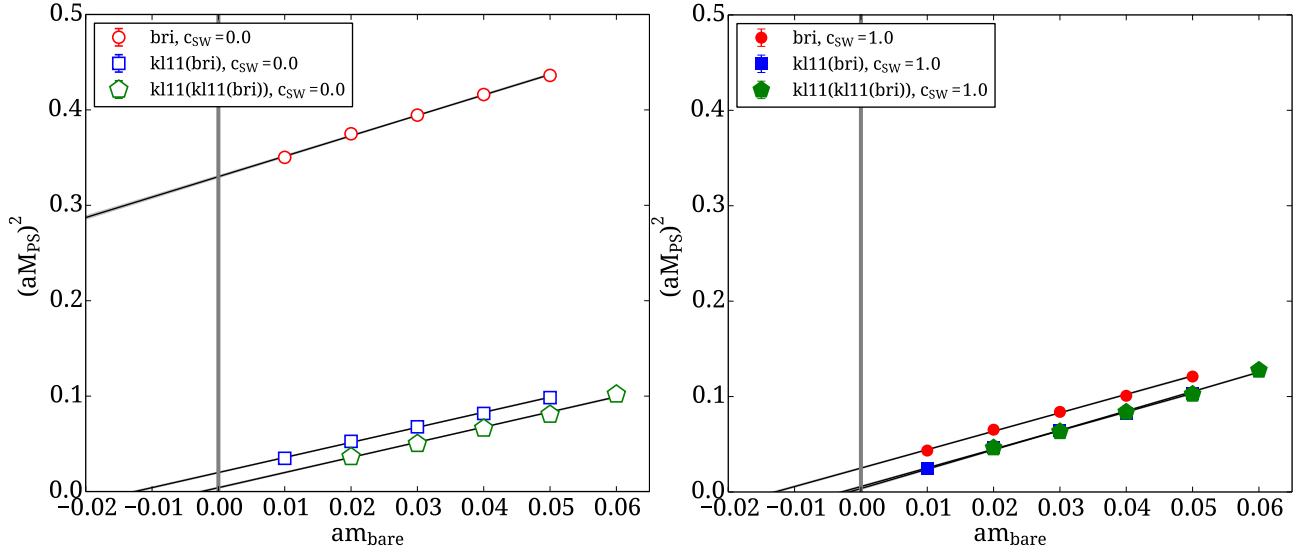


Figure 15:  *$(aM_\pi)^2$  versus  $am$  for  $f_{1,1}$  and  $f_{4,4} = f_{1,1}^{(2)}$  overlap fermions based on the Brillouin kernel with  $c_{\text{SW}} = 0$  (left) and  $c_{\text{SW}} = 1$  (right). Throughout 3 levels of APE smearing are used.*

In the literature on approximate overlap fermions it is common practice to determine a “residual mass”, i.e. an effective fermion mass evaluated at  $am = 0$ . In case of domain-wall fermions typically a version is used which explicitly refers to 5 dimensions [56]. Since this is not an option for us, we choose the PCAC quark mass, employing the definition which is standard for Wilson fermions [57]. The result is shown in Fig. 16, where the intercepts of the gray bands

$am$	$aM_\pi$	$M_\pi$ [MeV]	$n_{\text{iter}}$	time [sec]	$n_{\text{nodes}}$
0.004	0.162(1)	430	3433.9	865.2	320
0.010	0.192(1)	520	2314.3	592.9	320
0.020	0.227(2)	620	1311.9	320.9	320
0.035	0.271(2)	730	878.7	215.6	320
0.050	0.313(1)	850	652.8	178.6	320
0.01	0.155(2)	420	13175.2	3469.8	320
0.02	0.191(2)	520	6996.2	3858.8	200
0.03	0.222(2)	600	3828.0	1570.6	160
0.04	0.249(2)	670	2524.0	1656.6	160
0.05	0.278(2)	750	2166.9	1463.3	160
0.06	0.307(2)	830	1549.2	661.7	160

Table 4: Overview of the pion mass  $aM_\pi$  as a function of the quark mass  $am$  for the  $f_{1,1}$  Brillouin overlap fermion with  $c_{\text{SW}} = 0$  (top part) and  $c_{\text{SW}} = 1$  (bottom part) in the kernel. In addition, we give the average number of BiCGstab iterations and the average time per right-hand-side on  $n_{\text{nodes}}$  of the commodity cluster JUROPA.

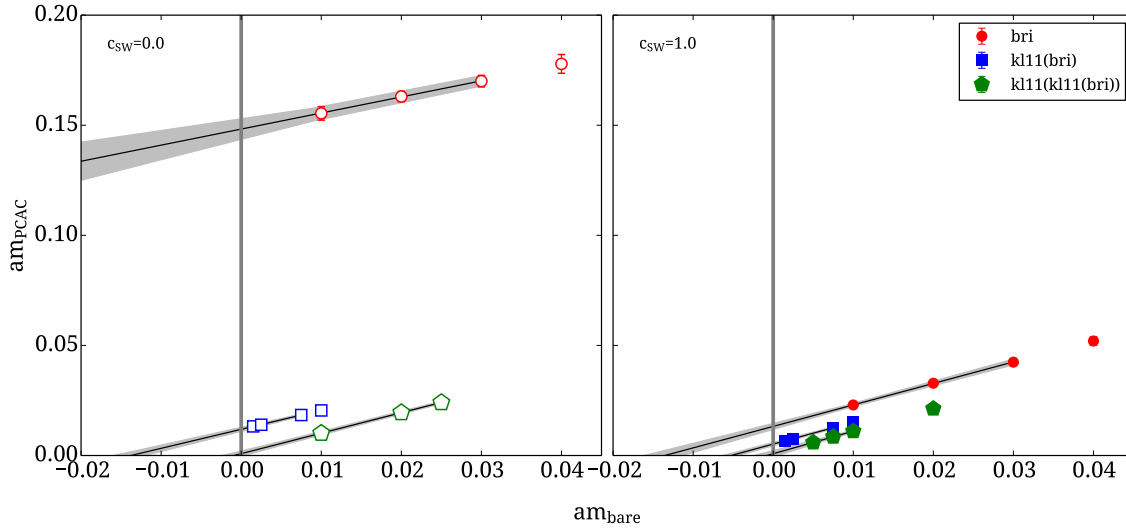


Figure 16:  $PCAC$  versus bare quark mass for the Brillouin operator, its first and second  $KL11$  iterates, with  $c_{\text{SW}} = 0$  (left) and  $c_{\text{SW}} = 1$  (right) on the  $40^3 \times 64$  ensemble by QCDSF.

with the  $y$ -axis represent our residual quark masses. The overall picture looks similar to the one in Fig. 15, except that this time also the  $f_{1,1}$  overlap action with  $c_{\text{SW}} = 1$  kernel shows a residual mass which is clearly non-zero. But with  $f_{4,4}$  the intercept is zero within errors, regardless of the value of  $c_{\text{SW}}$  in the kernel.

## 8 Cascaded preconditioning

In our opinion a dedicated research effort is needed to identify good preconditioners for the repeated inversions of the type  $(D_{-\rho}^{\text{ke} \dagger} D_{-\rho}^{\text{ke}} + \sigma)x = b$  which occur in the “inner” (CG-type) solver employed in the evaluation of partial fraction representation (17). For  $n = 1$  the shift is



$\sigma = 1/3$ , for  $n = 4$  the smallest shift is  $\sigma \simeq 0.0311$ ; see App. B for details.

There is, however, a simple preconditioning strategy for the “outer” (BiCGstab-type) solver which is particularly convenient with a Brillouin kernel. It builds on the relative proximity of  $D^{\text{ke}} - 1$  on one hand and various Kenney-Laub iterates of this combination on the other hand, see Fig. 5. Suppose we wish to invert the operator defined by  $f_{1,1}^{(3)} = f_{13,13}$ . One may then use

$$D_{-1}^{\text{ke}} + (1 + \frac{\tilde{m}}{\rho}) = D^{\text{ke}} + \frac{\tilde{m}}{\rho} \quad (34)$$

as a preconditioner to the operator (with  $A_{-\rho}^{(k)} = D_{-\rho}^{(k)\dagger} D_{-\rho}^{(k)}$  and  $D_{-\rho}^{(0)} = D_{-\rho}^{\text{ke}}$ )

$$D_{-\rho}^{(1)} + (1 + \frac{\tilde{m}}{\rho}) \equiv D_{-\rho}^{\text{ke}} \left( \frac{1}{3} + \frac{8/9}{A_{-\rho}^{\text{ke}} + 1/3} \right) + (1 + \frac{\tilde{m}}{\rho}) \quad (35)$$

which in turn is used as a preconditioner to the operator defined by  $f_{1,1}^{(2)} = f_{4,4}$

$$D_{-\rho}^{(2)} + (1 + \frac{\tilde{m}}{\rho}) \equiv D_{-\rho}^{(1)} \left( \frac{1}{3} + \frac{8/9}{A_{-\rho}^{(1)} + 1/3} \right) + (1 + \frac{\tilde{m}}{\rho}) \quad (36)$$

$$= D_{-\rho}^{\text{ke}} \left( \frac{1}{9} + \frac{0.229...}{A_{-\rho}^{\text{ke}} + 0.0311...} + \frac{0.296...}{A_{-\rho}^{\text{ke}} + 0.333...} + \frac{0.538...}{A_{-\rho}^{\text{ke}} + 1.42...} + \frac{1.90...}{A_{-\rho}^{\text{ke}} + 7.55...} \right) + (1 + \frac{\tilde{m}}{\rho})$$

where the representation in terms of  $A^{\text{ke}}$  uses the coefficients in the partial fraction expansion given in App. B. The latter operator is used as a preconditioner to solve the equation

$$\left[ D_{-\rho}^{(3)} + (1 + \frac{\tilde{m}}{\rho}) \right] x = \frac{\tilde{b}}{\rho} \quad \text{where} \quad D_{-\rho}^{(3)} \equiv D_{-\rho}^{\text{ke}} \left( \frac{1}{27} + \frac{0.0743...}{A_{-\rho}^{\text{ke}} + 0.00339...} + \dots + \frac{5.50...}{A_{-\rho}^{\text{ke}} + 73.2...} \right)$$

for  $x$ , where  $\tilde{m}, \tilde{b}$  are given in (30), and the full set of coefficients is again found in App. B.

In Fig. 17 and Tab. 5 we show that this concept works over one and two steps (i.e. for  $D^{(1)}$  and  $D^{(2)}$  implementing one and two iterations of  $f_{1,1}$ , respectively). It seems that significant savings can be achieved even if none of the preconditioner masses is tuned (we use the same bare mass in all operators, even in the Brillouin action). Without preconditioning the  $f_{4,4}$ -based approximant  $D^{(2)}$  to the overlap action is significantly more expensive than the  $f_{1,1}$ -based operator  $D^{(1)}$ , but with cascaded preconditioning the extra cost of the better approximation becomes more tolerable.

To the best of our knowledge, preconditioning of an overlap operator by its kernel was first tried in Ref. [58], and a more elaborate version of this (with tuning of the preconditioner mass) was presented in Ref. [59]. Both of these references use the Wilson action as a preconditioner to the Wilson overlap. Our Fig. 17 demonstrates that the same concept works for Brillouin fermions, too, but obviously there is much room for optimization, still, on our side.

## 9 Outlook on dynamical Brillouin overlap simulations

We close with a brief outlook on how our findings fit into the perspective of carrying out dynamical overlap simulations based on the Brillouin kernel and the Hybrid Monte Carlo (HMC) algorithm [60]. For dynamical overlap simulations with a Wilson kernel see e.g. Refs. [61–66].

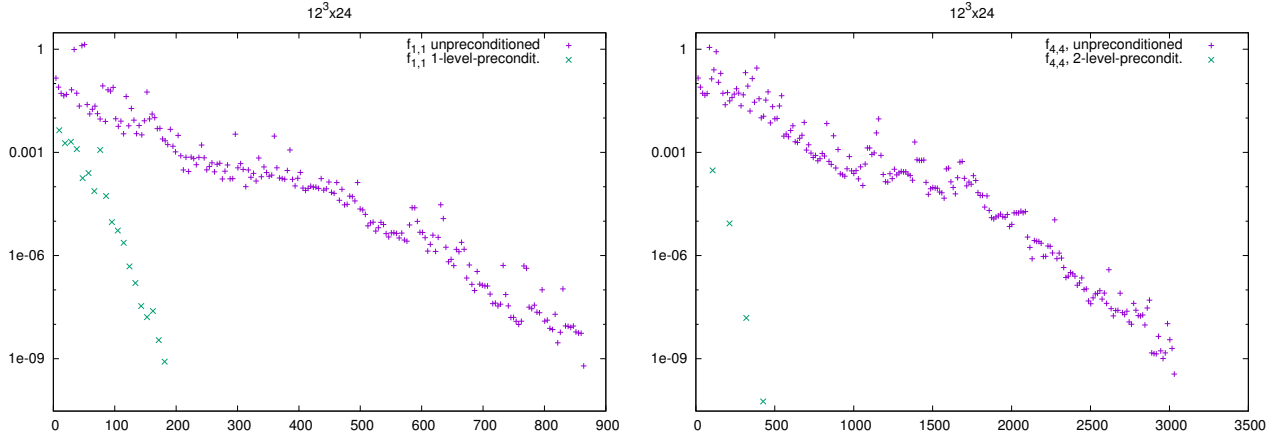


Figure 17: Convergence of the BiCGstab solver for the Brillouin overlap with  $\rho = 1$ ,  $c_{\text{sw}} = 0$ ,  $n_{\text{smear}} = 7$  and  $am = 0.03$  on a pure glue  $12^3 \times 24$  lattice at  $\beta = 6.0$ . We use (35) as outer operator and no or one level of preconditioning (left) or (36) as outer operator and no or two levels of preconditioning (right). We handle 12 right-hand sides simultaneously (cf. App. C), and the worst relative residual is shown. The x-axis denotes wall-clock time in seconds.

Volume	$f_{n,n}$	precond.	$n_{\text{iter}}$	time [sec]
$12^3 \times 24$	1,1	none	204	864
	1,1	1-level	19	181
	4,4	none	212	3031
	4,4	2-level	4	426
$16^3 \times 32$	1,1	none	249	4265
	1,1	1-level	22	732
	4,4	none	285	16706
	4,4	2-level	10	3397
$24^3 \times 48$	1,1	none	264	27490
	1,1	1-level	25	5040
	4,4	none	292	132975
	4,4	2-level	9	19835

Table 5: Timings of the BiCGstab solver for the Brillouin overlap action with  $\rho = 1$ ,  $c_{\text{sw}} = 0$ ,  $n_{\text{smear}} = 7$  and  $am = 0.03$  on pure glue lattices at  $\beta = 6.0$ . We use no or 1-level preconditioning for the operator defined through  $f_{1,1}$ , and no or 2-level preconditioning with  $f_{4,4}$ . The exit criterion was set to  $\|r\|/\|b\| \leq 10^{-9}$ . The number of iterations is for the outer solver, the timings are per right-hand-side on one (4-core) node. Details for  $12^3 \times 24$  are shown in Fig. 17.

The HMC algorithm is governed by the molecular dynamics time evolution which, in turn, builds on the HMC force. Let  $D_m$  be an undoubled fermion operator, which implicitly depends on the “thin” gauge links  $U$ . The pseudo-fermion action for  $N_f$  degenerate fermions is

$$S_{\text{pf}} = \langle \phi | A_m^{-N_f/2} | \phi \rangle = \int \phi^\dagger(x) A_m^{-N_f/2}(x, y) \phi(y) d^4x d^4y, \quad A_m = D_m^\dagger D_m > 0 \quad (37)$$

where  $\phi$  denotes a boson field with the spinor/color components of a standard Dirac fermion. The HMC force is defined as minus the derivative of  $S_{\text{pf}}$  with respect to the thin gauge links [60].

For  $N_f = 1$  one exploits the fact that the  $p$ -th order diagonal rational approximation of  $x^{-1/2}$  over the relevant spectral range admits a partial fraction representation [60]

$$x^{-1/2} \simeq \alpha_0 + \sum_{k=1}^p \frac{\alpha_k}{x + \beta_k} \quad (38)$$

with  $\alpha_k > 0$  for  $1 \leq k \leq p$  and  $0 < \beta_1 < \dots < \beta_p$ . For  $N_f = 1$  the pseudo-fermion force is thus

$$F_{\text{pf}} = -S'_{\text{pf}} = \sum_{k=1}^p \alpha_k \langle \phi | (A_m + \beta_k)^{-1} A'_m (A_m + \beta_k)^{-1} | \phi \rangle \quad (39)$$

where the prime denotes the derivative with respect to the gauge field element  $A_\mu^a(x + \hat{\mu}/2)$ , defined as a Gell-Mann component of  $\log(U_\mu(x))$ . For  $N_f \in 2\mathbf{N}$  no rational representation is needed (though it still might be favorable for efficiency reasons [60]), and one gets away with products of powers of  $A_m^{-1}$  and factor  $A'_m$ . The bottom line is that in both cases – even or odd  $N_f$  – the derivative  $A'_m$  with respect to the thin gauge field is required to work out the pseudo-fermion contribution to the molecular dynamics force.

This  $A'_m$  is straightforward to write down in case  $D_m$  is an ultralocal operator. On the other hand,  $A'_m$  is more involved with an overlap action. The attentive reader will have noticed that we advocate using a fixed-order rational approximation to the matrix sign function (see Sec. 3), regardless of the spectral properties of  $H_{-\rho}^{\text{ke}} = \gamma_5 D_{-\rho}^{\text{ke}}$  on the current gauge background  $U$  (which is common practice with the domain-wall setup [56, 64]). This is convenient, since the order of the rational approximation to the sign function does not increase as the belly-mode of the underlying hermitean kernel changes sign (in an attempt of the HMC algorithm to increase/decrease the global topological charge by one unit). However, with a fixed-order approximation to the matrix sign function in the definition of  $D_m$ , it is still fairly easy to work out the inner derivative  $A'_m$ . From the definition (11) we obtain

$$g'_\ell(x) = 2\ell \frac{(1+x)^\ell(1-x)^{\ell-1} + (1+x)^{\ell-1}(1-x)^\ell}{[(1+x)^\ell + (1-x)^\ell]^2} = 4\ell \frac{(1-x^2)^{\ell-1}}{[(1+x)^\ell + (1-x)^\ell]^2} \quad (40)$$

which implies  $g'_\ell(-x) = g'_\ell(x)$  for any  $\ell \in \mathbf{N}$ , along with  $g'_\ell(x) \geq 0$  for odd  $\ell$ . In other words, the diagonal Kenney-Laub approximations  $f_{n,n}$  to the matrix sign function grow monotonically on  $]0, \infty[$  and strictly monotonically on the intervals  $]0, 1[$  and  $]1, \infty[$ . This is in marked distinction to the situation with optimal rational approximations which show, within the accessible interval  $[\lambda_{\min}, \lambda_{\max}]$ , many “wiggles”, i.e. small-scale oscillations, in particular close to the endpoints. It remains to be seen whether this peculiar property of the diagonal Kenney-Laub approximants has any impact on e.g. the topological tunneling rate at a fixed lattice spacing  $a$ .

In the end the computation of various inverses of  $D_m$  and  $A_m$  is required, and for this purpose the cascaded preconditioning technique as discussed in Sec. 8 will prove useful. However, in the HMC algorithm there is still more room for optimization. In principle, the pseudo-fermion force can be calculated with any fermion operator, as long as one includes the difference between the used and the desired pseudo-fermion action in the final accept-reject step. For instance, in Ref. [67] it was proposed to use the staggered action as a driving engine in the molecular dynamics evolution for HMC simulations of the Wilson overlap action. Of course, the further the two actions involved and the longer the trajectory length the lower the acceptance rate. However, in two space-time dimensions such games have been played successfully [68], and we feel optimistic that the relative proximity of the Brillouin kernel to the Brillouin overlap action will allow for significant savings in four space-time dimensions, too.

## 10 Summary

We summarize the main results of our investigation as follows:

1. The free-field dispersion relation of the Brillouin overlap action with generic  $\rho$  deviates from the continuum relation  $(aE)^2 - (a\mathbf{p})^2 = (am)^2$  through a term proportional to  $(am)^4$ . With the “magic” value  $\rho = (3 - \sqrt{3})/2$  the leading discretization error is lifted to order  $(am)^6$ . We hope that this feature proves useful to compute properties of systems with charm quarks, perhaps in a further perspective even with bottom quarks.
2. We advocate using any of the diagonal Kenney-Laub approximants  $g_{2n+1} = f_{n,n}$  to the matrix sign function at fixed order  $n$  in the definition of the overlap action. This is close in spirit to what is done in the domain-wall setup, except that then a five-dimensional framework is used, and the effective four-dimensional operator corresponds to an element  $f_{n-1,n} = g_{2n}$  of the Kenney-Laub family of matrix iterations. In both cases the partial-fraction expansion involved in the definition of the overlap action does not depend on the gauge field details and can be constructed beforehand. This, in turn, makes it easier to demonstrate good strong-scaling properties on massively parallel architectures.
3. We advocate defining the massive overlap action  $D_m$  through (29) rather than (27). Effectively this has been done in the past through the “extra prescription” of dressing external currents or densities with a “chiral symmetry ensuring factor”  $(1 - aD/[2\rho])$ , where  $D$  is the zero-mass overlap operator. Hence our proposal adds to the “piece of mind”, since the danger of missing an important ingredient in a later step of the calculation is bypassed.
4. We checked that the eigenvalue spectra of both the non-hermitean and the hermitean kernel operator show promising features for the extraction of the unique unitary part of  $D_{-\rho}^{\text{ke}}$  through an application of  $f_{n,n} = g_{2n+1}$  with sufficiently high  $n$ . With the Wilson and the Brillouin kernel link smearing helps to open the “eye” in the eigenvalue spectrum and hence to reduce the appropriate value of  $n$ . On the other hand a clover term in the kernel was neither found to bring a clear advantage nor a clear disadvantage.
5. In terms of physics properties the second important advantage of the Brillouin kernel over the Wilson kernel is the improved locality of the resulting overlap operator. Whenever locality is important and the CPU cost scale with a high power of  $a$  (as is the case in the study of bulk thermodynamics properties), a complete study with a valid continuum limit may be cheaper with the Brillouin overlap than with the Wilson overlap action.
6. The proximity of the non-chiral Brillouin kernel to the (exact) Brillouin overlap operator (and the diagonal Kenney-Laub approximants to the latter) makes the Brillouin overlap action particularly susceptible to cascaded preconditioning strategies. Even without any tuning effort, significant savings were found for a BiCGstab solver of the Brillouin overlap operator with one and two levels of preconditioning.

### Acknowledgments:

The authors gratefully acknowledge the computing time granted by Forschungszentrum Jülich GmbH and provided on the supercomputer JUROPA at Jülich Supercomputing Centre (JSC). This work was in parts supported by DFG through the program SFB-TR-55.

## A Details of quark-level dispersion relations

In this appendix we shall use the boson momentum  $\hat{p}$  as well as the std-fermion and iso-fermion momenta  $\bar{p}$  and  $\tilde{p}$  which are defined such that  $\nabla_\mu^{\text{std}} = i\bar{p}_\mu$  and  $\nabla_\mu^{\text{iso}} = i\tilde{p}_\mu$  [3], viz.

$$\hat{p}_\mu = \frac{2}{a} \sin\left(\frac{ap_\mu}{2}\right), \quad \bar{p}_\mu = \frac{1}{a} \sin(ap_\mu), \quad \tilde{p}_\mu = \frac{1}{27a} \sin(ap_\mu) \prod_{\nu \neq \mu} \{\cos(ap_\nu) + 2\}. \quad (41)$$

Furthermore, we need the momentum space representations of the Laplacians [3]

$$\Delta^{\text{std}} = -\frac{4}{a^2} \sum_\mu \sin^2\left(\frac{ap_\mu}{2}\right) = \frac{2}{a^2} \sum_\mu \cos(ap_\mu) - \frac{8}{a^2} = -\sum_\mu \hat{p}_\mu^2 \equiv -\hat{p}^2 \quad (42)$$

$$\Delta^{\text{bri}} = \frac{4}{a^2} \prod_\mu \cos^2\left(\frac{ap_\mu}{2}\right) - \frac{4}{a^2} = \frac{1}{4a^2} \prod_\mu \{\cos(ap_\mu) + 1\} - \frac{4}{a^2} \equiv -\check{p}^2 \quad (43)$$

and we caution that (unlike  $\hat{p}^2$ ) the quantity  $\check{p}^2$  is not a sum of squares.

### A.1 Dispersion relation for Wilson operator

The Green's function of the Wilson operator at mass  $am$  and  $r = 1$  follows as

$$\begin{aligned} D_{W,m} &= \nabla_\mu^{\text{std}} \gamma_\mu - \frac{a}{2} \Delta^{\text{std}} + m = i\bar{p}_\mu \gamma_\mu + \frac{a}{2} \hat{p}^2 + m \\ G_{W,m} &= \frac{-i\bar{p}_\sigma \gamma_\sigma + \frac{a}{2} \hat{p}^2 + m}{(i\bar{p}_\mu \gamma_\mu + \frac{a}{2} \hat{p}^2 + m)(-i\bar{p}_\nu \gamma_\nu + \frac{a}{2} \hat{p}^2 + m)} = \frac{-i\bar{p}_\sigma \gamma_\sigma + \frac{a}{2} \hat{p}^2 + m}{\hat{p}^2 + (\frac{a}{2} \hat{p}^2 + m)^2}. \end{aligned} \quad (44)$$

Searching for a zero of the denominator with  $\frac{a}{2} \hat{p}^2 = -\frac{1}{a} \sum_\mu \cos(ap_\mu) + \frac{4}{a}$  and  $p_4 \rightarrow iE$  yields

$$\sinh^2(aE) - \sum_i \sin^2(ap_i) = \cosh^2(aE) + 2 \cosh(aE) [\sum_i \cos(ap_i) - 4 - am] + [\dots]^2 \quad (45)$$

and upon using  $\cosh^2 - \sinh^2 = 1$  this turns into a *linear equation* in  $\cosh(aE)$  which gives (4). For the sake of a check we note that (45) simplifies to  $1 + [1 + am]^2 = \cosh \cdot [2 + 2am]$  at  $\mathbf{ap} = \mathbf{0}$ .

### A.2 Dispersion relation for Brillouin operator

The Green's function of the Brillouin operator at mass  $am$  and  $r = 1$  follows as

$$\begin{aligned} D_{B,m} &= \nabla_\mu^{\text{iso}} \gamma_\mu - \frac{a}{2} \Delta^{\text{bri}} + m = i\tilde{p}_\mu \gamma_\mu + \frac{a}{2} \check{p}^2 + m \\ G_{B,m} &= \frac{-i\tilde{p}_\sigma \gamma_\sigma + \frac{a}{2} \check{p}^2 + m}{(i\tilde{p}_\mu \gamma_\mu + \frac{a}{2} \check{p}^2 + m)(-i\tilde{p}_\nu \gamma_\nu + \frac{a}{2} \check{p}^2 + m)} = \frac{-i\tilde{p}_\sigma \gamma_\sigma + \frac{a}{2} \check{p}^2 + m}{\check{p}^2 + (\frac{a}{2} \check{p}^2 + m)^2}. \end{aligned} \quad (46)$$

Searching for a zero of the denominator with  $\frac{a}{2} \check{p}^2 = \frac{2}{a} - \frac{1}{8a} \prod_\mu \{c_\mu + 1\}$  yields

$$a^2 \sum_\mu \tilde{p}_\mu^2 + \frac{1}{64} \prod_\mu \{c_\mu + 1\}^2 - \frac{1}{4} \prod_\mu \{c_\mu + 1\} [2 + am] + [2 + am]^2 = 0$$

with  $a^2 \tilde{p}^2 = \frac{1}{729} \sum_\mu s_\mu^2 \prod_{\nu \neq \mu} \{c_\nu + 2\}^2$  and  $c_\mu \equiv \cos(ap_\mu)$ ,  $s_\mu \equiv \sin(ap_\mu)$ . Next,  $p_4 \rightarrow iE$  leads to

$$\begin{aligned} &\frac{1}{729} \sum_i s_i^2 \prod_{j \neq i} \{c_j + 2\}^2 \{\cosh^2 + 4 \cosh + 4\} - \frac{1}{729} \prod_i \{c_i + 2\}^2 \sinh^2 \\ &+ \frac{1}{64} \prod_i \{c_i + 1\}^2 \{\cosh^2 + 2 \cosh + 1\} - \frac{1}{4} \prod_i \{c_i + 1\} \{\cosh + 1\} [2 + am] + [2 + am]^2 = 0 \end{aligned} \quad (47)$$

and upon using  $\cosh^2 - \sinh^2 = 1$  this turns into a *quadratic equation* in  $\cosh(aE)$ . We note that it simplifies to  $\{1 - \cosh^2\} + \{\cosh^2 + 2\cosh + 1\} - \{\cosh + 1\}[4 + 2am] + [2 + am]^2 = 0$  at  $a\mathbf{p} = \mathbf{0}$ , which agrees with the respective (linear) expression in the Wilson case, as it must be [9]. For  $a\mathbf{p} \neq \mathbf{0}$  we have  $A \cosh^2 + B \cosh + C = 0$  with

$$\begin{aligned} \frac{-B - \sqrt{B^2 - 4AC}}{2A} &= 1 + \frac{1}{2}(am)^2 - \frac{1}{2}(am)^3 + \frac{1}{2}(am)^4 - \frac{1}{2}(am)^5 \\ &+ \left[ \frac{1}{2} + \frac{1}{12}(am)^2 - \frac{1}{24}(am)^3 \right] (a\mathbf{p})^2 \\ &+ \frac{1}{12} \sum_{i < j} a^4 (p_i^2 p_j^2) + \left[ \frac{1}{24} + \frac{1}{24} am \right] \sum_i (ap_i)^4 + O(a^6) \end{aligned} \quad (48)$$

and upon taking  $\text{arcosh}(\cdot)$  one produces the dispersion relation (6).

### A.3 Dispersion relation for overlap operator with Wilson kernel

To establish the dispersion relation of the overlap operator with Wilson kernel one starts from

$$D_{W,-\rho} = i\bar{p}_\mu \gamma_\mu + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a}$$

$$D_{W,-\rho}^\dagger D_{W,-\rho} = (-i\bar{p}_\mu \gamma_\mu + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a})(i\bar{p}_\nu \gamma_\nu + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a}) = \bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2$$

as this yields the the free-field form of the desired operator and of its inverse:

$$\begin{aligned} D_{NW,m} &= (1 - \frac{am}{2\rho}) D_{NW} + m \quad \text{with} \quad D_{NW} = \frac{\rho}{a} \{1 + D_{W,-\rho} [D_{W,-\rho}^\dagger D_{W,-\rho}]^{-1/2}\} \\ D_{NW,m} &= \underbrace{\left( \frac{\rho}{a} + \frac{m}{2} \right)}_{\equiv c} + \underbrace{\left( \frac{\rho}{a} - \frac{m}{2} \right)}_{\equiv d} (i\bar{p}_\mu \gamma_\mu + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a}) [\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{-1/2} \\ G_{NW,m} &= \frac{c + d [\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{-1/2} (-i\bar{p}_\sigma \gamma_\sigma + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a})}{\{c + d [\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{-1/2} (-i\bar{p}_\mu \gamma_\mu + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a})\} \{c + d (i\bar{p}_\nu \gamma_\nu + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a}) [\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{-1/2}\}} \\ &= \frac{c + d [\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{-1/2} (-i\bar{p}_\sigma \gamma_\sigma + \frac{a}{2} \hat{p}^2 - \frac{\rho}{a})}{c^2 + 2cd(\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})[\dots]^{-1/2} + d^2[\dots]^{-1/2}[\dots]^{-1/2}} \cdot \end{aligned} \quad (49)$$

Hence, one ends up searching for zero in  $c^2 + 2cd(\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})[\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{-1/2} + d^2 = 0$  with

$$\bar{p}^2 = \frac{1}{a^2} \sum_\mu s_\mu^2 \quad \text{and} \quad \hat{p}^2 = \frac{8}{a^2} - \frac{2}{a^2} \sum_\mu c_\mu$$

or equivalently for a zero in  $(c^2 + d^2)[\bar{p}^2 + (\frac{a}{2} \hat{p}^2 - \frac{\rho}{a})^2]^{1/2} + 2cd(\frac{a}{2} \hat{p}^2 - \frac{\rho}{a}) = 0$ , and the square root makes this a *transcendental equation* in  $\cosh(aE)$ . The result through  $O(a^5)$  is given in (7).

### A.4 Dispersion relation for overlap operator with Brillouin kernel

To establish the dispersion relation of the overlap operator with Brillouin kernel one starts from

$$D_{B,-\rho} = i\tilde{p}_\mu \gamma_\mu + \frac{a}{2} \tilde{p}^2 - \frac{\rho}{a}$$

$$D_{B,-\rho}^\dagger D_{B,-\rho} = (-i\tilde{p}_\mu \gamma_\mu + \frac{a}{2} \tilde{p}^2 - \frac{\rho}{a})(i\tilde{p}_\nu \gamma_\nu + \frac{a}{2} \tilde{p}^2 - \frac{\rho}{a}) = \tilde{p}^2 + (\frac{a}{2} \tilde{p}^2 - \frac{\rho}{a})^2$$

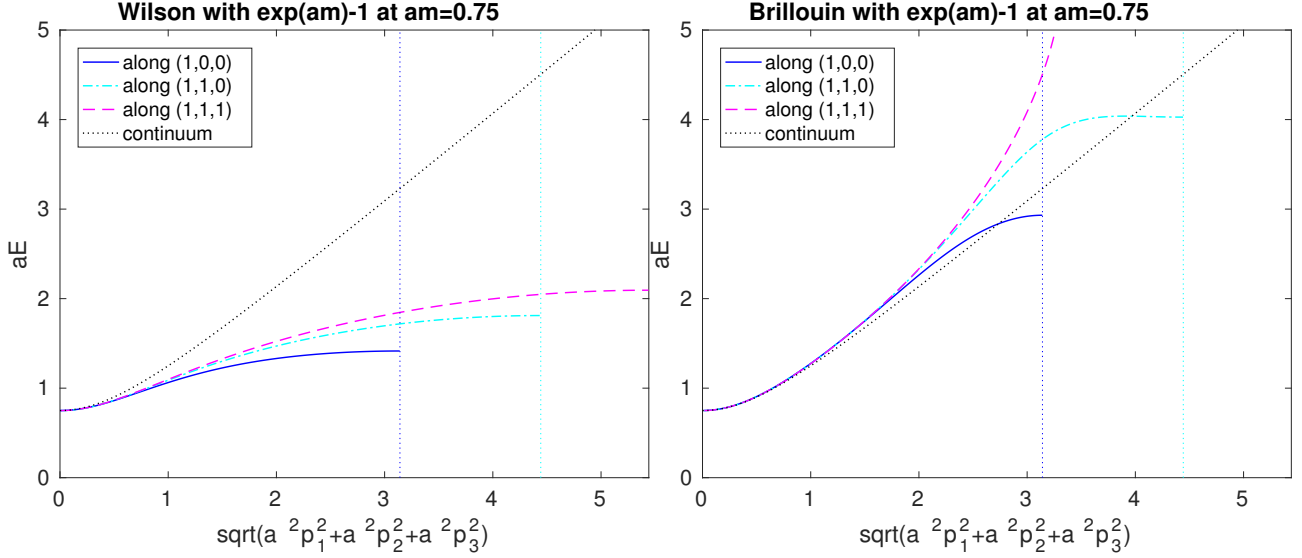


Figure 18: Same as the bottom part of Fig. 1 but with the replacement  $am \rightarrow \exp(am) - 1$  in the bare Wilson (left) and Brillouin (right) actions. Note that these dispersion relations extend over a larger portion of the Brillouin zone than the massive Wilson overlap and Brillouin overlap dispersion relations shown in the bottom part of Fig. 2.

as this yields the the free-field form of the desired operator and of its inverse:

$$\begin{aligned}
D_{\text{NB},m} &= (1 - \frac{am}{2\rho})D_{\text{NB}} + m \quad \text{with} \quad D_{\text{NB}} = \frac{\rho}{a}\{1 + D_{\text{B},-\rho}[D_{\text{B},-\rho}^\dagger D_{\text{B},-\rho}]^{-1/2}\} \\
D_{\text{NB},m} &= \underbrace{(\frac{\rho}{a} + \frac{m}{2})}_{\equiv c} + \underbrace{(\frac{\rho}{a} - \frac{m}{2})}_{\equiv d} (i\tilde{p}_\mu \gamma_\mu + \frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2} \\
G_{\text{NB},m} &= \frac{c+d[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2}(-i\tilde{p}_\sigma \gamma_\sigma + \frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})}{\{c+d[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2}(-i\tilde{p}_\mu \gamma_\mu + \frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})\}\{c+d(i\tilde{p}_\nu \gamma_\nu + \frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2}\}} \\
&= \frac{c+d[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2}(-i\tilde{p}_\sigma \gamma_\sigma + \frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})}{c^2 + 2cd(\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})[\dots]^{-1/2} + d^2[\dots]^{-1/2}[\dots]^{-1/2}} \quad (50)
\end{aligned}$$

Hence, one ends up searching for zero in  $c^2 + 2cd(\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{-1/2} + d^2 = 0$  with

$$\tilde{p}^2 = \frac{1}{729a^2} \sum_\mu s_\mu^2 \prod_{\nu \neq \mu} \{c_\nu + 2\}^2 \quad \text{and} \quad \tilde{p}^2 = \frac{4}{a^2} - \frac{1}{4a^2} \prod_\mu \{c_\mu + 1\}$$

or equivalently for a zero in  $(c^2 + d^2)[\tilde{p}^2 + (\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a})^2]^{1/2} + 2cd(\frac{a}{2}\tilde{p}^2 - \frac{\rho}{a}) = 0$ , and the square root renders this a *transcendental equation* in  $\cosh(aE)$ . The result through  $O(a^5)$  is given in (8).

## A.5 Dispersion relations with exponential quark mass trick

The observation that the first line in (4, 6) is just an expansion of  $\log^2(1 + am)$  [9] suggests that one might try the replacement  $am \rightarrow \exp(am) - 1$  in both the Wilson and the Brillouin actions (without the overlap procedure). With this substitution we find

$$\begin{aligned}
(aE)^2 - (a\mathbf{p})^2 &= (am)^2 \\
&+ \left[ -\frac{2}{3}(am)^2 + \frac{1}{2}(am)^3 \right] (a\mathbf{p})^2 \\
&+ \left[ -\frac{2}{3} + \frac{am}{2} \right] \left( \sum_{i < j} a^4 p_i^2 p_j^2 + \sum_i (ap_i)^4 \right) + O(a^6) \quad (51)
\end{aligned}$$

for the (heavy-quark) Wilson operator and

$$\begin{aligned}
(aE)^2 - (a\mathbf{p})^2 &= (am)^2 \\
&+ \left[0 + \frac{1}{12}(am)^3\right](a\mathbf{p})^2 \\
&+ \left[0 + \frac{am}{12}\right]\left(\sum_{i<j} a^4 p_i^2 p_j^2 + \sum_i (ap_i)^4\right) + O(a^6)
\end{aligned} \tag{52}$$

for the (heavy-quark) Brillouin operator. A comparison with (4, 6) shows that indeed the first lines (the momentum independent parts) now match the continuum behavior, while the coefficients in the second and third lines are at most as large (in magnitude) as before.

A similar conclusion is suggested by plotting the dispersion relation of either the Wilson operator or the Brillouin operator with this substitution, as done in Fig. 18. The cut-off effects at  $a\mathbf{p} = \mathbf{0}$  are gone (by construction), but even for non-zero momenta the free-field dispersion relation looks better than for the Wilson overlap and Brillouin overlap actions, respectively (see Fig. 2). Similarly to what was said about the Wilson overlap and Brillouin overlap actions, one might caution that such a behavior is only a necessary requirement for such a substitution to be useful in heavy-quark physics. But this time the investigation was carried out long ago, since this observation was the starting point for the development of the Fermilab action [12].

## B Details of diagonal Kenney-Laub iterates

In this appendix we give details of the partial fraction expansion of some of the diagonal elements of Kenney-Laub mappings as defined in Tabs. 1, 2.

The diagonal elements  $f_{n,n}$  for  $n = 1, \dots, 8$  are given in partial fraction form in Tab. 6. One notices that the smallest shift (fourth column) decreases with increasing  $n$ . Furthermore, for any fixed  $n$ , the weight (third column) is a monotonic (and positive) function of the shift. This means that the stopping criterion in the CG solver can be relaxed for smaller shifts. In fact, given the hierarchy among the shifts for a fixed  $f_{n,n}$ , it is clear that the cost of the numerical inversion is dominated by the cost of the smallest shift (or the smallest few shifts).

The elements  $f_{n-1,n}$  for  $n = 1, \dots, 3$  of the first upper diagonal are given in partial fraction form in Tab. 7. One notices that the majority of the features discussed in the previous paragraph persist, except that there is no constant contribution any more. Nonetheless, some of the properties of the overall function are quite different (see Sec. 3).

## C Flop count and memory traffic considerations

For an efficient implementation of the Brillouin operator it is vital to precompute the off-axis links that are implicitly used in the covariant derivative  $\nabla^{\text{iso}}$  and the covariant Laplacian  $\Delta^{\text{bri}}$ . The underlying reason is that, in order to maintain  $\gamma_5$ -hermiticity, one must average, within any  $k$ -hop contribution, over the  $k!$  shortest paths [with optional backprojection to  $SU(N_c)$  in a quenched setting, but we favor an average]. Moreover, each individual  $k$ -hop path requires  $k - 1$  matrix multiplications in color space.



KL(1,1)	1/3	0.8888888888888889	0.3333333333333333
KL(2,2)	1/5	0.4422291236000336	0.1055728090000841
		1.157770876399966	1.894427190999916
KL(3,3)	1/7	0.3005985953147677	0.05209508360168703
		0.4674182302787388	0.6359638059755859
		1.517697460120779	4.311941110422727
KL(4,4)	1/9	0.2291313786946141	0.03109120412576338
		0.2962962962962963	0.3333333333333333
		0.5378392501024903	1.420276625461206
		1.899696037869562	7.548632170413030
KL(5,5)	1/11	0.1855767632407464	0.02067219782410498
		0.2197383478725930	0.2085609132992616
		0.3183328723857196	0.7508307981214580
		0.6220419130221986	2.421230521622092
		2.290673739842379	11.59870556913308
KL(6,6)	1/13	0.1561143535399426	0.01474329800962692
		0.1759739298223624	0.1438305438453555
		0.2271454286305450	0.4764452860985422
		0.3498637189117062	1.274114172926091
		0.7123574780333906	3.630323607217039
		2.686237398754361	16.46054309190335
KL(13,13)	1/27	0.07432535480401804	0.003392289854243521
		0.07637712623153803	0.03109120412576338
		0.08071322905386415	0.08962859222716607
		0.08785700982653639	0.1860696326582413
		0.09876543209876543	0.3333333333333333
		0.1151288280773304	0.5542391790439610
		0.1400074395304560	0.8901004336611559
		0.1792797500341634	1.420276625461206
		0.2453107874470617	2.311695630535332
		0.3677579938327626	3.964732916742294
		0.6332320126231874	7.548632170413030
		1.392797081723151	17.80276060326253
		5.496102275704820	73.19738072201507

Table 6: Partial fraction form of the diagonal functions  $f_{n,n}(x) = x p_{n,n}(x^2)/q_{n,n}(x^2)$  for  $n = 1, \dots, 6, 13$  from the Kenney-Laub family (9) for the matrix sign function. The second column gives the constant contribution; the third and fourth columns give the numerator and the shift in the rational contribution. For instance  $f_{2,2}(x) = x(\frac{1}{5} + \frac{0.4422291236000336}{x^2+0.1055728090000841} + \frac{1.157770876399966}{x^2+1.894427190999916})$ .

KL(0,1)	2.000000000000000	1.000000000000000
KL(1,2)	0.5857864376269050	0.1715728752538099
	3.414213562373095	5.828427124746190
KL(2,3)	0.3572655899081636	0.07179676972449083
	0.6666666666666667	1.000000000000000
	4.976067743425180	13.92820323027551

Table 7: Partial fraction form of the first upper diagonal functions  $f_{n-1,n}(x)$  for  $n = 1, \dots, 3$  from the Kenney-Laub family (9) for the matrix sign function. The meaning of the columns is as in Tab. 6, except that this time there is no constant contribution.

## C.1 Brillouin operator flop count

Let  $U$  and  $V$  be objects which hold the original and smeared gauge fields, respectively. In Fortran-style languages (which use column-major memory layout) they may be defined as rank 7 arrays, e.g. `V(1:Nc,1:Nc,1:4,1:Nx,1:Ny,1:Nz,1:Nt)`. Here `1:4` in the third slot limits the values that the direction index  $\mu$  may take,  $N_c$  is the number of colors, and the box size is  $N_x \times N_y \times N_z \times N_t$ . Next, number the 81 elements in the  $[-1 : 1]^4$  hypercube around a given position  $n$  such that directions  $\nu \in \{1, \dots, 81\}$  and  $82 - \nu$  are opposite; in particular  $\nu = 41$  corresponds to the 0-hop movement. Since  $W_\nu(n)$  and  $W_{82-\nu}(n + \hat{\nu})$  relate to each other through hermitean conjugation, it suffices to store the first 40 off-axis links (constructed from  $V$ ) in the rank 7 array `W(1:Nc,1:Nc,1:40,1:Nx,1:Ny,1:Nz,1:Nt)`. In C-style languages (which use row-major memory layout) the ordering must be reversed, such that `W[[.]][[.]][[.]][[.]][[.]][[0:Nc-1]][[0:Nc-1]]` represents a  $N_c \times N_c$  matrix which occupies a contiguous space in memory. In the following we assume that  $W$  is ready for use, and we ignore this kind of set-up cost, since on the overall scale it is negligible.

We now discuss the structure of the matrix-times-vector routine which constructs, for a given source vector  $x$ , the target vector  $y = D^{\text{bril}}x$ . The source and sink vectors may be represented by rank 3 arrays, e.g. `x(1:Nc,1:4,1:Nx*Ny*Nz*Nt)` in Fortran-style languages. This routine consists of an outer loop (or set of four loops) which runs over the position  $n$  of the target  $y$ , and an inner loop (or set of four loops) which runs over the 81 elements of the hypercube around  $n$  and thus over the positions  $m$  of the source  $x$  which contribute to  $y(n)$ . In 80 out of the 81 cases the  $N_c \times 4$  matrix  $x(:, :, m)$  must be parallel transported through a left-multiplication with  $W(:, :, \nu, n)$  or  $W(:, :, \nu, m)^\dagger$ . In addition, the result (which is still a  $N_c \times 4$  matrix) must be right-multiplied with 0 ( $\nu = 41$ ) to 4 ( $\nu$  pointing to any of the 16 edges of the hypercube) elements of the set  $\{\gamma_1^t, \dots, \gamma_4^t\}$ , where  $t$  means transposition. In the chiral representation any  $\gamma$ -matrix contains one of the elements  $\pm 1, \pm i$  in each row and column, and the right-multiplication amounts to a re-ordering of the columns of this  $N_c \times 4$  matrix (times factors of  $\pm i$  which again implies reorderings of real and imaginary parts). Since such reorderings can be done on the fly, we assume that the right-multiplication is for free, and we take only the left-multiplication into account in our cost estimate.

With this input, the flop count of the Brillouin matrix-times-vector routine is as follows:

- (i)  $SU(N_c)$ -multiply the  $N_c \times 4$  block for each non-trivial direction. A complex-times-complex multiplication takes 6 flops, a complex-plus-complex addition takes 2 flops, there are  $N_c$  multiplications and  $N_c - 1$  additions per site, and there are 80 directions. Overall, this takes  $N_c \cdot 4 \cdot (6N_c + 2N_c - 2) \cdot 80$  flops; hence 21120 flops for  $SU(3)$ .
- (ii) Multiply the resulting  $N_c \times 4$  matrix with the correct weight factor as given by the isotropic derivative and the hypercubic Laplacian. These weight factors are real, and for each  $\nabla_\mu^{\text{iso}}$  non-zero only for 54 out of the 81 directions. The mass term may be incorporated into the 0-hop (i.e.  $\nu = 41$ ) contribution of the Laplacian. Overall, this takes  $N_c \cdot 8 \cdot (4 \cdot 54 + 81)$  flops; hence 7128 flops for  $SU(3)$ .
- (iii) Accumulate the 81 contributions to the out-spinor. Overall, this takes  $N_c \cdot 8 \cdot 80$  flops; hence 1920 flops for  $SU(3)$ .

All together we arrive at a grand total of 30168 flops per site for  $SU(3)$ .

## C.2 Wilson operator flop count

For reference, let us give a brief account how such a flop count looks for the Wilson operator. Here, the main difference is that for each one of the 8 directions  $\pm\mu$  the  $N_c \times 4$  block  $x(:, :, m)$  is right-multiplied by  $\frac{1}{2}(1 \pm \gamma_\mu)^t$ , and the latter operator is a projector whose eigenvectors can be precomputed. In consequence, the block is shrunk into  $N_c \times 2$  format before the left-multiplication with  $V_\mu(n)$  or  $V_\mu(m)^\dagger$  takes place, and expanded afterwards.

With this input, the flop count of the Wilson matrix-times-vector routine is as follows:

- (i) Spin project (from 4 to 2 components) the  $N_c \times 4$  matrix for each direction. Overall, this takes  $N_c \cdot 4 \cdot 8$  flops; hence 96 flops for  $SU(3)$ .
- (ii)  $SU(N_c)$ -multiply the  $N_c \times 2$  block for each direction, and expand back to  $N_c \times 4$  format (for free). Overall, this takes  $N_c \cdot 2 \cdot (6N_c + 2N_c - 2) \cdot 8$  flops; hence 1056 flops for  $SU(3)$ .
- (iii) Accumulate these 8 directions, as well as the 0-hop contribution which uses the precomputed factor  $(4 + m)$ . Overall, this takes  $N_c \cdot 8 \cdot 9$  flops; hence 216 flops for  $SU(3)$ .

All together we arrive at a grand total of 1368 flops per site for  $SU(3)$ .

## C.3 Brillouin operator memory traffic

The memory traffic of the Brillouin matrix-times-vector routine is as follows:

- (a) Read one color-spinor block for each direction. Overall, this amounts to  $N_c \cdot 8 \cdot 81$  floats; hence 1944 floats for  $SU(3)$ .
- (b) Read one gauge link  $W_\nu$  for each non-trivial direction. Overall, this amounts to  $N_c^2 \cdot 2 \cdot 80$  floats; hence 1440 floats for  $SU(3)$ .
- (c) Write one color-spinor block back into memory. Overall, this amounts to  $N_c \cdot 8$  floats; hence 24 floats for  $SU(3)$ .

All together we arrive at a grand total of 3408 floats per site for  $SU(3)$ , i.e. 13632 bytes if everything is in single-precision, and twice as much in double-precision. Here we assume that everything is to be read afresh, i.e. nothing is in cache. By handling  $N_v$  vectors simultaneously, the contribution (b) per vector is reduced by a factor  $N_v$ . For instance for  $N_c = 3$  and  $N_v = 12$  the grand total is  $1940 + 120 + 24 = 2088$  floats from/to memory per vector and site.

## C.4 Wilson operator memory traffic

The memory traffic of the Wilson matrix-times-vector routine is as follows:

- (a) Read one color-spinor block for each direction. Overall, this amounts to  $N_c \cdot 8 \cdot 9$  floats; hence 216 floats for  $SU(3)$ .
- (b) Read one gauge link  $V_\nu$  for each direction. Overall, this amounts to  $N_c^2 \cdot 2 \cdot 8$  floats; hence 144 floats for  $SU(3)$ .
- (c) Write one color-spinor block back into memory. Overall, this amounts to  $N_c \cdot 8$  floats; hence 24 floats for  $SU(3)$ .

All together we arrive at a grand total of 384 floats per site for  $SU(3)$ , i.e. 1536 bytes if everything is in single-precision, and twice as much in double-precision. Here we assume that everything is to be read afresh, i.e. nothing is in cache. By handling  $N_v$  vectors simultaneously, the contribution (b) per vector is reduced by a factor  $N_v$ . For instance for  $N_c = 3$  and  $N_v = 12$  the grand total is  $216 + 12 + 24 = 252$  floats from/to memory per vector and site.

## C.5 Technical summary

The Brillouin operator flop count exceeds the Wilson flop count by a factor  $30168/1368 \simeq 22$  at  $N_c = 3$ . In the large- $N_c$  limit the Brillouin flop count scales as  $N_c^2 \cdot 2560$ , while the Wilson flop count scales as  $N_c^2 \cdot 128$ . This means that in the large- $N_c$  limit this ratio approaches 20.

The Brillouin memory traffic exceeds the Wilson memory traffic by a factor  $3408/384 \simeq 8.9$  at  $N_c = 3$ , if one right-hand-side is handled at a time. In the large- $N_c$  limit the Brillouin memory traffic scales as  $N_c^2 \cdot 160$ , while the Wilson traffic scales as  $N_c^2 \cdot 16$ . This means that in the large- $N_c$  limit this ratio approaches 10.

At any  $N_c$  the memory traffic per site and right-hand-side can be reduced by handling  $N_v$  vectors simultaneously. Overall, this brings an extra factor  $N_v$  under (a) and (c), but no change under (b), for either operator. On a per-vector basis this means that the traffic under (b) is reduced by a factor  $N_v$ , while (a) and (c) remain constant. In other words, whenever memory bandwidth is the main bottleneck in an actual computation (which on highly parallel architectures is usually true) handling  $N_v$  right-hand-sides simultaneously is an efficient means to speed up both the Wilson and the Brillouin matrix-times-vector performance.

## References

- [1] K.G. Wilson, Phys. Rev. D **10**, 2445 (1974).
- [2] K.G. Wilson, *New Phenomena In Subnuclear Physics. Part A. Proceedings of the First Half of the 1975 International School of Subnuclear Physics, Erice, Sicily, July 11 - August 1, 1975*, ed. A. Zichichi, Plenum Press, New York, 1977, p. 69, CLNS-321.
- [3] S. Dürr and G. Koutsou, Phys. Rev. D **83**, 114512 (2011) [arXiv:1012.3615].
- [4] B. Sheikholeslami and R. Wohlert, Nucl. Phys. B **259**, 572 (1985).
- [5] G. Heatlie, G. Martinelli, C. Pittori, G.C. Rossi and C.T. Sachrajda, Nucl. Phys. B **352**, 266 (1991).
- [6] M. Lüscher, S. Sint, R. Sommer and P. Weisz, Nucl. Phys. B **478**, 365 (1996) [hep-lat/9605038].
- [7] M. Lüscher, S. Sint, R. Sommer, P. Weisz and U. Wolff, Nucl. Phys. B **491**, 323 (1997) [hep-lat/9609035].
- [8] S. Dürr, G. Koutsou and T. Lippert, Phys. Rev. D **86**, 114514 (2012) [arXiv:1208.6270].
- [9] Y. G. Cho, S. Hashimoto, A. Jüttner, T. Kaneko, M. Marinkovic, J. I. Noaki and J. T. Tsang, JHEP **1505**, 072 (2015) [arXiv:1504.01630].
- [10] L. Del Debbio, L. Giusti, M. Lüscher, R. Petronzio and N. Tantalo, JHEP **0602**, 011 (2006) [hep-lat/0512021].
- [11] S. Dürr *et al.* [BMW Collab.], JHEP **1108**, 148 (2011) [arXiv:1011.2711].

- [12] A. X. El-Khadra, A. S. Kronfeld and P. B. Mackenzie, Phys. Rev. D **55**, 3933 (1997) [hep-lat/9604004].
- [13] M. B. Oktay and A. S. Kronfeld, Phys. Rev. D **78**, 014504 (2008) [arXiv:0803.0523].
- [14] P. H. Ginsparg and K. G. Wilson, Phys. Rev. D **25**, 2649 (1982).
- [15] P. Hasenfratz, Nucl. Phys. Proc. Suppl. **63**, 53 (1998) [hep-lat/9709110].
- [16] P. Hasenfratz, Nucl. Phys. B **525**, 401 (1998) [hep-lat/9802007].
- [17] M. Lüscher, Phys. Lett. B **428**, 342 (1998) [hep-lat/9802011].
- [18] D. B. Kaplan, Phys. Lett. B **288**, 342 (1992) [hep-lat/9206013].
- [19] Y. Shamir, Nucl. Phys. B **406**, 90 (1993) [hep-lat/9303005].
- [20] V. Furman and Y. Shamir, Nucl. Phys. B **439**, 54 (1995) [hep-lat/9405004].
- [21] H. Neuberger, Phys. Lett. B **417**, 141 (1998) [hep-lat/9707022].
- [22] H. Neuberger, Phys. Lett. B **427**, 353 (1998) [hep-lat/9801031].
- [23] F. Niedermayer, Nucl. Phys. Proc. Suppl. **73**, 105 (1999) [hep-lat/9810026].
- [24] Y. B. Yang *et al.*, Phys. Rev. D **92**, no. 3, 034517 (2015) [arXiv:1410.3343].
- [25] B. Fahy *et al.* [JLQCD Collaboration], PoS LATTICE **2015**, 074 (2016) [arXiv:1512.08599].
- [26] P. Boyle, L. Del Debbio, A. Jüttner, A. Khamseh, F. Sanfilippo, J. T. Tsang and O. Witzel, arXiv:1611.06804 [hep-lat].
- [27] S. Dürr and G. Koutsou, arXiv:1610.06798 [hep-lat].
- [28] C.S. Kenney and A.J. Laub, SIAM J. Matrix Anal. Appl. **12**, 273 (1991).
- [29] Nicolas J. Higham, “Functions of matrices: theory and computation”, SIAM, 2008.
- [30] H. Neuberger, Phys. Rev. Lett. **81**, 4060 (1998) [hep-lat/9806025].
- [31] R. G. Edwards, U. M. Heller and R. Narayanan, Nucl. Phys. B **540**, 457 (1999) [hep-lat/9807017].
- [32] J. van den Eshof, A. Frommer, T. Lippert, K. Schilling and H. A. van der Vorst, Comput. Phys. Commun. **146**, 203 (2002) [hep-lat/0202025].
- [33] A. D. Kennedy, arXiv:hep-lat/0607038.
- [34] A. Frommer, B. Nockel, S. Gusken, T. Lippert and K. Schilling, Int. J. Mod. Phys. C **6**, 627 (1995) [hep-lat/9504020].
- [35] B. Jegerlehner, hep-lat/9612014.
- [36] R. C. Brower, H. Neff and K. Orginos, Nucl. Phys. Proc. Suppl. **140**, 686 (2005) [hep-lat/0409118].
- [37] T.-W. Chiu and S. V. Zenkin, Phys. Rev. D **59**, 074501 (1999) [hep-lat/9806019].
- [38] Y. Kikukawa and T. Noguchi, hep-lat/9902022.
- [39] S. Capitani, M. Göckeler, R. Horsley, P. E. L. Rakow and G. Schierholz, Phys. Lett. B **468**, 150 (1999) [hep-lat/9908029].
- [40] K.-F. Liu and S.J. Dong, Int. J. Mod. Phys. A **20**, 7241 (2005) [hep-lat/0206002].
- [41] P. Hasenfratz and F. Niedermayer, Nucl. Phys. B **414**, 785 (1994) [hep-lat/9308004].
- [42] T. A. DeGrand, A. Hasenfratz, P. Hasenfratz and F. Niedermayer, Nucl. Phys. B **454**, 587 (1995) [hep-lat/9506030].
- [43] W. Bietenholz and U.J. Wiese, Nucl. Phys. B **464**, 319 (1996) [hep-lat/9510026].

- [44] P. Hasenfratz, S. Hauswirth, K. Holland, T. Jörg, F. Niedermayer and U. Wenger, *Int. J. Mod. Phys. C* **12**, 691 (2001) [hep-lat/0003013].
- [45] P. Hasenfratz, S. Hauswirth, T. Jörg, F. Niedermayer and K. Holland, *Nucl. Phys. B* **643**, 280 (2002) [hep-lat/0205010].
- [46] P. Hernandez, K. Jansen and M. Lüscher, *Nucl. Phys. B* **552**, 363 (1999) [hep-lat/9808010].
- [47] W. Bietenholz, *Nucl. Phys. B* **644**, 223-247 (2002) [hep-lat/0204016].
- [48] T. A. DeGrand [MILC Collaboration], *Phys. Rev. D* **63**, 034503 (2000) [hep-lat/0007046].
- [49] T.G. Kovacs, *Phys. Rev. D* **67**, 094501 (2003) [hep-lat/0209125].
- [50] S. Dürr, C. Hoelbling and U. Wenger, *JHEP* **0509**, 030 (2005) [hep-lat/0506027].
- [51] M. Golterman and Y. Shamir, *Phys. Rev. D* **68**, 074501 (2003) [hep-lat/0306002].
- [52] M. Golterman, Y. Shamir and B. Svetitsky, *Phys. Rev. D* **71**, 071502 (2005) [hep-lat/0407021].
- [53] P. Hegde, F. Karsch, E. Laermann and S. Shcheredin, *Eur. Phys. J. C* **55**, 423 (2008) [arXiv:0801.4883].
- [54] M. Göckeler *et al.* [QCDSF Collab.], *Phys. Rev. D* **73**, 054508 (2006) [hep-lat/0601004].
- [55] G. S. Bali *et al.* [QCDSF Collaboration], *Prog. Part. Nucl. Phys.* **67**, 467 (2012) [arXiv:1112.0024].
- [56] T. Blum *et al.* [RBC and UKQCD Collaborations], *Phys. Rev. D* **93**, no. 7, 074505 (2016) [arXiv:1411.7017].
- [57] T. Bhattacharya, R. Gupta, W. Lee, S. R. Sharpe and J. M. S. Wu, *Phys. Rev. D* **73**, 034504 (2006) [hep-lat/0511014].
- [58] N. Cundy, J. van den Eshof, A. Frommer, S. Krieg, T. Lippert and K. Schäfer, *Comput. Phys. Commun.* **165**, 221 (2005) [hep-lat/0405003].
- [59] J. Brannick, A. Frommer, K. Kahl, B. Leder, M. Rottmann and A. Strebel, *Numer. Math.* (2015) [arXiv:1410.7170].
- [60] M. A. Clark, *PoS LAT* **2006**, 004 (2006) [hep-lat/0610048].
- [61] Z. Fodor, S. D. Katz, K. K. Szabo, *JHEP* **0408**, 003 (2004) [hep-lat/0311010].
- [62] T. A. DeGrand and S. Schaefer, *Phys. Rev. D* **71**, 034507 (2005) [hep-lat/0412005].
- [63] N. Cundy, S. Krieg, T. Lippert, A. Schäfer, *Comput. Phys. Commun.* **180**, 201-208 (2009) [arXiv:0803.0294].
- [64] C. Allton *et al.* [RBC-UKQCD Collaboration], *Phys. Rev. D* **78**, 114509 (2008) [arXiv:0804.0473].
- [65] J. Noaki *et al.* [JLQCD and TWQCD Collaborations], *Phys. Rev. Lett.* **101**, 202004 (2008) [arXiv:0806.0894].
- [66] S. Borsanyi, Y. Delgado, S. Dürr, Z. Fodor, S. D. Katz, S. Krieg, T. Lippert and D. Negradi *et al.*, *Phys. Lett. B* **713**, 342 (2012) [arXiv:1204.4089].
- [67] S. Dürr, C. Hoelbling, U. Wenger, *Phys. Rev. D* **70**, 094502 (2004) [hep-lat/0406027].
- [68] W. Bietenholz, I. Hip, S. Shcheredin and J. Volkholz, *Eur. Phys. J. C* **72**, 1938 (2012), [arXiv:1109.2649].