

# Hybridization expansion Monte Carlo simulation of multi-orbital quantum impurity problems: matrix product formalism and improved sampling

Hiroshi Shinaoka<sup>1</sup>, Michele Dolfi<sup>1</sup>, Matthias Troyer<sup>1</sup> and Philipp Werner<sup>2</sup>

<sup>1</sup> Theoretische Physik, ETH Zürich, 8093 Zürich, Switzerland

<sup>2</sup> Department of Physics, University of Fribourg, 1700 Fribourg, Switzerland  
E-mail: [shinaoka@itp.phys.ethz.ch](mailto:shinaoka@itp.phys.ethz.ch)

**Abstract.** We explore two complementary modifications of the hybridization-expansion continuous-time Monte Carlo method, aiming at large multi-orbital quantum impurity problems. One idea is to compute the imaginary-time propagation using a matrix product state representation. We show that bond dimensions considerably smaller than the dimension of the Hilbert space are sufficient to obtain accurate results and that this approach scales polynomially, rather than exponentially with the number of orbitals. Based on scaling analyses, we conclude that a matrix product state implementation will outperform the exact-diagonalization based method for quantum impurity problems with more than 12 orbitals. The second idea is an improved Monte Carlo sampling scheme which is applicable to all variants of the hybridization expansion method. We show that this so-called sliding window sampling scheme speeds up the simulation by at least an order of magnitude for a broad range of model parameters, with the largest improvements at low temperature.

**Keywords:** quantum Monte Carlo simulations

**ArXiv ePrint:** 1404.1259

---

## Contents

<b>1. Introduction</b>	<b>3</b>
<b>2. Hybridization expansion algorithm</b>	<b>4</b>
<b>3. Imaginary time evolution with the Krylov subspace method</b>	<b>6</b>
<b>4. Quantum impurity model</b>	<b>7</b>
<b>5. Trace calculation with matrix product states</b>	<b>7</b>
5.1. Matrix product state formalism . . . . .	7
5.1.1. Matrix product states (MPS) . . . . .	8
5.1.2. Compressing MPS . . . . .	8
5.1.3. Matrix product operators (MPO) . . . . .	8
5.1.4. Linear algebra with MPS and MPO . . . . .	9
5.2. Numerical details. . . . .	9
5.3. Accuracy of the MPS method . . . . .	10
5.4. Performance of the MPS method . . . . .	10
5.5. Discussion and future perspectives . . . . .	12
<b>6. Improved Monte Carlo sampling</b>	<b>13</b>
6.1. Conventional method . . . . .	14
6.2. Sliding-window approach . . . . .	14
6.3. Benchmark setup. . . . .	15
6.4. Benchmark results . . . . .	17
6.4.1. Insulating region: $U = 6$ and $\beta = 50$ . . . . .	17
6.4.2. Temperature and $U$ dependence. . . . .	17
6.5. Discussion and future perspectives . . . . .	19
<b>7. Summary</b>	<b>19</b>
<b>Acknowledgments</b>	<b>19</b>
<b>Appendix A. MPO for a model with uniform all-to-all interactions</b>	<b>20</b>
<b>Appendix B. MPO for general interactions</b>	<b>21</b>
<b>Appendix C. Ergodicity of the sliding-window approach</b>	<b>22</b>
<b>Appendix D. Detailed Monte Carlo update procedure</b>	<b>23</b>
<b>References</b>	<b>24</b>

---

## 1. Introduction

Quantum impurity models appear in various contexts in condensed matter physics. An important example is the dynamical mean-field theory (DMFT) [1] for strongly correlated electron systems. In a DMFT calculation, a correlated lattice model is mapped to an impurity problem whose bath degrees of freedom are self-consistently determined. Although the DMFT formalism was originally proposed for the single-band Hubbard model, it can be extended to multi-orbital systems and cluster-type impurities [2]. Furthermore, DMFT can be combined with density functional theory based *ab initio* calculations, to describe strongly correlated materials such as transition metal oxides [3]. For these applications, it is important to develop efficient algorithms to solve quantum impurity problems with multiple orbitals or sites.

In recent years, two complementary types of continuous-time quantum Monte Carlo (MC) impurity solvers have been developed, which are based on a stochastic sampling of perturbation expansions: the weak-coupling method [4] and the hybridization expansion method [5, 6]. The former approach is based on a perturbation expansion in powers of the Coulomb interaction terms, while the latter one treats the local Coulomb interactions exactly and instead expands the partition function in the coupling between the impurity and the bath. For describing strongly correlated materials, the latter approach is typically favored because of its ability to treat general interactions such as spin flips and because the average perturbation order of the hybridization expansion is relatively low in the strongly correlated regime. The algorithm was further extended to treat retarded interactions [7], which has recently been used in an extended DMFT study of the effects of long-range interactions [8].

A drawback of the hybridization expansion approach is that the computational effort scales exponentially with the number of sites or orbitals, because the dimension of the Hilbert space grows exponentially. Without additional approximations, this limits the application to small impurity models with up to five orbitals, even if one uses an implementation based on sparse-matrix exact-diagonalization techniques [9].

On the other hand, various wavefunction based theories have been developed for interacting fermionic lattice models. In particular, the ground states of one-dimensional (1D) systems can be described essentially exactly by the formalism of matrix product states (MPS) [10] with reasonable computational effort. The MPS formalism is known to be equivalent to the density matrix renormalization group (DMRG) [11, 12]. It has also been used to solve impurity problems [13–17]. In such MPS based calculations, the bath is represented by a 1D chain (or 1D chains) attached to the impurity, which results in an exponential growth of the computational cost with the number of sites or orbitals in the impurity. Furthermore, it is not trivial to extend the formalism to a non-diagonal coupling between the impurity and the bath, or to retarded interactions.

A possible direction for the development of flexible impurity solvers for large multi-orbital systems may be to combine these two approaches, i.e., the hybridization expansion and the MPS formalism. In this paper, we propose and test such a combined approach, in which the local interaction is treated using an MPS representation. More specifically, we perform the imaginary time evolution, which is given by the local impurity Hamiltonian, using the MPS formalism. We test the accuracy of the imaginary

time evolution and compare its performance with that of the exact approach using a sparse-matrix exact-diagonalization technique.

Another direction of research is to develop a more efficient MC sampling algorithm. In the continuous-time MC method based on the hybridization expansion, one stochastically samples configurations represented by creation and annihilation operators of the local degree of freedoms on the imaginary time interval. In estimating the weight of a configuration, the most costly part in multiorbital cases is evaluating the trace of a matrix product over the local degrees of freedom of the quantum impurity. This matrix product consists of imaginary-time evolution operators as well as creation and annihilation operators. The cost of evaluating the trace grows as temperatures is lowered, because the expansion order increases.

The trace can be evaluated either by the matrix formalism [6, 18], by sparse-matrix exact-diagonalization techniques (Krylov method) [9] or by an MPS version of the Krylov method. In the former formalism, all operators are represented by matrices in the eigenbasis of the local Hamiltonian and the matrix product is computed by multiplying the matrices one by one. In the latter formalism, the trace is computed by performing the imaginary-time evolution starting from eigenstates using the basis in which operators are represented as sparse matrices. In this paper, we call this the Krylov method or Krylov-sparse-matrix method. It was shown that the Krylov method is superior in performance for impurity problems involving more than 4 orbitals as local degrees of freedom [9].

For the matrix formalism, an efficient MC sampling scheme based on a tree structure has been proposed to suppress the growth of the computational cost at low temperatures [19]. Instead of recomputing the matrix product from scratch at each MC step, one reuses partial products of matrices that have been previously computed and stored. By using a tree data structure, the cost can then be reduced from  $O(\beta)$  to  $O(\log \beta)$ , where  $\beta$  is the inverse temperature. However, these ideas based on storing matrix products cannot be applied to the Krylov method. Thus, an alternative efficient MC sampling algorithm needs to be developed for the Krylov method.

The rest of the paper is organized as follows. In section 2, we describe the hybridization expansion algorithm. The Krylov method is described in section 3. The quantum impurity models used for the present study are defined in section 4. In section 5, we propose a combined approach of the Krylov method and the matrix-product formalism. We propose an improved MC sampling algorithm for the Krylov method in section 6. A summary is given in section 7

## 2. Hybridization expansion algorithm

A fermionic quantum impurity model is defined by the following Hamiltonian:

$$\mathcal{H} = \mathcal{H}_{\text{loc}} + \mathcal{H}_{\text{mix}} + \mathcal{H}_{\text{bath}}, \quad (1)$$

where

$$\mathcal{H}_{\text{loc}} = \sum_{\alpha, \beta} t_{\alpha, \beta} \hat{c}_{\alpha}^{\dagger} \hat{c}_{\beta} + \sum_{\alpha, \beta, \gamma, \delta} U^{\alpha, \beta, \gamma, \delta} \hat{c}_{\alpha}^{\dagger} \hat{c}_{\beta}^{\dagger} \hat{c}_{\gamma} \hat{c}_{\delta}, \quad (2)$$

$$\mathcal{H}_{\text{bath}} = \sum_{k,\alpha} \epsilon_{k,\alpha} \hat{a}_{k,\alpha}^\dagger \hat{a}_{k,\alpha}, \quad (3)$$

$$\mathcal{H}_{\text{mix}} = \sum_{k,\alpha,\beta} V_k^{\alpha,\beta} \hat{a}_{k,\alpha}^\dagger \hat{c}_\beta + \text{h.c.} \quad (4)$$

The term  $\mathcal{H}_{\text{loc}}$  describes an impurity with chemical potentials, intra-orbital hoppings and two-body interactions, where  $\alpha$  and  $\beta$  are combined orbital and spin indices. (We call the combined index of spin and orbital a flavor.)  $\mathcal{H}_{\text{bath}}$  describes a non-interacting bath with quantum numbers  $k$  and flavor  $\alpha$ . The hybridization term  $\mathcal{H}_{\text{mix}}$  describes the exchange of electrons between the impurity and the bath.

In the hybridization expansion impurity solver, one expands the partition function  $Z = \text{Tr}[e^{-\beta\mathcal{H}}]$  with respect to the hybridization term  $\mathcal{H}_{\text{mix}}$  as

$$\begin{aligned} Z &= \text{Tr}[e^{-\beta\mathcal{H}}] = \text{Tr}\left[e^{-\beta\mathcal{H}_1} T e^{-\int_0^\beta d\tau \mathcal{H}_2(\tau)}\right] \\ &= \sum_{n=0}^{\infty} \int_0^\beta d\tau_1 \cdots \int_{\tau_{n-1}}^\beta d\tau_n (-1)^n \text{Tr}\left[e^{-(\beta-\tau_n)\mathcal{H}_1} \mathcal{H}_2 e^{-(\tau_n-\tau_{n-1})\mathcal{H}_1} \cdots \mathcal{H}_2 e^{-\tau_1\mathcal{H}_1}\right], \end{aligned} \quad (5)$$

where  $\mathcal{H}_1 = \mathcal{H}_{\text{loc}} + \mathcal{H}_{\text{bath}}$  and  $\mathcal{H}_2 = \mathcal{H}_{\text{mix}}$  and we employed the interaction picture.

In equation (5), the partition function  $Z$  is represented as the sum of all configurations  $c = \{\tau_1, \dots, \tau_n\}$  with weight

$$w_c = (-d\tau)^n \text{Tr}\left[e^{-(\beta-\tau_n)\mathcal{H}_1} \mathcal{H}_2 e^{-(\tau_n-\tau_{n-1})\mathcal{H}_1} \cdots \mathcal{H}_2 e^{-\tau_1\mathcal{H}_1}\right] d\tau^n. \quad (6)$$

The weight can be simplified further by exploiting the fact that the time evolution of the impurity and the bath are not coupled by  $\mathcal{H}_2$ . By tracing out the bath degrees of freedom, one obtains

$$\begin{aligned} w_{\tilde{c}} &= Z_{\text{bath}} \text{Tr}_{\text{loc}} \left[ e^{-\beta\mathcal{H}_{\text{loc}}} T \hat{c}_{\alpha_n}(\tau_n) \hat{c}_{\alpha_1}^\dagger(\tau_n') \cdots \hat{c}_{\alpha_1}(\tau_1) \hat{c}_{\alpha_n}^\dagger(\tau_1') \right] \\ &\quad \times \det M^{-1}(\{\tau_1, \alpha_1\}, \dots, \{\tau_n, \alpha_n\}; \{\tau_1', \alpha_1'\}, \dots, \{\tau_n', \alpha_n'\}) (d\tau)^{2n}. \end{aligned} \quad (7)$$

Here,  $\tilde{c}$  represents a configuration with annihilation operators at  $\tau_1 < \dots < \tau_n$  with flavor  $\alpha_1, \dots, \alpha_n$  and creation operators at  $\tau_1' < \dots < \tau_n'$  with flavor  $\alpha_1', \dots, \alpha_n'$ . The matrix element of  $M^{-1}$  at  $(i, j)$  is given by the hybridization function  $\Delta_{\alpha_i', \alpha_j}(\tau_i' - \tau_j)$  defined in terms of  $\epsilon_{k,\alpha}$  and  $V_k^{\alpha,b}$ . The trace in equation (7) reduces to the form

$$\begin{aligned} \text{Tr}_{\text{loc}} \left[ e^{-(\beta-\tau_n)\mathcal{H}_{\text{loc}}} \hat{O}_{2n} e^{-(\tau_n-\tau_{n-1})\mathcal{H}_{\text{loc}}} \hat{O}_{2n-1} \cdots \hat{O}_1 e^{-\tau_1\mathcal{H}_{\text{loc}}} \right] \\ = \sum_m \left\langle \Psi_m \left| e^{-(\beta-\tau_n)\mathcal{H}_{\text{loc}}} \hat{O}_{2n} e^{-(\tau_n-\tau_{n-1})\mathcal{H}_{\text{loc}}} \hat{O}_{2n-1} \cdots \hat{O}_1 e^{-\tau_1\mathcal{H}_{\text{loc}}} \right| \Psi_m \right\rangle, \end{aligned} \quad (8)$$

where  $\hat{O}_1, \dots, \hat{O}_{2n}$  are time-ordered creation and annihilation operators appearing in equation (7).  $|\Psi_m\rangle$  denotes an eigenstate of  $\mathcal{H}_{\text{loc}}$  and the sum is over all eigenstates.

The contributions of the configurations  $\tilde{c}$  are stochastically sampled in the Monte Carlo simulation with the weight  $w_{\tilde{c}}$ . When  $\mathcal{H}_{\text{loc}}$  contains only chemical potentials and density-density interactions, the occupation number basis is an eigensystem of  $\mathcal{H}_{\text{loc}}$ . In this case, equation (8) can be evaluated efficiently. Otherwise, the evaluation of equation (8) is exponentially costly with respect to the number of orbitals in the impurity.

In [9], it was shown that the sum over eigenstates can be restricted to ground states at low enough temperature. It was also proposed to evaluate the trace using the so-called Krylov subspace method described in the next section.

### 3. Imaginary time evolution with the Krylov subspace method

In evaluating the trace in equation (8), we perform an imaginary time evolution

$$e^{-\tau\mathcal{H}}v \quad (9)$$

in each time-interval between creation/annihilation operators. We employ the Krylov subspace method in the same manner as in [9].

For a given Hamiltonian  $\mathcal{H}$  and vector  $v$ , the Krylov subspace is defined as

$$\mathcal{K}_p = \text{span}\{v, \mathcal{H}v, \dots, \mathcal{H}^{p-1}v\}, \quad (10)$$

where  $p$  is the dimension of the subspace. Then, the full matrix exponential  $e^{-\tau\mathcal{H}}v$  is approximated by the matrix exponential of the Hamiltonian projected onto the Krylov space.

We construct an orthonormal basis for the Krylov subspace that tridiagonalizes  $\mathcal{H}$  as

$$U^\dagger H U = T = \begin{pmatrix} \alpha_1 & \beta_1 & 0 & \cdots \\ \beta_1 & \alpha_2 & \beta_2 & \ddots \\ 0 & \beta_2 & \alpha_3 & \ddots \\ \vdots & \ddots & \ddots & \ddots \end{pmatrix} \quad (11)$$

by using the Lanczos method. Here,  $\alpha_i$  and  $\beta_i$  are real numbers. The column vectors of  $U$  are orthonormal basis vectors  $\{u_i\}$  with  $u_1 = v/\|v\|$ .

The basis vectors  $u_i$  and the matrix elements  $\alpha_i, \beta_i$  are obtained step by step for  $i = 1, 2, 3, \dots$  as follows:

$$\alpha_i = u_i^\dagger H u_i, \quad (12)$$

$$v_{i+1} = \begin{cases} H u_i - \alpha_i u_i & (i=1) \\ H u_i - \beta_{i-1} u_{i-1} - \alpha_i u_i & (i>1) \end{cases} \quad (13)$$

$$\beta_i = \|v_{i+1}\|, \quad (14)$$

$$u_{i+1} = v_{i+1} / \beta_i. \quad (15)$$

Convergence of the result is checked at each Lanczos step between equations (12) and (13) by evaluating the matrix exponential as

$$e^{-\tau H} v = \beta_0 e^{-\tau H} u_1 = \beta_0 \sum_{i=1}^p (e^{-\tau T})_{i1} u_i, \quad (16)$$

where  $\beta_0 = \|v\|$ . The matrix exponential  $e^{-\tau T}$  can be evaluated by a direct diagonalization because of the small dimension of the Krylov subspace. In the following calculations, we use the criterion  $|(e^{-\tau T})_{m1}/(e^{-\tau T})_{11}| < \epsilon$  with the tolerance  $\epsilon = 10^{-5}$ .

#### 4. Quantum impurity model

Throughout this paper, we consider an  $N$ -orbital impurity model with a ‘‘Slater-Kanamori’’ interaction. The Hamiltonian is

$$\begin{aligned} \mathcal{H}_{\text{loc}} = & \sum_i U \hat{n}_{i\uparrow} \hat{n}_{i\downarrow} - \mu \sum_i \hat{n}_i + \sum_{i>j,\sigma} [U' \hat{n}_{i\sigma} \hat{n}_{j-\sigma} + (U' - J) \hat{n}_{i\sigma} \hat{n}_{j\sigma}] \\ & - \sum_{i \neq j} J (\hat{c}_{i\downarrow}^\dagger \hat{c}_{j\uparrow}^\dagger \hat{c}_{j\downarrow} \hat{c}_{i\uparrow} + \hat{c}_{j\uparrow}^\dagger \hat{c}_{j\downarrow}^\dagger \hat{c}_{i\uparrow} \hat{c}_{i\downarrow}), \end{aligned} \quad (17)$$

where  $\hat{c}_i^\dagger$  and  $\hat{c}_i$  are creation/annihilation operators of an electron at site  $i$  and  $\hat{n}_i \equiv \hat{c}_i^\dagger \hat{c}_i$ . We take  $U' = U - 2J$  and  $J = U/6$ . The chemical potential is chosen such that the system is at half filling:  $\mu = \left(n - \frac{1}{2}\right) U - (n-1) \frac{5}{2} J$ . We consider an orbital-diagonal hybridization function corresponding to a noninteracting model with semicircular density of states of bandwidth 4.

While the interaction terms in equation (17) may not correspond to a rotationally invariant interaction for  $N > 3$ , we use this Hamiltonian for the purpose of benchmark calculations. We do not take into account the special conserved quantities [20] which enable a particularly efficient sampling of the Slater-Kanamori Hamiltonian. None of the procedures discussed in the following sections depend on a specific form of the Hamiltonian.

#### 5. Trace calculation with matrix product states

In this section, we investigate the accuracy and efficiency of a combined Krylov and MPS approach. A brief introduction of the MPS formalism is given in section 5.1. In section 5.2, we describe the details of benchmark calculations. In section 5.3 we discuss the accuracy of the method, while the performance of the method is investigated in section 5.4. Future perspectives are given in section 5.5.

##### 5.1. Matrix product state formalism

Here we provide a very brief overview of the MPS formalism. For details see the review by Schollwöck [21].

*5.1.1. Matrix product states (MPS).* Let us consider a one-dimensional lattice of length  $L$ , with a local Hilbert space of dimension  $d$  at each site. Hereafter, the dimension  $d$  is referred to as the local dimension. For instance, Hubbard models with  $S = 1/2$  electrons have a local dimension  $d = 4$ : The local Hilbert space at site  $i$  can be spanned by  $|0\rangle, \hat{c}_{i\downarrow}^\dagger|0\rangle, \hat{c}_{i\uparrow}^\dagger|0\rangle, \hat{c}_{i\uparrow}^\dagger\hat{c}_{i\downarrow}^\dagger|0\rangle$ .

Any pure state can be represented in the form

$$\begin{aligned} |\Psi\rangle &= \sum_{\sigma_1, \dots, \sigma_L} c_{\sigma_1, \dots, \sigma_L} |\sigma_1, \dots, \sigma_L\rangle = \sum_{\sigma_1, \dots, \sigma_L} \left( \sum_{b_1, \dots, b_L}^{r_1, \dots, r_L} M_{1, b_1}^{\sigma_1} M_{b_1, b_2}^{\sigma_2} \dots M_{b_{L-1}, b_L}^{\sigma_L} \right) \times |\sigma_1 \dots \sigma_L\rangle \\ &= \sum_{\sigma_1, \dots, \sigma_L} \mathbf{M}^{\sigma_1} \mathbf{M}^{\sigma_2} \dots \mathbf{M}^{\sigma_L} |\sigma_1 \dots \sigma_L\rangle, \end{aligned} \quad (18)$$

where  $\mathbf{M}^{\sigma_l}$  ( $l = 1, \dots, L$ ) are rank-3 tensors of dimension  $d \times r_{l-1} \times r_l$ . At the left ( $l = 1$ ) and right ( $l = L$ ) edges, we take  $r_0 = r_{L+1} = 1$ . The maximum value of  $b_l$  is referred to as the bond dimension of the MPS.

The MPS formalism is the underlying variational approximation made by the DMRG algorithm [21]. For a non-critical 1D system with short-range interactions, the ground state can be described very accurately by an MPS with a small bond dimension of  $O(1)$ . Note that the exponentially large tensor  $c_{\sigma_1, \dots, \sigma_L}$  is reduced to a product of small tensors of size  $O(1)$  because the entanglement entropy of the ground state is  $O(1)$  with respect to the system length.

*5.1.2. Compressing MPS.* An important remark is that MPS with a fixed bond dimension do not form a vector space. For example, the sum of two MPS results in a larger bond dimension as discussed later in section 5.1.4. In general, an MPS with a larger bond dimension can contain more information. Thus, to keep the bond dimension bounded, one may have to reduce the bond dimension after an operation, while keeping the loss of accuracy as small as possible. This can be done by an algorithm based on the so-called singular value decomposition (SVD). A truncation of the bond dimension from  $D'$  to  $D$  costs  $O(dD'^3L)$  for  $D' \gg D$ .

*5.1.3. Matrix product operators (MPO).* Matrix product operators are a natural generalization of the MPS concept to operators. Let us consider an arbitrary operator  $\hat{O}$ :

$$\hat{O} = \sum_{\sigma, \sigma'} O_{\sigma, \sigma'} |\sigma\rangle \langle \sigma'|. \quad (19)$$

The idea of MPS is directly applicable to operators by regarding  $(\sigma_l \sigma'_l)$  as one big index at each site. That is, the coefficients are represented as a product of local tensors as follows:

$$O_{\sigma, \sigma'} = \sum_{b_1, \dots, b_L}^{r_1, \dots, r_L} W_{1, b_1}^{\sigma_1 \sigma'_1} W_{b_1, b_2}^{\sigma_2 \sigma'_2} \dots W_{b_{L-1}, b_L}^{\sigma_L \sigma'_L} = \mathbf{W}^{\sigma_1 \sigma'_1} \mathbf{W}^{\sigma_2 \sigma'_2} \dots \mathbf{W}^{\sigma_L \sigma'_L}, \quad (20)$$

where the  $\mathbf{W}$ 's are now rank-4 tensors. The maximum value of  $b_l$  is referred to as the bond dimension of the MPO. We discuss how to construct an MPO for a given Hamiltonian in section 5.2.



5.1.4. *Linear algebra with MPS and MPO.* We can perform fundamental operations in quantum mechanics in the framework of MPS and MPO. One of the simplest examples is the summation of two wavefunctions  $|\varphi_1\rangle$  and  $|\varphi_2\rangle$ , as is required in equations (13). The sum of two MPSs with bond dimensions  $D_1$  and  $D_2$ , respectively, has a bond dimension of  $D' \leq D_1 + D_2$ . This can be understood by considering the sum of two MPS with bond dimension one:

$$|\phi_1\rangle = \sum_{\sigma} A^{\sigma_1} \dots A^{\sigma_L} |\sigma\rangle, \quad (21)$$

$$|\phi_2\rangle = \sum_{\sigma} B^{\sigma_1} \dots B^{\sigma_L} |\sigma\rangle. \quad (22)$$

One can easily see that the sum is given by

$$\sum_{\sigma} (A^{\sigma_1} B^{\sigma_1}) \begin{pmatrix} A^{\sigma_2} & 0 \\ 0 & B^{\sigma_2} \end{pmatrix} \times \dots \times \begin{pmatrix} A^{\sigma_{L-1}} & 0 \\ 0 & B^{\sigma_{L-1}} \end{pmatrix} \begin{pmatrix} A^{\sigma_L} \\ B^{\sigma_L} \end{pmatrix} |\sigma\rangle, \quad (23)$$

with a bond dimension of two. This can be extended to larger bond dimensions in a straightforward way. A sum of two MPS of bond dimension  $D$  requires only  $O(dD^2L)$  operations. However, it may be necessary to compress the resulting MPS to keep the bond dimension bounded at  $D$ . This cost dominates over the summation for  $D \gg 1$  because the compression is  $O(dD^3L)$ .

Another important operation is applying an operator  $\hat{O}$  to a wavefunction  $|\varphi\rangle$ , such as applying the Hamiltonian to a wavefunction in equation (13). Let us consider an MPO of bond dimension  $D_W$  and an MPS of bond dimension  $D$ . In this paper, we adopt an iterative approach which minimizes the residual  $\| |\tilde{\phi}\rangle - \hat{O}|\phi\rangle \|^2$  with respect to  $|\tilde{\phi}\rangle$  for a fixed bond dimension  $D$ . This algorithm scales as  $O(LD^3D_Wd)$  for  $1 \ll D_W \ll D$  [21].

## 5.2. Numerical details

The simulations in this section are carried out for the impurity model given in section 4. We take  $U = 6$  and  $J = U/6$  and  $\beta = 50$ . The Hamiltonian (17) can be represented by an MPO with a bond dimension of  $D_W \propto N_{\text{orb}}^2$  because the MPO for each term in equation (17) has a bond dimension of 1. This means that the computational effort scales polynomially with  $N_{\text{orb}}$  as  $O(D^3N_{\text{orb}}^3)$ . A further speed-up can be achieved by compressing the MPO. As will be explained in appendix A, the Hamiltonian can be represented by a more compact MPO with bond dimension eight irrespective of  $N_{\text{orb}}$  because the two-body interactions are homogeneous. Using this compact MPO, the computational effort now scales as  $O(D^3N_{\text{orb}})$ .

Although we consider the Slater–Kanamori interaction in this paper, the approach can be applied to any impurity model including general one- and two-body interactions like intra-orbital hopping and correlated hopping. As explained in appendix B, any one- and two-body interaction term can be represented by an MPO with bond dimension one. The compression of the MPO for the Hamiltonian is also possible for general one- and two-body interactions.

The following calculations were performed on a 2GHz Intel Core i7 CPU (Ivy Bridge), without parallelization. We used the Intel C++ Compiler v13.0 and the Math Kernel Library. The imaginary-time evolution was implemented with the MAQUIS/DMRG code [22]. The following results were obtained without exploiting good quantum numbers such as the total electron number. We found that exploiting conserved quantum numbers does not reduce the computational cost for the small bond dimensions  $D \leq 50$  used in this study.

We measure the timings and the accuracy using the Krylov-sparse-matrix and Krylov-MPS methods as follows. First, we perform Monte Carlo simulations with an exact solver (Krylov-sparse-matrix solver) in the same way as in [9]. After thermalization, we randomly select several configurations and measure timings. Calculations with the MPS method are then repeated for the same configurations using the MAQUIS/DMRG code. In the following, we measure the timings and the accuracy of the imaginary-time evolution for the ground state of the largest subspace with  $(N_\uparrow, N_\downarrow) = (3, 2), (4, 3), (4, 4), (5, 5)$  for  $N_{\text{orb}} = 5, 6, 7, 8$  and 10, respectively.

### 5.3. Accuracy of the MPS method

First, we discuss the accuracy of the MPS formalism. In figure 1, we show the convergence of the calculated value of the trace with respect to the bond dimension  $D$ . The impurity sizes are  $N_{\text{orb}} = 5, 7, 8$  and 10. The relative error is defined as  $|(t(D) - t_{\text{exact}}) / t_{\text{exact}}|$ , where  $t(D)$  and  $t_{\text{exact}}$  are the values of the trace calculated by the MPS formalism and the exact solver, respectively. The expansion order per flavor  $N_{\text{exp}}$  is 4.5–5 for  $N_{\text{orb}} = 5$  and 3–3.5 for  $N_{\text{orb}} = 7, 8$  and 10, respectively.

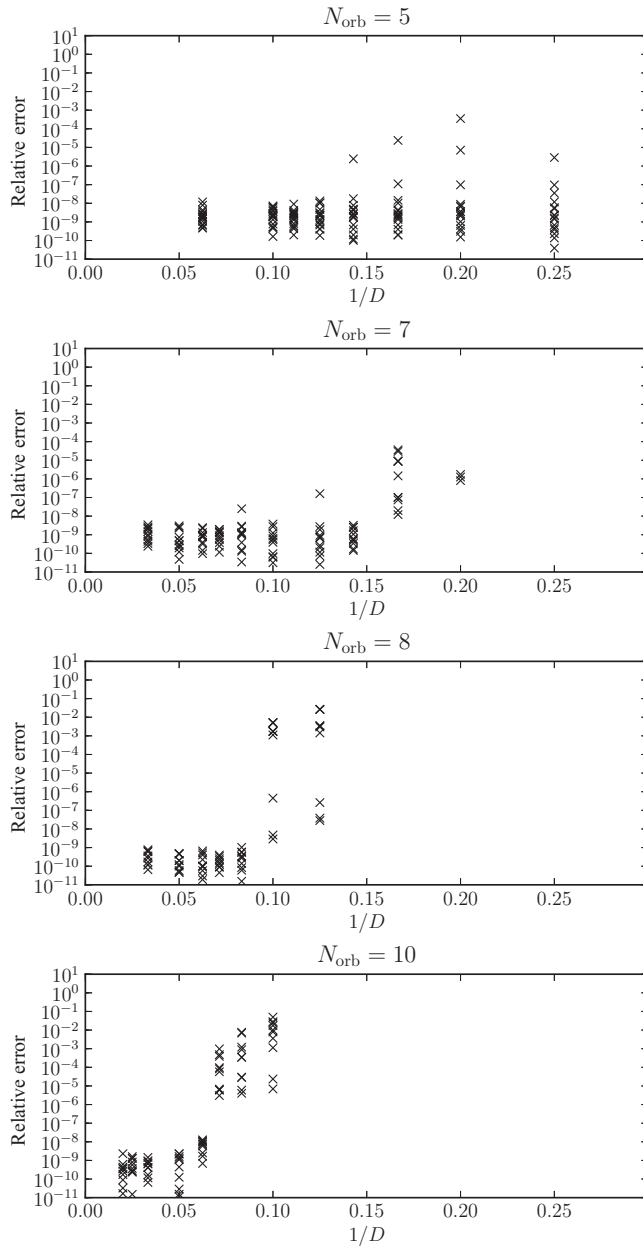
As seen in figure 1, the relative error of the MPS method decreases rapidly as  $D$  increases. For  $N_{\text{orb}} = 5$ , the results are already converged at  $D = 8$  for all sets of  $\{\widehat{O}(\tau_i)\}$ . For the largest system, i.e.,  $N_{\text{orb}} = 10$ , the relative error is well converged (and below  $10^{-7}$ ) at  $D = 16$ , even though the dimension of the Hilbert space is  $\binom{10}{5} = 63\,504$ . These results show that the MPS formalism yields accurate results even with a bond dimension considerably smaller than the dimension of the Hilbert space.

### 5.4. Performance of the MPS method

Next, we compare the performance of the two methods. Figure 2 shows the timing for an imaginary time evolution in the interval  $[0, \beta]$ . It is clearly seen that the timing for the exact solver increases exponentially with  $N_{\text{orb}}$ . The red broken line in figure 2 is a fit by

$$CN_{\text{orb}}^2 4^{N_{\text{orb}}}, \quad (24)$$

where  $C$  is a positive constant. Equation (24) is derived as follows. The most costly operation in the imaginary time evolution is applying a sparse matrix  $\mathcal{H}$  to a dense vector in equations (12) and (13). Each such operation costs  $O(N_{\text{orb}}^2 D_{\text{Hilbert}})$ , where the dimension of the largest subspace  $D_{\text{Hilbert}}$  is given by  $\binom{N_{\text{orb}}}{N_{\text{orb}}/2} \propto 4^{N_{\text{orb}}} / N_{\text{orb}}$ . Assuming that the expansion order per orbital is  $O(1)$ , we immediately arrive at equation (24). As shown in figure 2, the data are well fitted by equation (24) for  $N_{\text{orb}} \geq 7$  with  $C = 2.5 \times 10^{-8}$ .

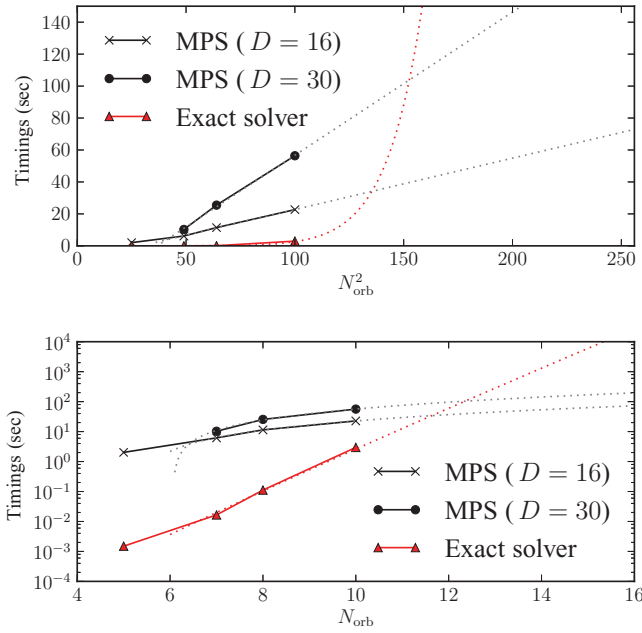


**Figure 1.** Convergence of the value of the trace with respect to the bond dimension  $D$  for  $N_{\text{orb}} = 5, 7, 8, 10$ , respectively.

On the other hand, the timing for the MPS formalism is expected to scale as  $O(D^3 N_{\text{orb}}^2)$  for a fixed  $D$ . This comes from the fact that applying  $\mathcal{H}$  to an MPS costs  $O(D^3 N_{\text{orb}})$ . As seen in figure 2, the data are indeed well fitted by the expected scaling

$$a(N_{\text{orb}}^2 - b) \tag{25}$$

with  $a$  and  $b$  positive constants. We note that the estimated value of  $a$  increases only slightly from 0.322 to 0.896 as  $D$  increases from  $D = 16$  to  $D = 30$ , though one expects a  $(30/16)^3 (\simeq 6.59)$  time increase. This may be due to overhead in treating many small matrices for small  $D$ . This can be seen more explicitly when we plot the timings as



**Figure 2.**  $N_{\text{orb}}$  dependence of the timings for the imaginary time evolution in the interval  $[0, \beta]$ . The data are averaged over 10 different operator configurations  $\{\widehat{O}(\tau_i)\}$  for each  $N_{\text{orb}}$ . The red broken line and the dotted black line are the fit by equations (24) and (25), respectively.

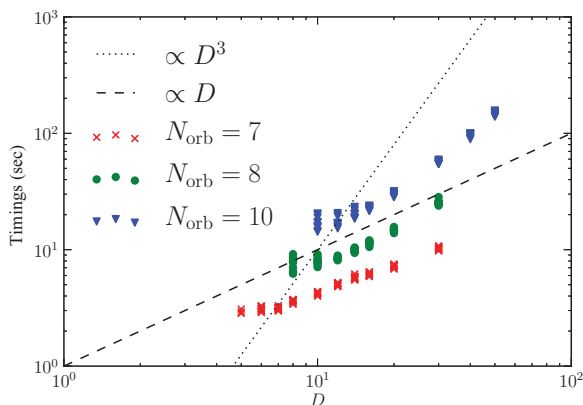
a function of  $D$  for each  $N_{\text{orb}}$  in figure 3. It is obvious that the timing increases more slowly than the expected asymptotic scaling  $O(D^3)$  for  $D \leq 50$ .

Even for the largest  $N_{\text{orb}}$  considered ( $N_{\text{orb}} = 10$ ), the MPS formalism with  $D = 16$  runs about 10 times slower than the exact solver. However, the MPS formalism is expected to become more efficient than the exact solver for larger  $N_{\text{orb}}$ . Extrapolating the timings of the two methods using equations (24) and (25), the crossover point is estimated to be  $N_{\text{orb}} = 12\text{--}13$ , with only a slight dependence on the value of  $D$  (see the lower panel of figure 3).

### 5.5. Discussion and future perspectives

Our results show that the Krylov-MPS formalism can be potentially superior to the exact Krylov-sparse-matrix solver for a large number of orbitals  $N_{\text{orb}} \gtrsim 12$ . Impurity problems with  $N_{\text{orb}} \geq 12$  are relevant for example for cluster-type DMFT calculations of multi-orbital Hubbard models. However, the MC simulation of such large impurity problems is not feasible at the moment with our present code. (To date, simulations with hybridization expansion solvers have been restricted to at most 7 orbitals.) Thus, in this section, we discuss how the performance might be improved.

In a MC simulation, we update  $\{\widehat{O}(\tau_i)\}$  by an elementary update such as inserting or removing a pair of annihilation and creation operators. Each operator must be updated before the MC sampling loses its memory of the original configuration. Thus, the autocorrelation time  $\tau_{\text{auto}}$  is expected to be roughly  $2N_{\text{orb}}N_{\text{exp}}/p_{\text{acc}}$  in units of elementary updates ( $N_{\text{exp}}$  is the expansion order per flavor). The acceptance rate  $p_{\text{acc}}$  depends on the system and on parameters such as  $\beta$ . It is typically on the order



**Figure 3.** Bond-dimension  $D$  dependence of the timings for an imaginary time evolution in the interval  $[0, \beta]$ . Data for  $N_{\text{orb}} = 7, 8$  and  $10$  are shown. The different points represent data for different operator configurations  $\{\widehat{O}(\tau_i)\}$ .

of 0.01–0.1. Assuming  $p_{\text{acc}} = 0.1$  and  $N_{\text{exp}} = 3$  (a typical value in the strongly correlated regime and for temperatures of about 1% of the bandwidth) for  $N_{\text{orb}} = 12$ , we obtain  $\tau_{\text{auto}} \simeq 700$  elementary updates. Recalling that the timing for evaluating the trace is  $O(10^2)$  seconds (see figure 2), the autocorrelation time is estimated to be  $O(10^5)$  s or 30 h. In a weakly correlated metal, where the perturbation order is higher, the autocorrelation time is on the order of a week. This is too long for practical DMFT calculations.

There are possible ways to reduce the autocorrelation time. First, we can increase the acceptance rate  $p_{\text{acc}}$  by proposing several candidates at each MC update. Evaluating their weights can be assigned to different nodes. By using the heat bath algorithm or a better algorithm [23], the acceptance rate  $p_{\text{acc}}$  can be increased to almost 1. Another factor of 10 can be gained by using the improved MC sampling introduced in section 6, which avoids recomputing the full imaginary-time evolution from scratch. By using these two tricks, the autocorrelation time  $\tau_{\text{auto}}$  can be reduced to  $O(10^3)$  s, which is still not short enough for practical applications.

Another possible way is to speed up the imaginary time evolution by parallel computing. However, this is not trivial because the bond dimensions of MPS and MPO tensors are quite small in the case of quantum impurity problems.

## 6. Improved Monte Carlo sampling

In this section, we propose an improved Monte Carlo sampling procedure which significantly reduces autocorrelation times for multi-orbital impurity problems. This sampling strategy can be used both in the Krylov-sparse-matrix method and the Krylov-MPS method. After reviewing previously proposed improved sampling strategies in section 6.1, we describe our new MC sampling scheme in section 6.2. We explain the details of benchmark calculations in section 6.3. The benchmark results are shown in section 6.4 and future perspectives are discussed in section 6.5.

## 6.1. Conventional method

In a hybridization-expansion continuous-time Monte Carlo simulation, one updates a current configuration by proposing a new configuration which is slightly different from the current one: For example, one tries to insert or remove a pair of creation and annihilation operators. Then, the new configuration is accepted stochastically according to the ratio of the weights of the current and new configurations. Naively, one might evaluate the weight of the new configuration by performing an imaginary time evolution in the full interval  $[0, \beta]$ . However, the computational cost of such a calculation grows linearly with the expansion order  $N_{\text{exp}}$ . This is costly at low temperature or in a metallic phase, where the expansion order is large.

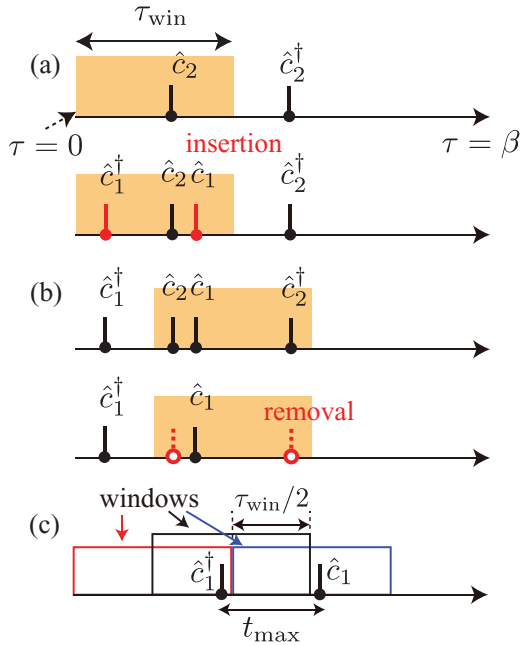
For the matrix formalism, in which the trace is evaluated using matrix products, Haule proposed a trick to improve the efficiency of the Monte Carlo sampling [18]. Here, we introduce a related idea in the context of the Krylov method. The improved sampling strategy proposed in [18] is based on the observation that the insertion or removal of pairs of operators is predominantly a local (in imaginary time) process. In other words, the acceptance rate for an insertion or removal of a pair of operators with a large time difference is very low. It was furthermore proposed to do the time evolution from both sides, storing the resultant matrix products at several intermediate  $\tau$  points. Then, when trying to insert or remove a pair of operators with a short time difference, the evaluation of the trace only requires the time evolution in a short time interval. Although this makes trial steps cheaper, one has to recompute the intermediate results once an update is accepted. Since this requires the time evolution from both sides, which typically costs  $O(N_{\text{exp}})$ , the method does not change the scaling with respect to  $N_{\text{exp}}$ . Thus, the method gives a significant improvement in performance only when the acceptance rate is quite low and  $N_{\text{exp}}$  is not so large.

## 6.2. Sliding-window approach

We now propose an improved update scheme in which the computational cost of an elementary update stays constant with respect to the expansion order  $N_{\text{exp}}$ . Although the exponential scaling with the number of orbitals is not affected, this method substantially reduces the prefactor of the scaling at low temperatures. Although we explain the idea in the context of the Krylov algorithm, it can be applied to the matrix formalism as well.

First, to make the maximum use of the locality in the imaginary time, we introduce an upper bound  $t_{\text{max}}$  on the time difference between the two operators which we try to insert or remove. As mentioned in the previous study,  $t_{\text{max}}$  can be almost independent of  $\beta$  and  $N_{\text{exp}}$ . In addition, we introduce an imaginary-time window in which updates are allowed (see figure 4(a)). The window width  $\tau_{\text{win}} = \beta / N_{\text{win}}$  is taken to be larger than (but on the order of)  $t_{\text{max}}$ . Now, similarly to [18], one performs the time evolution from both sides and stores the results at the end-points of the window. This allows us to evaluate the trace for a new configuration at constant cost.

After several updates, the window is moved to the next position by  $\tau_{\text{win}}/2$  (see figure 4(b)). Concurrently, one updates the wave vectors at the end-points, which again costs only  $O(\tau_{\text{win}}) = O(1)$ . This procedure is repeated so that the window moves back and forth in the whole interval  $[0, \beta]$ . This procedure is ergodic because we can produce operator pairs with arbitrary separation by inserting one with a short separation and then gradually



**Figure 4.** (a) Inserting a pair of operators in the window. Updates are allowed only in the window. (b) The window is moved by  $\tau_{\text{win}}/2$  from (a). Now, the pair  $\{\hat{c}_2^\dagger, \hat{c}_2\}$  can be removed. (c) The pair does not fit into any of the three windows. This can happen if  $\tau_{\text{win}} < 2t_{\text{max}}$ .

increasing the separation through MC updates. We refer to appendix C for a proof of ergodicity. The advantage of our algorithm is that the cost of each MC step is independent of  $\beta$  and reperforming the time evolution over the full time interval is not required.

By definition,  $\tau_{\text{win}}$  needs to be larger than  $t_{\text{max}}$ . In practical simulations, however,  $\tau_{\text{win}}$  has a stricter lower bound:

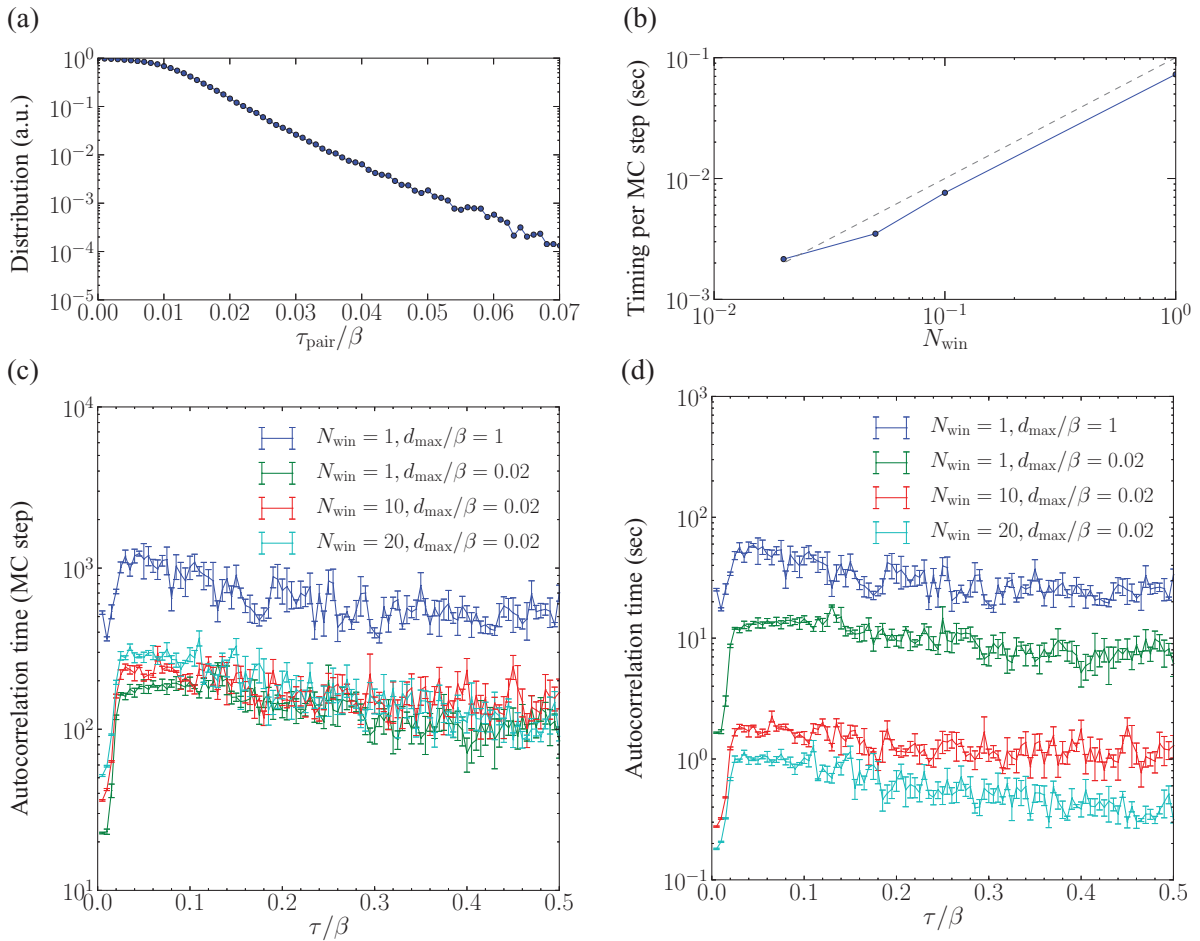
$$\tau_{\text{win}} > 2t_{\text{max}}. \quad (26)$$

Let us consider a pair of operators with a time difference of  $t_{\text{max}}$  as shown in figure 4(c). When  $\tau_{\text{win}} < 2t_{\text{max}}$ , this pair does not fit into any of the local windows shown there and thus cannot be removed by a single elementary update. Although this does not break the ergodicity, the autocorrelation time may increase. This problem can be avoided by taking  $\tau_{\text{win}} > 2t_{\text{max}}$ .

If the window moves from one side to the other fast enough compared to the autocorrelation time, the sequential sweep should not badly affect the autocorrelation time. As discussed in section 5.5, the autocorrelation time is  $O(N_{\text{exp}}N_{\text{orb}}/p_{\text{acc}})$ . On the other hand, moving the window from one side to the other takes  $2N_{\text{win}} = O(N_{\text{orb}}N_{\text{exp}})$  MC steps. Because these two time scales are always of the same order, the autocorrelation time should not be severely affected.

### 6.3. Benchmark setup

In this section, we show benchmark results for a 5-orbital impurity problem. At each position of the window, we try to insert or remove a pair  $N_f$  times, where  $N_f$  is the

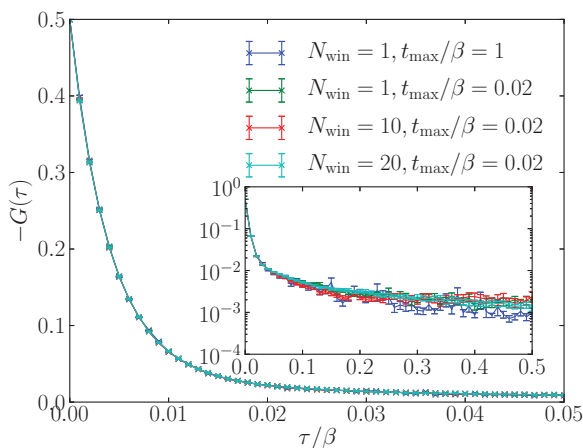


**Figure 5.** Benchmark results for a 5-orbital model with  $U = 6$  and  $J = 1$ . (a) Distribution of the time difference of a successfully removed or inserted pair of operators in the MC sampling performed with  $t_{\text{max}}/\beta = 1$  and  $N_{\text{win}} = 1$ . (b) Timing per MC step. The broken line represents the relation  $\text{timing} \propto \tau_{\text{win}}$ . (c), (d) Autocorrelation time of  $G(\tau)$  in units of MC steps (c) and seconds (d).

number of flavors. When trying to insert a pair, the time difference between the operators is chosen randomly in the interval  $[0, t_{\text{max}}]$ . Correspondingly, when we try to remove a pair in the window, we first list all pairs of creation and annihilation operators with a time difference equal to or less than  $t_{\text{max}}$ . Then, we try to remove one of them. The detailed procedure is described in appendix D.

The following simulations were done by the Krylov algorithm based on sparse-matrix techniques on one CPU core of AMD Opteron 6174 (2.2 GHz). We divide the Hilbert space into sectors by using the total particle number  $\hat{N}$  and magnetization  $\hat{S}_z$  as conserved quantum numbers. All data are averaged over four independent MC runs of fixed  $1.28 \times 10^7$  steps. The diagonal Green's function  $G(\tau)$  is measured on 1001 points on the imaginary-time interval. We symmetrize  $G(\tau)$  by using the particle-hole symmetry. The autocorrelation time of  $G(\tau)$  is estimated by a binning analysis for each time point using bins of 16384 MC steps. Simulations were performed using the ALPS libraries [24].





**Figure 6.**  $G(\tau)$  computed at  $U = 6$  and  $\beta = 50$ . The inset shows a log-scale plot.

## 6.4. Benchmark results

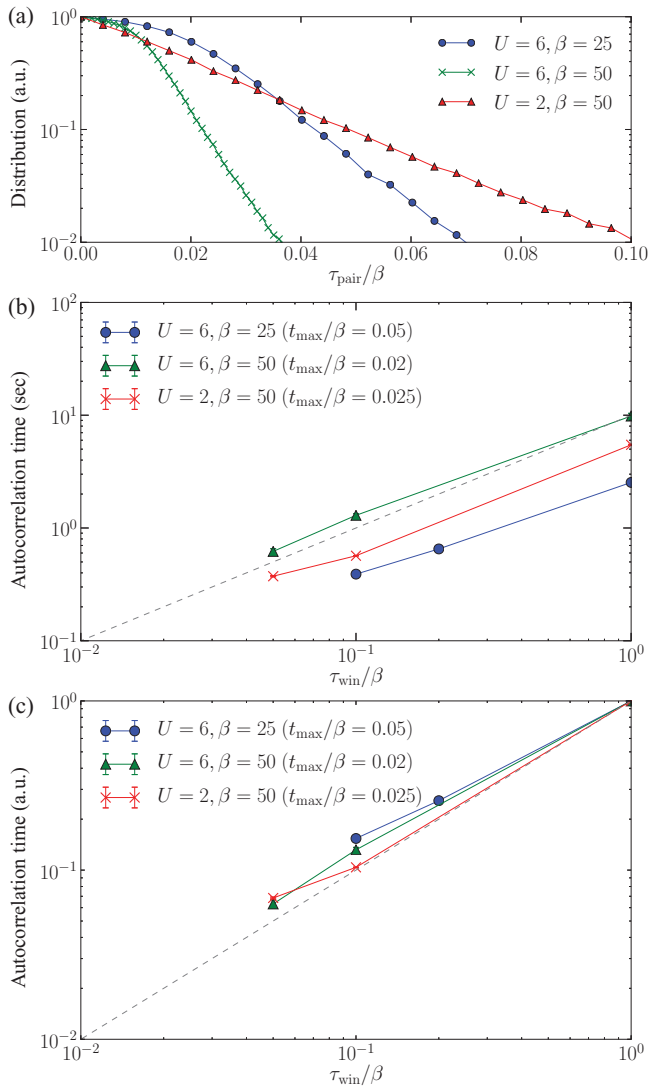
*6.4.1. Insulating region:  $U = 6$  and  $\beta = 50$ .* In figure 5(a), we show the distribution of the time difference between pairs of operators successfully removed or inserted in the Monte Carlo sampling performed with  $t_{\max}/\beta = 1$  and  $N_{\text{win}} = 1$ . As expected, we see that the accepted updates are local in imaginary time. The distribution decreases exponentially for large time differences. We found that the range  $\tau_{\text{pair}}/\beta \leq 0.02$  accounts for almost 94% of the successful updates. Considering the condition in equation (26), we take  $t_{\max}/\beta = 0.02$  and  $N_{\text{win}} \leq 20$ . Since  $N_{\text{exp}} \simeq 5.6$ , the window contains on average 5.6 operators for  $N_{\text{win}} = 20$ .

In figure 6, we show the Green's function  $G(\tau)$  computed using  $t_{\max}/\beta = 0.02$  for different values of  $\tau_{\text{win}}$ . We also present data obtained for  $t_{\max}/\beta = 1$  for comparison. All the data shown are consistent within error bars, indicating our algorithm works correctly. However, we found that  $G(\tau)$  for  $N_{\text{win}} = 1$  and  $t_{\max} = 1$  is systematically smaller than the others. This may be because the autocorrelation time is too long for the MC simulation to be thermalized.

In figure 5(b), we show the  $N_{\text{win}}$  dependence of the timing per MC step for  $t_{\max}/\beta = 0.02$ . As expected, the timing decreases linearly with the window size  $\tau_{\text{win}}$ . The estimated autocorrelation time is shown in figures 5(c) in units of MC steps. For  $N_{\text{win}} = 1$ , the autocorrelation time is shorter by one order of magnitude for  $t_{\max}/\beta = 0.02$  compared to that for  $t_{\max}/\beta = 1$ . This is consistent with the increase in the acceptance rate from 0.022 to 0.34 by introducing the cutoff. Now, we discuss the  $N_{\text{win}}$  dependence. Around  $\tau \simeq \beta/2$ , the autocorrelation time is not affected badly by introducing the window, consistent with the above argument. Although the autocorrelation is affected around  $\tau = 0.01$ , the increase is considerably smaller than the reduction in the CPU time.

Figure 5(d) shows the autocorrelation time in units of seconds. It is clearly seen that the autocorrelation becomes shorter in the entire  $\tau$  region as  $N_{\text{win}}$  increases up to  $N_{\text{win}} = 20$ . The improvement is as much as two orders of magnitude from the most naive approach ( $N_{\text{win}} = 1$  and  $t_{\max}/\beta = 1$ ) to the best case ( $N_{\text{win}} = 20$  and  $t_{\max}/\beta = 0.02$ ).

*6.4.2. Temperature and  $U$  dependence.* Figure 7(a) shows the distribution function of the length of successfully inserted and removed pairs of operators for different values of



**Figure 7.** (a) Distribution of the time difference of a pair of operators successfully removed or inserted in the MC sampling. (b), (c)  $\tau_{\text{win}}$  dependence of the autocorrelation time in CPU time. The autocorrelation time is averaged over the interval of  $0 < \tau < \beta$ . In (c), the data are normalized by the timings for  $\tau_{\text{win}}/\beta = 1$ .

$\beta$  and  $U$ . The weakly correlated metallic region corresponds to  $U \lesssim 2$ . First, we discuss the temperature dependence for  $U = 6$ . Comparing the data for  $\beta = 25$  and  $\beta = 50$ , one can see that the distribution becomes more localized at low temperatures. This may be because the hybridization function  $\Delta(\tau)$  decays more rapidly with  $\tau$  similarly to  $G(\tau)$  at low temperatures. This result indicates that our improved MC sampling works even better at low temperatures. On the other hand, although the distribution becomes broader at smaller  $U$ , the distribution still decays exponentially at long distances. In figures 7(b) and (c), we plot the  $\tau_{\text{win}}$  dependence of the autocorrelation time averaged over the interval  $0 < \tau < \beta$ . We took  $t_{\text{max}} = 0.05, 0.02$  and  $0.025$  for  $(U = 6, \beta = 25)$ ,  $(U = 6, \beta = 50)$  and  $(U = 2, \beta = 50)$ . It is clearly seen that the autocorrelation time scales linearly with  $\tau_{\text{win}}$  down to the lowest  $\tau_{\text{win}}$  for all parameter sets. This indicates the robustness of the sliding window approach.

## 6.5. Discussion and future perspectives

A simple way to choose the window size is to measure the distribution function of the distance between successfully inserted or removed pairs of operators during the thermalization process. Then, one can choose a reasonable cutoff  $t_{\max}$  such that most of the distribution, say 95%, is contained within the cutoff. The window size  $\tau_{\text{win}} = \beta / N_{\text{win}}$  is then given by the minimum size that satisfies the lower bound given in equation (26).

Further improvements of the efficiency may be possible by using the heat-bath algorithm or a better algorithm [23] where we propose several candidates at each update. This allows to increase the acceptance rate and reduce the autocorrelation time.

There are other kinds of local updates with acceptance rates higher than inserting/removing pairs of operators. Examples include shifting an operator on the imaginary time axis or swapping two nearest neighboring operators. Introducing such efficient updates helps in practical calculations.

## 7. Summary

In this paper, we discussed two complementary approaches based on the hybridization-expansion continuous-time Monte Carlo method for multi-orbital systems. First, we proposed to combine the Krylov approach with the MPS/MPO representation of states and operators. We found that highly accurate results can be obtained by using bond dimensions considerably smaller than the dimension of the whole Hilbert space. Based on a scaling analysis, we showed that the performance becomes superior to the conventional method for quantum impurity problems involving more than 12 orbitals.

Second, we proposed an improved Monte Carlo sampling algorithm for the hybridization expansion Monte Carlo method. Detailed benchmark tests were carried out for a 5-orbital impurity model. We showed that the new algorithm works robustly for a broad range of on-site repulsions and temperatures. In particular, we confirmed that the “sliding window” approach works particularly efficiently at low temperatures and we expect that it will be useful in the study of phenomena emerging at low temperatures. The sampling scheme is easy to implement in existing Monte Carlo codes and applies to any variant of the hybridization expansion method.

## Acknowledgments

We thank Iztok Pizorn for discussions on matrix product states. We also thank Jakub Imriska, Hidemaro Suwa, Hugo Strand, Lei Wang and Li Huang for useful comments on the improved Monte Carlo sampling. HS and PW acknowledge support from the DFG via FOR 1346 and SNF Grant 200021E-149122. This project was supported by ERC grant SIMCOFE. Simulations were performed using the ALPS libraries [24].

## Appendix A. MPO for a model with uniform all-to-all interactions

Let us consider a Hamiltonian with uniform all-to-all interactions:

$$\mathcal{H} = \sum_{n=1}^{N_{\text{op}}} \sum_{i \geq 1, j \geq 2, i < j}^L \widehat{A}_i^{(n)} \widehat{B}_j^{(n)} + \sum_{i=1}^L \widehat{O}_i, \quad (\text{A.1})$$

where  $\widehat{A}_i^{(n)}$  and  $\widehat{B}_j^{(n)}$  are operators acting on the local Hilbert spaces on sites  $i$  and  $j$ , respectively.  $\widehat{O}_i$  is an operator acting on site  $i$ . A compressed MPO can be explicitly constructed for this kind of model with all-to-all uniform interactions.

The Hamiltonian in equation (A.1) may be written in the form

$$\mathcal{H} = \mathbf{W}_1 \mathbf{W}_2 \cdots \mathbf{W}_L, \quad (\text{A.2})$$

where  $\mathbf{W}_i$  is a matrix whose elements are operators acting on the local Hilbert space at site  $i$ . The  $\mathbf{W}_i$  are given as follows:

$$\mathbf{W}_1 = \left( \widehat{I} \quad \widehat{A}^{(1)} \quad \cdots \quad \widehat{A}^{(N_{\text{op}})} \quad \widehat{O} \right), \quad (\text{A.3})$$

$$\mathbf{W}_i = \begin{pmatrix} \widehat{I} & \widehat{A}^{(1)} & \cdots & \widehat{A}^{(N_{\text{op}})} & \widehat{O} \\ 0 & \widehat{I} & 0 & 0 & \widehat{B}^{(1)} \\ 0 & 0 & \ddots & 0 & \vdots \\ 0 & 0 & 0 & \widehat{I} & \widehat{B}^{(N_{\text{op}})} \\ 0 & 0 & 0 & 0 & \widehat{I} \end{pmatrix}, \quad (\text{A.4})$$

$$\mathbf{W}_L = \begin{pmatrix} \widehat{O} \\ \widehat{B}^{(1)} \\ \vdots \\ \widehat{B}^{(N_{\text{op}})} \\ \widehat{I} \end{pmatrix}, \quad (\text{A.5})$$

where  $1 < i < L$  and  $\widehat{I}$  denotes the identity operator. One can see that equation (A.2) is in the MPO form with bond dimension  $N_{\text{op}} + 2$  when each element in  $\mathbf{W}_i$  is regarded as a  $4 \times 4$  matrix.

For the multi-orbital Hubbard model given in equation (17), one obtains an MPO of bond dimension eight by taking

$$\widehat{A}^{(1)} = n_{\uparrow}, \quad (\text{A.6})$$

$$\widehat{B}^{(1)} = (U' - J)n_{\uparrow} + U'n_{\downarrow}, \quad (\text{A.7})$$

$$\widehat{A}^{(2)} = n_{\downarrow}, \quad (\text{A.8})$$

$$\widehat{B}^{(2)} = (U' - J)n_{\downarrow} + U'n_{\uparrow}, \quad (\text{A.9})$$

$$\widehat{A}^{(3)} = S^+ (\equiv c_{\uparrow}^{\dagger} c_{\downarrow}), \quad (\text{A.10})$$

$$\widehat{B}^{(3)} = -JS^- \quad (\equiv -Jc_{\downarrow}^{\dagger}c_{\uparrow}), \quad (\text{A.11})$$

$$\widehat{A}^{(4)} = S^-, \quad (\text{A.12})$$

$$\widehat{B}^{(4)} = -JS^+, \quad (\text{A.13})$$

$$\widehat{A}^{(5)} = D^+ \quad (\equiv c_{\uparrow}^{\dagger}c_{\downarrow}^{\dagger}), \quad (\text{A.14})$$

$$\widehat{B}^{(5)} = -JD^- \quad (\equiv -Jc_{\uparrow}c_{\downarrow}), \quad (\text{A.15})$$

$$\widehat{A}^{(6)} = D^-, \quad (\text{A.16})$$

$$\widehat{B}^{(6)} = -JD^+, \quad (\text{A.17})$$

$$\widehat{O} = U\widehat{n}_{\uparrow}\widehat{n}_{\downarrow}. \quad (\text{A.18})$$

For the local Hilbert space spanned by  $|0\rangle, \widehat{c}_{i\downarrow}^{\dagger}|0\rangle, \widehat{c}_{i\uparrow}^{\dagger}|0\rangle, \widehat{c}_{i\uparrow}^{\dagger}\widehat{c}_{i\downarrow}^{\dagger}|0\rangle$ ,

$$c_{\uparrow}^{\dagger} = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad (\text{A.19})$$

$$c_{\downarrow}^{\dagger} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \quad (\text{A.20})$$

$$n_{\uparrow} = c_{\uparrow}^{\dagger}c_{\uparrow}, \quad (\text{A.21})$$

$$n_{\downarrow} = c_{\downarrow}^{\dagger}c_{\downarrow}. \quad (\text{A.22})$$

## Appendix B. MPO for general interactions

In this appendix, we show how the MPS formalism is extended to general interactions. Let us begin by showing the MPO representation of annihilation and creation operators. In the operator representation, they look like

$$\widehat{f}_1 \otimes \widehat{f}_2 \otimes \cdots \otimes \widehat{f}_{i-1} \otimes \widehat{O}_i \otimes \widehat{I}_{i+1} \cdots \otimes \widehat{I}_L, \quad (\text{B.1})$$

with the site index explicitly shown. We omit the spin index for simplicity. Here,  $\widehat{O}$  is the matrix representation of the annihilation or creation operators given in appendix A. The operator  $\widehat{f}_i$  counts the number of particles and is given by

$$f = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (\text{B.2})$$

in the local basis introduced in appendix A. Therefore, annihilation and creation operators are obviously represented by an MPO with bond dimension one. Since the product of two MPO with bond dimension one has bond dimension one, any product of annihilation and creation operators can be represented by an MPO with bond dimension one. For example, a correlated hopping term  $\hat{n}_1 \hat{c}_2^\dagger \hat{c}_4$  reads

$$\hat{n}_1 \otimes \hat{c}_2^\dagger \hat{f}_2 \otimes \hat{f}_3 \otimes \hat{f}_4 \hat{c}_4 \otimes \hat{I}_5 \otimes \dots \quad (\text{B.3})$$

The summation over the site index can be explicitly taken in a way similar to that in appendix A. Let us consider the sum of correlated hopping terms

$$\sum_{i \neq j \neq k} \hat{n}_i \hat{c}_j^\dagger \hat{c}_k \quad (\text{B.4})$$

as an example. For simplicity, we restrict ourselves to the case  $i < j < k$ . In this case, the sum is represented by the MPO with the following local tensors:

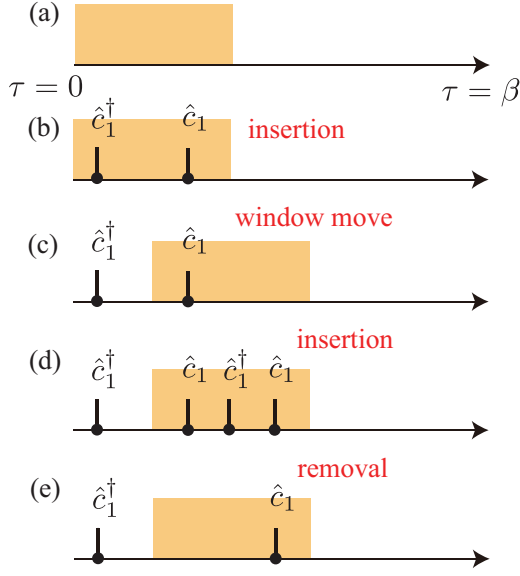
$$\mathbf{W}_1 = \begin{pmatrix} \hat{I} & \hat{n} & 0 & 0 \end{pmatrix}, \quad (\text{B.5})$$

$$\mathbf{W}_i = \begin{pmatrix} \hat{I} & \hat{n} & 0 & 0 \\ 0 & \hat{I} & \hat{c}^\dagger \hat{f} & 0 \\ 0 & 0 & \hat{f} & \hat{f} \hat{c} \\ 0 & 0 & 0 & \hat{I} \end{pmatrix} \quad (1 < i < L), \quad (\text{B.6})$$

$$\mathbf{W}_L = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \hat{I} \end{pmatrix}. \quad (\text{B.7})$$

### Appendix C. Ergodicity of the sliding-window approach

In this appendix, we show that the MC sampling based on the sliding window approach is ergodic. In particular, we show that a pair of operators with an arbitrary time difference can be inserted by repeated insertions and removals of pairs with a short time difference. The procedure is illustrated in figure B1. First, we insert a pair in the window as shown in figures B1(a) and (b). Then, the window is moved to the next position (figure B1(c)). As shown in figure B1(d), the distance between the operators can be increased by inserting a new pair and removing two operators in the middle because the two windows are overlapping each other. By repeating this procedure, one can create a pair with an arbitrary time difference. One can also remove any pair of operators, independent of the time difference, by reversing the above procedure. Therefore, it is obvious that one can transform any configuration into any other configuration by inserting and removing pairs within the sliding window.



**Figure B1.** How to insert a pair of operators with an arbitrary time difference.

## Appendix D. Detailed Monte Carlo update procedure

The local Monte Carlo update procedure has been described in section II B of [18]. In this appendix, we explain how this procedure is modified when the cutoff  $t_{\max}$  and the sliding window are introduced.

Let us consider an attempt to insert a pair of creation and annihilation operators of flavor  $f$  at  $\tau_c$  and  $\tau_a$ . More specifically, we first choose  $\tau_c$  randomly and uniformly in the window. Then,  $\tau_a$  is chosen randomly and uniformly in the window under the constraint  $|\tau_c - \tau_a| \leq t_{\max}$ . The reverse process of this update is removing one operator pair of flavor  $f$  whose length is equal to or less than  $t_{\max}$ .

We first discuss the case without a cutoff  $t_{\max}$ . The window is located on the interval  $[\tau_{\text{win}}^{\min}, \tau_{\text{win}}^{\max}]$  with  $\tau_{\text{win}} = \tau_{\text{win}}^{\max} - \tau_{\text{win}}^{\min}$ . The weights of the original and new configurations are denoted by  $w_{\text{org}}$  and  $w_{\text{new}}$ , respectively. The probability to accept this insertion is

$$P = \min \left[ 1, \left| \frac{w_{\text{new}}}{w_{\text{org}}} \right| \frac{\tau_{\text{win}}^2}{N_{\text{pair}}^{t_{\max}}} \right], \quad (\text{D.1})$$

where  $\tau_{\text{win}} = \tau_{\text{win}}^{\max} - \tau_{\text{win}}^{\min}$  is the size of the window and  $N_{\text{pair}}$  is the number of operator pairs of flavor  $f$  in the window after the insertion.

By introducing a cutoff  $t_{\max}$ , the probability is changed to

$$P = \min \left[ 1, \left| \frac{w_{\text{new}}}{w_{\text{org}}} \right| \frac{\tau_{\text{win}} \Delta \tau_a}{N_{\text{pair}}^{t_{\max}}} \right], \quad (\text{D.2})$$

where

$$\Delta \tau_a = \min(\tau_c + t_{\max}, \tau_{\text{win}}^{\max}) - \max(\tau_c - t_{\max}, \tau_{\text{win}}^{\min}) \quad (\text{D.3})$$

and  $N_{\text{pair}}^{t_{\text{max}}}$  is the number of operator pairs of flavor  $f$  whose length is equal to or less than  $t_{\text{max}}$  in the window after the insertion.

The probability to accept an attempt to remove a pair at  $\tau_c$  and  $\tau_a$  is correspondingly given by

$$P = \min \left[ 1, \left| \frac{w_{\text{new}}}{w_{\text{old}}} \right| \frac{N_{\text{pair}}^{t_{\text{max}}}}{\tau_{\text{win}} \Delta \tau_a} \right], \quad (\text{D.4})$$

where  $N_{\text{pair}}^{t_{\text{max}}}$  is the number of operator pairs of flavor  $f$  for the original configuration.

## References

- [1] Georges A, Kotliar G, Krauth W and Rozenberg M J 1996 *Rev. Mod. Phys.* **68** 13–125
- [2] Maier T, Jarrell M, Pruschke T and Hettler M H 2005 *Rev. Mod. Phys.* **77** 1027–80
- [3] Kotliar G, Savrasov S Y, Haule K, Oudovenko V S, Parcollet O and Marianetti C A 2006 *Rev. Mod. Phys.* **78** 865–951
- [4] Rubtsov A N, Savkin V V and Lichtenstein A I 2005 *Phys. Rev. B* **72** 035122
- [5] Werner P, Comanac A, de' Medici L, Troyer M and Millis A J 2006 *Phys. Rev. Lett.* **97** 076405
- [6] Werner P and Millis A J 2006 *Phys. Rev. B* **74** 155107
- [7] Werner P and Millis A J 2010 *Phys. Rev. Lett.* **104** 146401
- [8] Ayrat T, Biermann S and Werner P 2013 *Phys. Rev. B* **87** 125149
- [9] Läuchli A M and Werner P 2009 *Phys. Rev. B* **80** 235117
- [10] Östlund S and Rommer S 1995 *Phys. Rev. Lett.* **75** 3537–40
- [11] White S R 1992 *Phys. Rev. Lett.* **69** 2863–6
- [12] White S R and Noack R M 1992 *Phys. Rev. Lett.* **68** 3487–90
- [13] Nishimoto S and Jeckelmann E 2004 *J. Phys.: Condens. Matter* **16** 613
- [14] Raas C, Uhrig G S and Anders F B 2004 *Phys. Rev. B* **69** 041102
- [15] Raas C and Uhrig G S 2005 *Eur. Phys. J. B: Condens. Matter Complex Syst.* **45** 293–303
- [16] Nishimoto S, Pruschke T and Noack R M 2006 *J. Phys.: Condens. Matter* **18** 981
- [17] Peters R 2011 *Phys. Rev. B* **84** 075139
- [18] Haule K 2007 *Phys. Rev. B* **75** 155113
- [19] Gull E, Millis A J, Lichtenstein A I, Rubtsov A N, Troyer M and Werner P 2011 *Rev. Mod. Phys.* **83** 349–404
- [20] Parragh N, Toschi A, Held K and Sangiovanni G 2012 *Phys. Rev. B* **86** 155158
- [21] Schollwöck U 2011 *Ann. Phys.* **326** 96–192
- [22] Dolfi M *et al* 2014 in preparation
- [23] Suwa H and Todo S 2010 *Phys. Rev. Lett.* **105** 120603
- [24] Bauer B *et al* 2011 *J. Stat. Mech.* P05001