

Variational optimization algorithms for uniform matrix product states

V. Zauner-Stauber,¹ L. Vanderstraeten,² M.T. Fishman,³ F. Verstraete,^{1,2} and J. Haegeman²

¹*Vienna Center for Quantum Technology, University of Vienna, Boltzmanngasse 5, 1090 Wien, Austria*

²*Ghent University, Faculty of Physics, Krijgslaan 281, 9000 Gent, Belgium*

³*Institute for Quantum Information and Matter,
California Institute of Technology, Pasadena, California 91125, USA*

We combine the Density Matrix Renormalization Group (DMRG) with Matrix Product State tangent space concepts to construct a variational algorithm for finding ground states of one dimensional quantum lattices in the thermodynamic limit. A careful comparison of this variational uniform Matrix Product State algorithm (VUMPS) with infinite Density Matrix Renormalization Group (IDMRG) and with infinite Time Evolving Block Decimation (ITEBD) reveals substantial gains in convergence speed and precision. We also demonstrate that VUMPS works very efficiently for Hamiltonians with long range interactions. The new algorithm can be conveniently implemented as an extension of an already existing DMRG implementation.

I. INTRODUCTION

The strategy of renormalization group (RG) techniques to successively reduce a large number of microscopic degrees of freedom to a smaller set of effective degrees of freedom has led to powerful numerical and analytical methods to probe and understand the effective macroscopic behavior of both classical and quantum many body systems.^{1–4} However, it was not until the advent of White's celebrated *Density Matrix Renormalization Group* (DMRG)^{5,6} that variational RG methods reached unprecedented accuracy in numerically studying strongly correlated one dimensional quantum lattice systems at low temperature. The underlying variational ansatz of *Matrix Product States* (MPS)^{7–13} belongs to a class of ansatzes known as *Tensor Network States*.^{11,14,15} These variational classes encode the many body wavefunction in terms of virtual entanglement degrees of freedom living on the boundary and thus satisfy an area law scaling of entanglement entropy per construction. As such, they provide a natural parameterization for the physical corner of Hilbert space, where low energy states of quantum many body systems ought to live in.^{16,17} MPS in particular are especially fit for studying ground states of strongly correlated one dimensional quantum systems with local interactions.^{18–20}

The variational parameters in MPS are contained within local tensors associated with the individual sites of the lattice system. For homogeneous systems, the global wave function can then be captured using just a single (or a small number of) such tensors, independent of the system size. They consequently offer very natural access to the thermodynamic limit, providing a clear advantage over other numerical approaches such as Exact Diagonalization or Quantum Monte Carlo.

On finite lattices, (one-site) DMRG implements the variational principle (energy minimization) by exploiting that the quantum state is a multilinear function of the local tensors. By fixing all but one tensors, the global eigenvalue problem is transformed into an effective eigenvalue problem for the local tensor.^{5,6,12,21–23}

Using a translation invariant parameterization gives rise to an energy expectation value with a highly non-linear dependence on the tensor(s). Two different algorithms are widely used to obtain such an MPS in the thermodynamic limit. *Infinite system DMRG* (IDMRG)^{5,6,24} proceeds by performing regular DMRG on a successively growing lattice, inserting and optimizing over new tensors in the center of the lattice in each step only, effectively mimicking an infinite lattice by using a finite, albeit very large lattice. After convergence the most recently inserted tensors in the center are taken as a unit cell for an infinite MPS approximation of the ground state. An alternative approach is known as *infinite time evolving block decimation* (ITEBD).^{25,26} It works directly in the thermodynamic limit and is based on evolving an initial state in imaginary time by using a Trotter decomposition of the evolution operator.

We present a new variational algorithm, inspired by tangent space ideas,^{13,27,28} that combines the advantages of IDMRG and ITEBD and addresses some of their shortcomings. As such it is directly formulated in the thermodynamic limit, but at the same time optimizes the state by solving effective eigenvalue problems, rather than employing imaginary time evolution. We find that it leads to a significant increase in efficiency in all of our test cases. The following section introduces MPS notations and definitions and presents our variational algorithm, heuristically motivated from the perspective of finite size DMRG. Sec. III illustrates the performance of our algorithm on various test cases, and compares to conventional IDMRG and ITEBD results. After the conclusion in Sec. IV, we provide further technical details in the appendices. Appendix A contains additional theoretical background: we derive the self-consistent conditions that characterize the variational minimum and provide additional motivation for our algorithm from the perspective of the MPS tangent space. Appendix B presents a suitable strategy to expand the bond dimension of translation invariant MPS. Appendix C explains how to construct effective Hamiltonians in the thermodynamic limit. These involve infinite geometric sums of the transfer matrix, which are further studied in Appendix D.

II. A VARIATIONAL ALGORITHM FOR MATRIX PRODUCT STATES IN THE THERMODYNAMIC LIMIT

In this section we introduce a variational algorithm for optimizing MPS directly in the thermodynamic limit. Because the algorithm strongly resembles conventional DMRG, we explain it by describing a single iteration step from the viewpoint of DMRG and show that only a few additional ingredients are needed to arrive at our variational algorithm. We only briefly motivate these extra ingredients for the sake of readability, and refer to Appendix A for additional explanations and more rigorous theoretical motivations. As such, the new algorithm can easily be implemented as an extension to an already existing (I)DMRG implementation.

We start by considering a setting familiar from conventional DMRG: a finite homogeneous one dimensional quantum lattice system, where every site corresponds to a d level system. We label the sites by an integer n and thus have a basis $\{|s\rangle_n, s = 1, \dots, d\}$ for the local Hilbert space on site n . The total Hilbert space is spanned by the product basis $|\mathbf{s}\rangle = \bigotimes_n |s\rangle_n$. We assume the dynamics of the system to be governed by a translation invariant Hamiltonian H .

We further consider a variational parameterization of a ground state approximation of the system, for now in terms of a finite size (site dependent) MPS, but we will ultimately be interested in the thermodynamic limit. DMRG proceeds to find the best variational ground state approximation by employing an alternating least squares minimization: It starts from some initial state and successively optimizes each of the individual MPS tensor site by site by solving effective (Hamiltonian) eigenvalue problems, in a sweeping process through the lattice until convergence, where each iteration depends on already optimized tensors from previous iterations (see e.g. Refs 5, 6, 12, 21, and 23).

We are now however interested in the thermodynamic limit $n \in \mathbb{Z}$ (but will ignore the technical complications involving a rigorous definition of a Hilbert space in that limit). In that case the MPS ground state approximation will be given in terms of a translation invariant uniform MPS, described by a single MPS tensor (or a unit cell of N tensors), repeated on all sites. Two immediate difficulties arise: Firstly, conventional DMRG updates the variational state site by site, thus breaking translation invariance. Secondly, the effective Hamiltonian for a single-site optimization has to be constructed from an infinite environment.

After briefly introducing the variational class of uniform MPS and introducing necessary notation and conventions (for further details see Sec. A 2), we describe how the new algorithm modifies DMRG accordingly to exactly account for these two issues in order to arrive at a variational ground state algorithm directly formulated in the thermodynamic limit.

A. Uniform MPS

A uniform MPS (uMPS) of bond dimension D defined on an infinite translation invariant lattice is parameterized by a single collection of d matrices $A^s \in \mathbb{C}^{D \times D}$ for $s = 1, \dots, d$. The overall translation invariant variational state is then given by

$$|\Psi(A)\rangle = \sum_{\mathbf{s}} \left(\dots A^{s_{n-1}} A^{s_n} A^{s_{n+1}} \dots \right) |\mathbf{s}\rangle \quad (1)$$

and can be represented diagrammatically as

$$|\Psi(A)\rangle = \dots \text{---} \boxed{A} \text{---} \boxed{A} \text{---} \boxed{A} \text{---} \boxed{A} \text{---} \dots$$

Exploiting the invariance of (1) under local gauge transformations $A^s \rightarrow X A^s X^{-1}$, with $X \in \mathbb{C}^{D \times D}$ invertible, the state can be cast into certain favorable representations, among them the *left and right canonical representation*

$$\sum_s A_L^{s\dagger} A_L^s = \mathbb{1} \quad \sum_s A_L^s R A_L^{s\dagger} = R \quad (2a)$$

$$\sum_s A_R^s A_R^{s\dagger} = \mathbb{1} \quad \sum_s A_R^{s\dagger} L A_R^s = L, \quad (2b)$$

or diagrammatically

Here L and R correspond to the left and right reduced density matrices of a bipartition of the state respectively. We henceforth refer to A_L (A_R) as a left (right) isometric tensor, or just a left (right) isometry.

Defining the left and right transfer matrices

$$T_{L/R} = \sum_s \bar{A}_{L/R}^s \otimes A_{L/R}^s \quad (3)$$

and using the notation $\langle x|$ and $|x\rangle$ for vectorizations of a $D \times D$ matrix x in the D^2 dimensional “double layer” virtual space the transfer matrices act upon, the gauge conditions (2) are equivalent to

$$\langle \mathbb{1}| T_L = \langle \mathbb{1}| \quad T_L |R\rangle = |R\rangle \quad (4a)$$

$$T_R |\mathbb{1}\rangle = |\mathbb{1}\rangle \quad \langle L| T_R = \langle L|, \quad (4b)$$

i.e. $\mathbb{1}$ and R are the left and right dominant eigenvectors of T_L , while L and $\mathbb{1}$ are the left and right dominant eigenvectors of T_R .

In the case of nearest neighbor interactions, the action of H_{A_C} onto A_C splits up into four individual contributions, which follow from the decomposition $|\Psi\rangle = \sum_{\alpha,\beta,s} A_{C,(\alpha,\beta)}^s |\Psi_L^\alpha\rangle |s\rangle |\Psi_R^\beta\rangle$ (left block containing sites $n < 0$, center site $n = 0$, and right block containing sites $n > 0$). The action of H_{A_C} onto A_C is given by

$$A_C^s = \sum_{tk\ell} h_{k\ell}^{ts} A_L^t A_C^k A_C^\ell + h_{k\ell}^{st} A_C^k A_R^\ell A_R^t + H_L A_C^s + A_C^s H_R \quad (11)$$

where the first two terms correspond to the Hamiltonian terms $h_{-1,0}$ and $h_{0,1}$ coupling the center site to the left and right block, respectively, and H_L and H_R sum up the contributions of all the Hamiltonian terms $h_{n,n+1}$ acting strictly to the left and to the right of the center site.

The environments H_L and H_R are usually constructed iteratively while sweeping through the (finite) lattice in conventional DMRG, or grown successively in every iteration in IDMRG. In the thermodynamic limit, these terms consist of a diverging number of individual local interaction contributions $h_{n,n+1}$, and care needs to be taken in their construction.

Indeed, the k^{th} contribution to $(H_L|$ comes from the Hamiltonian term $h_{-k-1,-k}$ and is given by $(h_L|[T_L]^{k-1}$. Likewise, $[T_R]^{k-1}|h_R)$ is the k^{th} contribution to $|H_R)$ stemming from $h_{k,k+1}$. Here, we have used the definitions

$$\begin{aligned} h_L &= \sum_{stk\ell} h_{k\ell}^{st} A_L^t A_L^s A_C^k A_C^\ell \\ h_R &= \sum_{stk\ell} h_{k\ell}^{st} A_R^k A_R^\ell A_R^t A_R^s \end{aligned} \quad (12)$$

or diagrammatically

Summing up all such local contributions gives rise to infinite geometric sums of the transfer matrices $T_{L/R}$

$$(H_L| = (h_L| \sum_{k=0}^{\infty} [T_L]^k \quad |H_R) = \sum_{k=0}^{\infty} [T_R]^k |h_R), \quad (13)$$

where $(H_L|$ can be presented diagrammatically as

and likewise for $|H_R)$.

The transfer matrix T_L has a dominant eigenvalue of magnitude one, with corresponding left and right eigenvectors $(\mathbb{1}|$ and $|R)$. The projection $(h_L|[T_L]^k|R) = (h_L|R)$ is the energy density expectation value $e = \langle \Psi | h_{-k-1,-k} | \Psi \rangle$ and is independent of k . Subtracting the energy $\tilde{h} = h - e\mathbb{1}$ from the Hamiltonian, we can write $(h_L| = (\tilde{h}_L| + e\mathbb{1}|$. The second term is exactly proportional to the left eigenvector of eigenvalue 1 and therefore gives rise to a diverging contribution in the geometric sum, corresponding to the total energy of the left half infinite block. Since this contribution acts as the identity in the effective Hamiltonian H_{A_C} [Eq. (11)], we can however safely discard this diverging contribution without changing the eigenvectors of H_{A_C} . This corresponds to an overall energy shift of the left half infinite block such that $(H_L|R) = 0$. For the remaining part $(\tilde{h}_L|$ the geometric sum converges. With $|\tilde{h}_R) = |h_R) - e|\mathbb{1}$ the same comments apply to the construction of $|H_R)$.

We can evaluate H_L and H_R recursively as

$$\begin{aligned} (H_L^{[n+1]}| &= (H_L^{[n]}|[T_L + (\tilde{h}_L| \\ |H_R^{[n+1]}) &= T_R|H_R^{[n]} + |\tilde{h}_R) \end{aligned} \quad (14)$$

with initialization $(H_L^{[0]}| = (\tilde{h}_L|$ and $|H_R^{[0]}) = |\tilde{h}_R)$. We can repeat these recursions until e.g. $\|H_{L/R}^{[n+1]} - H_{L/R}^{[n]}\|$ drops below some desired accuracy ϵ_S . This strategy is conceptually simple and closely resembles the successive construction of the environments in the context of (1)DMRG, but is not very efficient, as its performance is comparable to that of a power method.

Algorithm 1 Explicit terms of effective Hamiltonians with nearest neighbor interactions and their application

Input: two-site Hamiltonian h , current uMPS tensors A_L, A_R in left and right gauge, left dominant eigenvector (L) of T_R , right dominant eigenvector (R) of T_L , desired precision ϵ_S for terms involving infinite geometric sums

Output: Explicit terms of effective Hamiltonians H_{AC} and H_C , updated A'_C and C'

- 1: **function** HEFFTERMS($H = h, A_L, A_R, L, R, \epsilon_S$) ▷ Calculates explicit terms of effective Hamiltonians
- 2: Calculate h_L and h_R from (12)
- 3: Calculate H_L and H_R by iteratively solving (14) or (preferably) (15), to precision ϵ_S
- 4: $H_{AC} \leftarrow \{h, A_L, A_R, H_L, H_R\}$
- 5: $H_C \leftarrow \{h, A_L, A_R, H_L, H_R\}$
- 6: **return** H_{AC}, H_C
- 7: **end function**
- 8: **function** APPLYHAC(A_C, H_{AC}) ▷ Terms of H_{AC} from HEFFTERMS($H, A_L, A_R, L, R, \epsilon_S$)
- 9: Calculate updated A'_C from (11)
- 10: **return** A'_C
- 11: **end function**
- 12: **function** APPLYHC(C, H_C) ▷ Terms of H_C from HEFFTERMS($H, A_L, A_R, L, R, \epsilon_S$)
- 13: Calculate updated C' from (16)
- 14: **return** C'
- 15: **end function**

Table I. Pseudocode for obtaining the explicit terms of the effective Hamiltonians H_{AC} and H_C for systems with nearest neighbor interactions and their applications onto a state.

A more efficient approach is to formally perform the geometric sums in (13) explicitly, and to iteratively solve the resulting two systems of equations

$$\begin{aligned} (H_L)[\mathbb{1} - T_L + |R\rangle\langle\mathbb{1}|] &= (h_L| - (h_L|R)\langle\mathbb{1}| \\ [\mathbb{1} - T_R + |\mathbb{1}\rangle\langle L|]H_R &= |h_R\rangle - |\mathbb{1}\rangle(L|h_R) \end{aligned} \quad (15)$$

for ($H_L|$ and $|H_R\rangle$) to precision ϵ_S , as explained in Appendix D.

So far, we have discussed the action of H_{AC} . The action of H_C onto C follows simply from (11) by projecting onto A_L or A_R . Using the defining property of H_L or H_R , the result simplifies to

$$C' = \sum_{stkl} h_{kl}^{st} A_L^s \dagger A_L^k C A_R^\ell A_R^t \dagger + H_L C + C H_R. \quad (16)$$

The first two terms of (11) can be applied in $\mathcal{O}(d^4 D^3)$ operations²⁹, and the last two terms in $\mathcal{O}(d D^3)$ operations. For (16) the first term can be applied in $\mathcal{O}(d^4 D^3)$ operations, and the last two terms in $\mathcal{O}(D^3)$ operations. To generate the necessary terms for (11) and (16) we have to iteratively evaluate two infinite geometric sums, involving $\mathcal{O}(D^3)$ operations (when iteratively solving (15) the solutions from the previous iteration can be used as starting vectors to speed up convergence). A pseudocode summary for obtaining the necessary explicit terms of H_{AC} and H_C and their applications onto a state is presented in Table I.

C. Updating the state

In DMRG, we would update the state by replacing A_C with the lowest eigenvector \tilde{A}_C of H_{AC} and then shift the center site to the right by computing an orthogonal factorization $\tilde{A}_C^s = \tilde{A}_L^s \tilde{C}_R$, or to the left by computing $\tilde{A}_C^s = \tilde{C}_L \tilde{A}_R^s$. As such, the state gets updated by only replacing the current site with \tilde{A}_L^s or \tilde{A}_R^s , while leaving all other sites untouched. However, applying this scheme in our setting would immediately destroy translation invariance after a single step.

We want to construct an alternative scheme that applies global updates in order to preserve translation invariance at any time. Global updates can most easily

be applied with an explicit uniform parameterization in terms of a single tensor A . On the other hand, DMRG experience teaches us that the stability is greatly enhanced when applying updates at the level of A_C and C , which are isometrically related to the full state.

We therefore calculate the lowest eigenvector \tilde{A}_C of H_{A_C} like in DMRG, but additionally also the lowest eigenvector \tilde{C} of H_C . We then globally update the state by finding new \tilde{A}_L and \tilde{A}_R as the left and right isometric tensors that minimize $\sum_s \|\tilde{A}_L^s \tilde{C} - \tilde{A}_C^s\|^2$ and $\sum_s \|\tilde{C} \tilde{A}_R^s - \tilde{A}_C^s\|^2$ respectively. These minimization problems can be solved directly (not iteratively) and without inverting \tilde{C} (see below). As shown in Appendix A, at the variational optimum the values of these objective functions go to zero, and current A_C and C are the lowest eigenvectors of H_{A_C} and H_C respectively.

For the remainder of this section we omit tildes and use the following matricizations of the 3-index tensors

$$\begin{aligned} \mathcal{A}_{L,(s\alpha,\beta)} &= A_{L,(\alpha,\beta)}^s \\ \mathcal{A}_{R,(\alpha,s\beta)} &= A_{R,(\alpha,\beta)}^s \\ \mathcal{A}_{C,(s\alpha,\beta)}^{[\ell]} &= \mathcal{A}_{C,(\alpha,s\beta)}^{[r]} = A_{C,(\alpha,\beta)}^s. \end{aligned} \quad (17)$$

We thus want to extract updated A_L and A_R from updated A_C and C by solving

$$\epsilon_L = \min_{\mathcal{A}_L \mathcal{A}_L^\dagger = \mathbb{1}} \|\mathcal{A}_C^{[\ell]} - \mathcal{A}_L C\|_2 \quad (18a)$$

$$\epsilon_R = \min_{\mathcal{A}_R \mathcal{A}_R^\dagger = \mathbb{1}} \|\mathcal{A}_C^{[r]} - C \mathcal{A}_R\|_2. \quad (18b)$$

In exact arithmetic, the solution of these minimization problems is known, namely \mathcal{A}_L will be the isometry in the polar decomposition of $\mathcal{A}_C^{[\ell]} C^\dagger$ (and similar for \mathcal{A}_R , see Thm. IX.7.2 in Ref. 30). Computing the singular value decompositions (SVD)

$$\mathcal{A}_C^{[\ell]} C^\dagger = U^{[\ell]} \Sigma^{[\ell]} V^{[\ell]\dagger} \quad C^\dagger \mathcal{A}_C^{[r]} = U^{[r]} \Sigma^{[r]} V^{[r]\dagger} \quad (19)$$

we thus obtain

$$\mathcal{A}_L = U^{[\ell]} V^{[\ell]\dagger} \quad \mathcal{A}_R = U^{[r]} V^{[r]\dagger}. \quad (20)$$

Notice that close to (or at) an exact solution $A_C^s = A_L^s C = C A_R^s$, the singular values contained in $\Sigma^{[\ell/r]}$ are the square of the singular values of C , and might well fall below machine precision. Consequently, in finite precision arithmetic, corresponding singular vectors will not be accurately computed.

An alternative that has proven to be robust and still close to optimal is given by directly using the following left and right polar decompositions

$$\mathcal{A}_C^{[\ell]} = U_{A_C}^{[\ell]} P_{A_C}^{[\ell]} \quad C = U_C^{[\ell]} P_C^{[\ell]} \quad (21a)$$

$$\mathcal{A}_C^{[r]} = P_{A_C}^{[r]} U_{A_C}^{[r]} \quad C = P_C^{[r]} U_C^{[r]} \quad (21b)$$

to obtain

$$\mathcal{A}_L = U_{A_C}^{[\ell]} U_C^{[\ell]\dagger} \quad \mathcal{A}_R = U_C^{[r]\dagger} U_{A_C}^{[r]}, \quad (22)$$

where matrices P are hermitian and positive. Alternative isometric decompositions might be considered in Eq. (21), though it is important that they are unique (e.g. QR with positive diagonal in R) in order to have $P_{A_C}^{[\ell/r]} \approx P_C^{[\ell/r]}$ close to convergence.

D. The Algorithm: VUMPS

We are now ready to formulate our *variational uniform MPS* (VUMPS) algorithm. As shown in Appendix A, a variational minimum (vanishing energy gradient) in the manifold of uMPS is characterized by tensors A_L , C and A_R satisfying the conditions

$$\mathbf{H}_{A_C} \mathbf{A}_C = E_{A_C} \mathbf{A}_C \quad (23a)$$

$$\mathbf{H}_C \mathbf{C} = E_C \mathbf{C} \quad (23b)$$

$$A_C^s = A_L^s C = C A_R^s. \quad (23c)$$

Here bold symbols denote vectorizations of the MPS tensors and matricizations of the effective Hamiltonians, and E_{A_C} and E_C are the lowest eigenvalues of the effective Hamiltonians.³¹

When iterating the steps outlined in the previous subsections, convergence is obtained when these conditions are satisfied. In particular, starting with a properly orthogonalized initial trial state $|\Psi(A)\rangle$ of some bond dimension D , we begin by solving the two eigenvalue problems for the effective Hamiltonians H_{A_C} and H_C . Since we are still far from the fixed point, the resulting lowest energy states \tilde{A}_C and \tilde{C} will in general not satisfy the gauge condition (23c) together with current $A_{L/R}$.

Following the procedure of the previous section we can however find optimal approximations \tilde{A}_L^s and \tilde{A}_R^s for (23c) to arrive at an updated uMPS. Conversely, \tilde{A}_C and \tilde{C} will not be the correct lowest energy eigenstates of the new effective Hamiltonians $H_{\tilde{A}_C}$ and $H_{\tilde{C}}$ generated from $\tilde{A}_{L/R}$. We then use the updated state and reiterate this process of alternately solving the effective eigenvalue problems, and finding optimal approximations for A_L and A_R to update the state. For a pseudocode summary of this algorithm, see Table II.

We now elaborate on the various steps in the VUMPS algorithm. Firstly, extracting new $\tilde{A}_{L/R}$ from updated \tilde{A}_C and \tilde{C} can be done using the theoretically optimal (but numerically often inaccurate) Eq. (20) or the more robust Eq. (22), depending on the magnitude of the smallest singular value in \tilde{C} . As a good uMPS approximation will always involve small singular values, Eq. (22) is preferable most of the time, except maybe during the first few iterations.

The maximum of the error quantities (18)

$$\epsilon_{\text{prec}} = \max(\epsilon_L, \epsilon_R) \quad (24)$$

provides an error measure for the fixed point condition in Eq. (23c) and is used as a global convergence criterion.

Algorithm 2 variational uMPS algorithm for single-site unit cells

Input: Hamiltonian H , initial uMPS A_L, A_R, C , convergence threshold ϵ

Output: uMPS approximation A_L, A_R, C of ground state of H , fulfilling fixed point relations (23a), (23b) and (23c) up to precision ϵ

```

1: procedure VUMPS( $H, A_L, A_R, C, \epsilon$ )
2:   initialize current precision  $\epsilon_{\text{prec}} > \epsilon$ 
3:   while  $\epsilon_{\text{prec}} > \epsilon$  do
4:     (optional) Dynamically adjust bond dimension following Appendix B
5:     Calculate explicit terms of effective Hamiltonians  $H_{A_C}, H_C \leftarrow \text{HEFFTERMS}(H, A_L, A_R, L, R, \epsilon_S \leq \epsilon_{\text{prec}})$  from
       Algorithm 1, 5 or 6
6:     Calculate ground state  $\tilde{A}_C$  of effective Hamiltonian  $H_{A_C}$  to precision  $\epsilon_H < \epsilon_{\text{prec}}$  using an iterative eigensolver,
       calling APPLYHAC( $A_C, H_{A_C}$ ) from Algorithm 1, 5 or 6
7:     Calculate ground state  $\tilde{C}$  of effective Hamiltonian  $H_C$  to precision  $\epsilon_H < \epsilon_{\text{prec}}$  using an iterative eigensolver,
       calling APPLYHC( $C, H_C$ ) from Algorithm 1, 5 or 6
8:     Calculate new  $\tilde{A}_L$  and  $\tilde{A}_R$  from  $\tilde{A}_C$  and  $\tilde{C}$  using (20) or (22), depending on singular values of  $\tilde{C}$ 
9:     Evaluate new  $\epsilon_L$  and  $\epsilon_R$  from (18)
10:    (optional) Calculate current expectation values
11:    Set  $\epsilon_{\text{prec}} \leftarrow \max(\epsilon_L, \epsilon_R)$  and replace  $A_L \leftarrow \tilde{A}_L, A_R \leftarrow \tilde{A}_R$  and  $C \leftarrow \tilde{C}$ 
12:  end while
13:  return  $A_L, A_R, C$ 
14: end procedure

```

Table II. Pseudocode of the VUMPS algorithm described in Sec. IID. Terms within step 5 involving the evaluation of infinite geometric sums usually require the left dominant eigenvector L of T_R and the right dominant eigenvector R of T_L , for which $L = C^\dagger C$ and $R = CC^\dagger$ with current C are a good enough approximation to current precision ϵ_{prec} (see main text). Notice that this algorithm is free of any possibly ill-conditioned inverses and therefore has no convergence issues in the presence of small Schmidt values. It also does not require expensive reorthogonalizations of the state at intermediate iterations.

It measures the precision of the current uMPS ground state approximation. Within every iteration, we use iterative methods (e.g. some variation of Lanczos) to find the eigenvectors \tilde{A}_C and \tilde{C} of the Hermitian operators H_{A_C} and H_C . As the goal is to drive the state towards the fixed point relations in Eqs. (23a) and (23b), it is not necessary to solve these eigenvalue problems to full machine precision. Rather, it is sufficient to use a tolerance ϵ_H chosen relative to ϵ_{prec} .³² A value of ϵ_H of the order of $\epsilon_{\text{prec}}/100$ has proven to work well in practice. It is also worthwhile to use tensors from the previous iteration as initial guess for the iterative solvers to speed up convergence.

As the main part of the algorithm works at fixed bond dimension (i.e. it is a single-site scheme in DMRG terminology), one might choose to increase the bond dimension D before starting a new iteration. We have developed a subspace expansion technique that works directly in the thermodynamic limit and is explained in Appendix B.

While the true comparison of this algorithm with IDMRG^{5,24} and ITEBD²⁶ will take place in Sec. III by gathering actual numerical simulation results, we can already compare the theoretical properties of these algorithms. Neither IDMRG or ITEBD is truly solving the variational problem in the sense of directly trying to satisfy the fixed point conditions Eqs. (23). IDMRG closely resembles regular DMRG on a successively growing lattice, as it inserts and optimizes over new tensors in the center of the lattice in each step. Tensors from previous steps are not updated, as this would render the cost

prohibitive. When this approach converges, the resulting fixed point tensors in the center can be assumed to specify the unit cell of an infinite MPS. VUMPS has the immediate advantage that i) it directly works in the thermodynamic limit at all iterations and ii) it completely replaces the entire state after every iteration, thus moving faster through the variational manifold. In contrast, IDMRG keeps memory of earlier iterations and cannot guarantee a monotonically decreasing energy that converges to an optimum associated with a translation invariant MPS in which the effects of the boundary have completely disappeared. The advantages of VUMPS come with a greater computational cost per iteration, as two eigenvalue problems (for A_C and for C) and – in the case of nearest neighbor interactions – two linear systems (for H_L and H_R) have to be solved. IDMRG only solves a single eigenvalue problem and builds H_L and H_R step by step in every iteration. The latter approach is analogous to a power method for eigenvalue problems and, while very cheap, is expected to require many iteration steps to converge, especially for systems with large correlation lengths (e.g. close to criticality).

ITEBD²⁶ is based on evolving an initial state in imaginary time by using a Trotter decomposition of the evolution operator. Like VUMPS, ITEBD works in the thermodynamic limit at any intermediate step, typically with a unit cell that depends on how the Hamiltonian was split into local terms in order to apply the Trotter decomposition. Furthermore, as every application of the evolution operator increases the virtual dimension of the

MPS, truncation steps are required to restore the original (or any suitable) value of the bond dimension. While VUMPS can take big steps through the variational space, time steps in ITEBD have to be chosen sufficiently small (especially in the final steps of the algorithm) to eliminate the Trotter error, which negatively affects the rate of convergence. (Ref. 33 however proposes a scheme to effectively obtain a larger time step). Furthermore, the Trotter splitting essentially limits the applicability of ITEBD to short-range interactions and dictates the size of the unit cell of the resulting MPS, e.g. in the most common case of nearest neighbor interactions a two-site unit cell is obtained. (The approach of Ref. 34 to obtain a translation invariant MPS is restricted to certain Hamiltonians, but see Ref. 35 for an alternative proposal that can in fact also deal with long range interactions.)

Finally, we can also compare VUMPS to the more recent *time dependent variational principle* (TDVP),²⁷ which was implemented as an alternative approach to simulate real and imaginary time evolution within the manifold of MPS by projecting the evolution direction onto the MPS tangent space. This approach can be applied to translation invariant MPS, independent of the type of Hamiltonian. When used to evolve in imaginary time, it can be identified as a covariant formulation of a gradient descent method, in that it evolves the state in the direction of the gradient of the energy functional, preconditioned with the metric of the manifold. As such, the energy decreases monotonically and at convergence, an exact (local) minimum is obtained, as characterized by the vanishing gradient. However, in its original formulation, TDVP was not formulated in a center site form and was therefore unstable and restricted to small time steps. For finite systems, a different formulation of the TDVP algorithm was provided in Ref. 28, which allows for taking the limit of the imaginary time step to infinite, and then becomes provably equivalent to the single-site DMRG algorithm. VUMPS can be motivated from these developments, as explained in Appendix A.

We conclude this section by elaborating on how to incorporate symmetries in the algorithm. The construction of uMPS that is explicitly invariant under onsite unitary symmetries is equivalent to (I)DMRG^{12,23} and (I)TEBD^{36,37}, and it is immediately clear that the various steps in VUMPS have a corresponding covariant formulation. The same comments apply to time reversal symmetry, in which case everything can be implemented in real arithmetic, or to reflection symmetry, in which case C and A_C^s will be symmetric matrices and $A_R^s = A_L^{sT}$ (which implies that H_L and H_R are also related). In all of these cases, the computational cost is reduced. However, explicitly imposing the symmetry in the MPS requires caution, as the physical system might have spontaneous symmetry breaking, or – more subtly – might be in a symmetry protected topological phase where the symmetries cannot be represented trivially on the MPS tensor.

In the case of spontaneous symmetry breaking, MPS

algorithms tend to converge to maximally symmetry broken states for which the entanglement is minimal. This is also the case for VUMPS. One can control which state the algorithm converges to by suitably biasing the initial state or by adding small perturbation terms to the Hamiltonian which explicitly break the symmetry, and which are switched off after a few iterations.

Explicit conservation of translation symmetry was the very first requirement in the construction of VUMPS. In the case of spontaneous breaking of translation symmetry down to N -site translation symmetry (as e.g. in the case of a state with antiferromagnetic order), enforcing one-site translation symmetry would result in a (non-injective) equal weight superposition of all symmetry broken uMPS ground state approximations. In order to reach an optimal accuracy with a given bond dimension, such a superposition of N states is however undesirable, as the effective bond dimension is reduced to D/N . In the case where this situation cannot be amended by a simple unitary transformation that restores one-site translation symmetry (such as e.g. flipping every second spin in the case of an antiferromagnet), it is preferable to choose an MPS ansatz with a N -site unit cell, such that the state can spontaneously break translation symmetry. The generalization of the algorithm to multi-site unit cells is described in the next section.

E. Multi Site Unit Cell Implementations

We now generalize the VUMPS algorithm of the previous section for one-site translation invariant uMPS to the setting of translation invariance over N sites. Such a uMPS ansatz is then parameterized by N independent tensors $A(k)^s \in \mathbb{C}^{D \times d \times D}$, $k = 1, \dots, N$, which define the unit cell tensor

$$\mathbb{A}^{\mathfrak{s}} = A(1)^{s_1} \dots A(N)^{s_N}, \quad (25)$$

where $\mathfrak{s} = (s_1, \dots, s_N)$ is a combined index. We can then write the variational state as

$$|\Psi(\mathbb{A})\rangle = \sum_{\mathfrak{s}} (\dots \mathbb{A}^{\mathfrak{s}_{n-1}} \mathbb{A}^{\mathfrak{s}_n} \mathbb{A}^{\mathfrak{s}_{n+1}} \dots) |\mathfrak{s}\rangle$$

and the left and right orthonormal forms are given by the relations

$$\begin{aligned} \sum_s A(k)_L^s \dagger A(k)_L^s &= \mathbb{1} \\ \sum_s A(k)_L^s R(k) A(k)_L^{s\dagger} &= R(k-1) \end{aligned} \quad (26a)$$

and

$$\begin{aligned} \sum_s A(k)_R^s A(k)_R^{s\dagger} &= \mathbb{1} \\ \sum_s A(k)_R^{s\dagger} L(k-1) A(k)_R^s &= L(k), \end{aligned} \quad (26b)$$

where it is understood that $N + 1 \equiv 1$ and $0 \equiv N$.

Defining the bond matrices $C(k)$ as the gauge transformation that relates left and right canonical form via $C(k-1)A(k)_R^s = A(k)_L^s C(k)$, we have $R(k) = C(k)C(k)^\dagger$ and $L(k) = C(k)^\dagger C(k)$. We can then also cast $|\Psi(\mathbb{A})\rangle$ in a mixed canonical form similar to (5a) with the center site tensor given by $A(k)_C^s = A(k)_L^s C(k) = C(k-1)A(k)_R^s$.

The variational minimum within this set of states is characterized by the following $3N$ fixed point relations

$$\mathbf{H}_{A(k)_C} \mathbf{A}(k)_C = E_{A(k)_C} \mathbf{A}(k)_C \quad (27a)$$

$$\mathbf{H}_{C(k)} \mathbf{C}(k) = E_{C(k)} \mathbf{C}(k) \quad (27b)$$

$$A(k)_C^s = A(k)_L^s C(k) = C(k-1)A(k)_R^s. \quad (27c)$$

Notice that due to (27c), the relations for different k are connected. There are several possible strategies for constructing algorithms which obtain states satisfying these conditions.

In the following we present two approaches which have shown good performance and stable convergence, which we shall term the ‘‘sequential’’ and ‘‘parallel’’ methods. But let us first elaborate on computing effective Hamiltonians for multi-site unit cells, which works similarly in both methods. We again restrict to the case of nearest neighbor interactions, such that the effective Hamiltonians are constructed similar as in Sec. II B. To construct e.g. the left block Hamiltonian H_L , we first collect all local contributions from a single unit cell in h_L , before performing the geometric series of the transfer matrix, which now mediates a translation over an entire unit cell.

1. Sequential Algorithm

The sequential algorithm is inspired by finite size DMRG, in that we sweep through the unit cell, successively optimizing one tensor at a time while keeping tensors on other sites fixed. Notice that at site k we however need *two* updated bond matrices $\tilde{C}(k)_L = \tilde{C}(k-1)$ and $\tilde{C}(k)_R = \tilde{C}(k)$, in order to calculate updated $\tilde{A}(k)_{L/R}^s$ from $\tilde{A}(k)_C^s \approx \tilde{A}(k)_L^s \tilde{C}(k)_R \approx \tilde{C}(k)_L \tilde{A}(k)_R^s$. We thus have to amend steps 5, 6 and 7 of the single-site algorithm in Table II by constructing and solving for *two* effective Hamiltonians $H_{C(k-1)}$ and $H_{C(k)}$ instead of a single one. The newly optimized tensors then get replaced in *all* unit cells of the infinite lattice, and contributions to the effective Hamiltonians have to be recalculated accordingly, before moving on to the next site. For a pseudocode summary see Algorithm 3 in Table III.

One could now try to argue, that e.g. in a left to right sweep it is enough at site k to calculate updated $\tilde{A}(k)_C$ and $\tilde{C}(k)_R = \tilde{C}(k)$ only, and to use $\tilde{C}(k-1)_R$ from the previous step at site $k-1$ as $\tilde{C}(k)_L$ for calculating $\tilde{A}(k)_R$. This approach however fails, as the effective Hamiltonian used for calculating $\tilde{A}(k)_C$ already contains updated $\tilde{A}(k-1)_{L/R}$, while the effective Hamiltonian used for calculating $\tilde{C}(k-1)_R$ does not, and we cannot

determine $\tilde{A}(k)_R$ from $\tilde{A}(k)_C$ and $\tilde{C}(k-1)_R$. Rather, $\tilde{C}(k)_L$ has to be recalculated using an *updated* effective Hamiltonian, which exactly leads to the sequential Algorithm 3.

There is an additional subtlety that needs to be considered, in order for all tensors to fulfill the gauge constraints (27c) to current precision. Bond matrices $\tilde{C}(k)$ are calculated as lowest energy eigenvectors of effective Hamiltonians $H_{C(k)}$ and are therefore only determined up to a phase. Consider $C(k)$ defined between sites k and $k+1$. At step k it is updated as $\tilde{C}(k)_R$ and used to calculate $\tilde{A}(k)_L^s$. In the next step $k+1$ however it is recalculated as $\tilde{C}(k+1)_L$ (with an updated effective Hamiltonian) and used to determine $\tilde{A}(k+1)_R^s$. At the fixed point we should then have $\tilde{C}(k)_R = \tilde{C}(k+1)_L = C(k)$, but this is only true if there is no phase ambiguity, which would also consequently lead to a phase mismatch between $\tilde{A}(k)_L$ and $\tilde{C}(k)$ after step $k+1$. This issue does not pose a problem for algorithm convergence (during calculations, matrices $C(k)$ always appear as products of the form $C(k)^\dagger C(k)$ or $C(k)C(k)^\dagger$ and mismatching phases thus cancel out), but can be easily circumvented by employing a phase convention when calculating updated $\tilde{C}(k)$.

2. Parallel Algorithm

In the parallel approach, we choose to update an entire unit cell at once, using effective Hamiltonians generated from the same current state. To that end, we first generate all terms necessary for all $H_{A(k)_C}$ and $H_{C(k)}$. For the case of nearest neighbor interactions, the contributions H_L and H_R to the left and right environment outside the unit cell can be shared, so that the corresponding geometric sum only needs to be computed once, and contributions inside the unit cell are obtained through successive applications of transfer matrices.

Next, we simultaneously and independently solve for the ground states $\tilde{A}(k)_C$ and $\tilde{C}(k)$ of all $2N$ effective Hamiltonians at once. Once these are obtained we again simultaneously and independently determine all updated $\tilde{A}(k)_L$ and $\tilde{A}(k)_R$, concluding one iteration for updating the entire unit cell. For a pseudocode summary see Algorithm 4 in Table III.

3. Juxtaposition of Both Approaches

Several comments on the two presented algorithms are in order. First, the parallel algorithm requires substantially less computational effort, since the construction of the different effective Hamiltonians $H_{A(k)_C}$ can recycle the calculation of the infinite geometric sum. Therefore, updating an entire unit cell only requires to evaluate two infinite geometric sums and $2N$ effective eigenvalue problems. In the sequential algorithm, updating the environment after every tensor update requires to reevaluate

Algorithm 3 *sequential* variational uMPS algorithm for multi-site unit cells

Input: Hamiltonian H , initial uMPS $\{A_L\}, \{A_R\}, \{C\}$ of an N -site unit cell, convergence threshold ϵ
Output: uMPS approximation $\{A_L\}, \{A_R\}, \{C\}$ of ground state of H , fulfilling fixed point relations (27a), (27b) and (27c) up to precision ϵ .

- 1: **procedure** VUMPSMULTISEQUENTIAL($H, \{A_L\}, \{A_R\}, \{C\}, \epsilon$)
- 2: initialize current precision $\epsilon_{\text{prec}} > \epsilon$
- 3: **while** $\epsilon_{\text{prec}} > \epsilon$ **do**
- 4: **for** $n = 1, \dots, N$ **do**
- 5: (optional) Dynamically adjust bond dimension following Appendix B
- 6: Calculate explicit terms of effective Hamiltonians from a multi-site version
 $H_{A(n)C}, H_{C(n-1)}, H_{C(n)} \leftarrow \text{HEFFTERMSMULTI}(H, \{A_L\}, \{A_R\}, \{L\}, \{R\}, \epsilon_S \leq \epsilon_{\text{prec}})$ of Algorithm 1, 5 or 6
- 7: Calculate ground state \tilde{A}_C of effective Hamiltonian $H_{A(n)C}$ to precision $\epsilon_H < \epsilon_{\text{prec}}$ using an iterative eigensolver, calling APPLYHAC($C, H_{A(n)C}$) from Algorithm 1, 5 or 6
- 8: Calculate ground state \tilde{C}_L of effective Hamiltonian $H_{C(n-1)}$ to precision $\epsilon_H < \epsilon_{\text{prec}}$ using an iterative eigensolver, calling APPLYHC($C, H_{C(n-1)}$) from Algorithm 1, 5 or 6 ▷ To ensure gauge consistency, employ a phase convention for \tilde{C}_L
- 9: Calculate ground state \tilde{C}_R of effective Hamiltonian $H_{C(n)}$ to precision $\epsilon_H < \epsilon_{\text{prec}}$ using an iterative eigensolver, calling APPLYHC($C, H_{C(n)}$) from Algorithm 1, 5 or 6 ▷ To ensure gauge consistency, employ a phase convention for \tilde{C}_R
- 10: Calculate new \tilde{A}_L from \tilde{A}_C and \tilde{C}_R using (20) or (22), depending on singular values of \tilde{C}_R
- 11: Calculate new \tilde{A}_R from \tilde{A}_C and \tilde{C}_L using (20) or (22), depending on singular values of \tilde{C}_L
- 12: Evaluate new $\epsilon_L(n)$ and $\epsilon_R(n)$ from (18a) and (18b)
- 13: Replace $A(n)_L \leftarrow \tilde{A}_L, A(n)_R \leftarrow \tilde{A}_R, C(n-1) \leftarrow \tilde{C}_L$ and $C(n) \leftarrow \tilde{C}_R$
- 14: **end for**
- 15: Set $\epsilon_{\text{prec}} \leftarrow \max(\{\epsilon_L\}, \{\epsilon_R\})$
- 16: (optional) Calculate current expectation values
- 17: **end while**
- 18: **return** $\{A_L\}, \{A_R\}, \{C\}$
- 19: **end procedure**

Algorithm 4 *parallel* variational uMPS algorithm for multi-site unit cells

Input: Hamiltonian H , initial uMPS $\{A_L\}, \{A_R\}, \{C\}$ of an N -site unit cell, convergence threshold ϵ
Output: uMPS approximation $\{A_L\}, \{A_R\}, \{C\}$ of ground state of H , fulfilling fixed point relations (27a), (27b) and (27c) up to precision ϵ .

- 1: **procedure** VUMPSMULTIPARALLEL($H, \{A_L\}, \{A_R\}, \{C\}, \epsilon$)
- 2: initialize current precision $\epsilon_{\text{prec}} > \epsilon$
- 3: **while** $\epsilon_{\text{prec}} > \epsilon$ **do**
- 4: (optional) Dynamically adjust bond dimension following Appendix B
- 5: **for** $n = 1, \dots, N$ **do**
- 6: Calculate explicit terms of effective Hamiltonians from a multi-site version
 $H_{A(n)C}, H_{C(n)} \leftarrow \text{HEFFTERMSMULTI}(H, \{A_L\}, \{A_R\}, \{L\}, \{R\}, \epsilon_S \leq \epsilon_{\text{prec}})$ of Algorithm 1, 5 or 6
- 7: Calculate ground state $\tilde{A}(n)_C$ of effective Hamiltonian $H_{A(n)C}$ to precision $\epsilon_H < \epsilon_{\text{prec}}$ using an iterative eigensolver, calling APPLYHAC($C, H_{A(n)C}$) from Algorithm 1, 5 or 6
- 8: Calculate ground state $\tilde{C}(n)$ of effective Hamiltonian $H_{C(n)}$ to precision $\epsilon_H < \epsilon_{\text{prec}}$ using an iterative eigensolver, calling APPLYHC($C, H_{C(n)}$) from Algorithm 1, 5 or 6
- 9: **end for**
- 10: **for** $n = 1, \dots, N$ **do**
- 11: Calculate new $\tilde{A}(n)_L$ from $\tilde{A}(n)_C$ and $\tilde{C}(n)$ using (20) or (22), depending on singular values of $\tilde{C}(n)$
- 12: Calculate new $\tilde{A}(n)_R$ from $\tilde{A}(n)_C$ and $\tilde{C}(n-1)$ using (20) or (22), depending on singular values of $\tilde{C}(n-1)$
- 13: Evaluate new $\epsilon_L(n)$ and $\epsilon_R(n)$ from (18a) and (18b)
- 14: **end for**
- 15: Replace $\{A_L\} \leftarrow \{\tilde{A}_L\}, \{A_R\} \leftarrow \{\tilde{A}_R\}$ and $\{C\} \leftarrow \{\tilde{C}\}$
- 16: (optional) Calculate current expectation values
- 17: Set $\epsilon_{\text{prec}} \leftarrow \max(\{\epsilon_L\}, \{\epsilon_R\})$
- 18: **end while**
- 19: **return** $\{A_L\}, \{A_R\}, \{C\}$
- 20: **end procedure**

Table III. Pseudocode for the two approaches for a multi-site unit cell implementation described in Sec. III E. Algorithm 3 sweeps through the unit cell and sequentially updates tensors site by site, replacing updated tensors in all unit cells before moving on to the next site. Algorithm 4 updates the entire unit cell at once by independently updating tensors on each site.

the geometric sum, thus leading to $2N$ infinite geometric sums and $3N$ effective eigenvalue problems for updating the complete unit cell. Additionally, the parallel approach offers the possibility of parallelizing the solution of all $2N$ eigenvalue problems in one iteration, while in the sequential approach only 3 eigenvalue problems can be solved in parallel for each site. However, while sweeping through the unit cell in the sequential approach, initial guesses for solving the infinite geometric sums can be generated easily from the previous iterations, and are usually much better than the initial guesses in the parallel algorithm. Equivalently, updated $\tilde{C}(k)$ obtained at site k is a very good initial guess for its recalculation with updated environment on site $k + 1$. Overall, the computational cost for the parallel update is still much cheaper, albeit less than expected.

On the other hand, state convergence in terms of iterations is generally substantially faster in the sequential approach. This seems reasonable, as the optimization on a current site takes into account all previous optimization steps, whereas in the parallel approach, the optimizations on different sites within one iteration are independent of each other. This effect gets amplified with increasing unit cell size N , and the performance of the parallel approach decreases, while the performance of the sequential approach seems more stable against increasing N .

In conclusion, while updating the entire unit cell is computationally cheaper in the parallel approach, the sequential algorithm usually requires a substantially smaller number of iterations due to faster convergence. While there are instances where one approach clearly outperforms the other by far, such cases are rare and strongly depend on initial conditions, and generally both approaches show comparable performance. For comparison benchmark results see Sec. III B 5.

III. TEST CASES AND COMPARISON

In this section we test the performance of the new algorithm on several paradigmatic strongly correlated lattice models in the thermodynamic limit, with nearest neighbor as well as long range interactions. In Sec. III A we introduce and discuss the models under considerations. In Sec. III B we first test the convergence and stability of the single and multi-site implementations of the new algorithm. Lastly, we compare its performance against established conventional MPS methods for ground state search in Sec. III C.

A. Models

As examples for spin chain models with nearest neighbor interactions we study the spin $S = 1/2$ *transverse*

field Ising (TFI) model

$$H_{\text{TFI}} = - \sum_j X_j X_{j+1} - h \sum_j Z_j \quad (28)$$

and the *XXZ* model for general spin S

$$H_{\text{XXZ}} = \sum_j X_j X_{j+1} + Y_j Y_{j+1} + \Delta Z_j Z_{j+1}. \quad (29)$$

Here X , Y and Z are spin S representations of the generators of $SU(2)$. The ground state energies are known exactly for the TFI model,³⁸ and for $S = 1/2$ also for the XXZ model.³⁹ For the $S = 1$ XXZ model we focus on the isotropic antiferromagnetic case $\Delta = 1$ and take the result of Ref. 27 for the ground state energy for $D = 1024$ as quasi-exact result.

As a further example for a system with nearest neighbor interactions we also study the *Fermi Hubbard* model

$$H_{\text{HUB}} = -t \sum_{\sigma,j} c_{\sigma,j} c_{\sigma,j+1}^\dagger - c_{\sigma,j}^\dagger c_{\sigma,j+1} + U \sum_j \left(n_{\uparrow,j} - \frac{1}{2} \right) \left(n_{\downarrow,j} - \frac{1}{2} \right), \quad (30)$$

where $c_{\sigma,j}$, $c_{\sigma,j}^\dagger$ are creation and annihilation operators of electrons of spin σ on site j , $n_{\sigma,j} = c_{\sigma,j}^\dagger c_{\sigma,j}$ and $n_j = n_{\uparrow,j} + n_{\downarrow,j}$ are the particle number operators. Again, the exact ground state energy is known.^{40,41}

Finally, as an example for an exactly solvable model with (algebraically decaying) long range interactions we consider the Haldane-Shastry model^{42,43}

$$H_{\text{HS}} = \sum_j \sum_{n>0} n^{-2} [X_j X_{j+n} + Y_j Y_{j+n} + Z_j Z_{j+n}], \quad (31)$$

where X , Y and Z are again spin $S = 1/2$ representations of the generators of $SU(2)$. In order to efficiently compute the terms of the effective Hamiltonian (see Appendix C 1), we expand the distance function $f(n) = n^{-2}$ in a sum of $K = 20$ exponentials, with maximum residual less than 10^{-6} for a fit over $N = 1000$ sites.

B. Performance Benchmarks

We performed convergence benchmarks for several instances of the models introduced in the previous section, using simple implementations of VUMPS for single or multi-site unit cells presented in Algorithm 2, 3 and 4, without explicitly exploiting any symmetries. Hereto, we consider firstly the error in the variational energy density

$$\Delta e = e - e_{\text{exact}} \quad (32)$$

as a function of the number of iterations. Here e_{exact} is the exact analytic (or quasi-exact numerical) ground state energy density of the model under consideration.

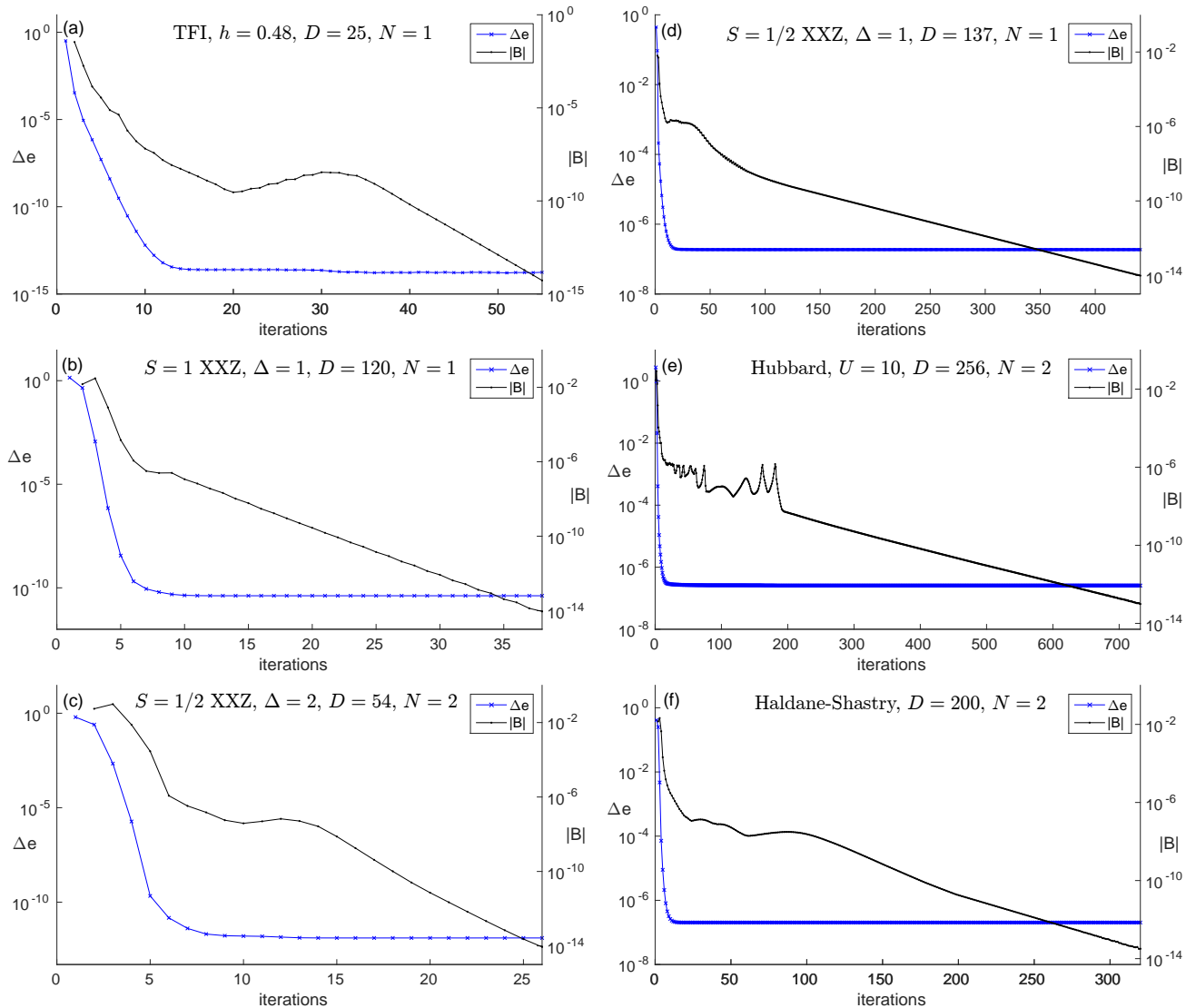


Figure 1. Plot of energy density error Δe and gradient norm $\|B\|$ for VUMPS with a single-site or N -site unit cell: (a) TFI model in the gapped symmetry broken phase at $h = 0.48$ and $D = 25$, (b) gapped isotropic $S = 1$ XXZ antiferromagnet and $D = 120$, (c) gapped $S = 1/2$ XXZ antiferromagnet at $\Delta = 2$ and $D = 54$, (d) critical isotropic $S = 1/2$ XXZ antiferromagnet at $\Delta = 1$ and $D = 137$, (e) critical Fermi Hubbard model at $U = 10$ and $D = 126$, and (f) critical $S = 1/2$ Haldane-Shastry model at $D = 200$. The uMPS ground state approximations of (c), (e) and (f) break translation invariance and have been obtained from Algorithm 3 with a two-site unit cell. Regardless of the criticality of the model, VUMPS converges exponentially fast in gradient norm $\|B\|$. Notice that at the point where the energy has already converged to machine precision, the gradient is still quite far from zero, and the state thus still some distance from the variational optimum.

VUMPS as formulated in Table II has its internal convergence measure used to determine when to stop the iteration loop, as well as to set the tolerance in the iterative solvers used within every single outer iteration. However, as a more objective quantity that measures the distance to the variational minimum, we also compute the norm of the energy gradient; it is an absolute measure of convergence which is independent of any prior iterations, as opposed to relative changes in e.g. the energy or Schmidt spectrum between iterations. We denote this quantity as $\|B\|$, the two-norm of a $D \times d \times D$ tensor B , which can

be worked out to be given by (see Appendix A 3)

$$B^s = A_C'^s - A_L^s C' \quad \text{or} \quad B^s = A_C'^s - C' A_R^s. \quad (33)$$

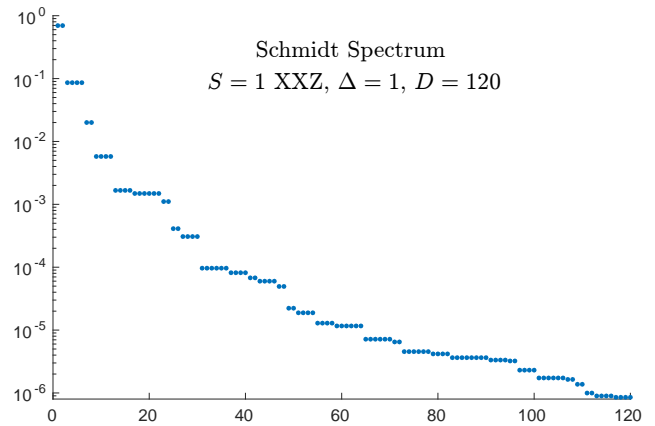
The efficient and accurate computation of the gradient norm is further discussed in Appendix A 4. To obtain the energy gradient $\|B\|$ of an N -site unit cell, it is equivalent to determine the gradients $B(k)$ for each site independently and to calculate the norm of the concatenation of all N gradients.

A well known-property of the variational principle is that the energy expectation value itself converges quadratically faster than the state. When the state has

converged to some accuracy $\|B\|$, the energy density has already converged to precision $\mathcal{O}(\|B\|^2)$ which can therefore be well beyond machine precision. The convergence measure $\|B\|$ does however dictate the convergence of other observables which are not diagonal in the energy eigenbasis. Note, however, that we are here referring to convergence towards the value at the variational optimum, not towards the exact value. The error between the variational optimum and the exact ground state can be quantified using e.g. the energy variance, or – in the context of DMRG – the truncation error. Both quantities are also discussed in Appendix A 4. We show results for the truncation error further down in Sec. III B 4, and when comparing VUMPS to IDMRG and ITEBD in Sec. III C.

We show results for examples of 3 gapped and 3 critical systems in Fig. 1. Specifically, as examples for gapped systems we considered (a) the TFI model (28) in the symmetry broken ferromagnetic phase at $h = 0.48$, (b) the isotropic $S = 1$ Heisenberg antiferromagnet, i.e. the $S = 1$ XXZ model (29) at $\Delta = 1$ and (c) the $S = 1/2$ XXZ model (29) in the symmetry broken antiferromagnetic phase at $\Delta = 2$. As examples for gapless systems we considered (d) the isotropic $S = 1/2$ Heisenberg antiferromagnet, i.e. the $S = 1/2$ XXZ model (29) at $\Delta = 1$, (e) the repulsive Fermi Hubbard model (30) at $U = 10$ and half filling and (f) the Haldane-Shastry model (31).

Out of the gapped systems, only the antiferromagnetic ground state of (c) physically breaks translation invariance by spontaneously breaking the \mathbb{Z}_2 spin-flip symmetry; we therefore choose a two-site unit cell in this case. The critical systems physically show no spontaneous symmetry breaking. However, for uMPS ground state approximations it is often energetically beneficial to artificially break symmetries (which are restored in the limit of infinite bond dimension). In all three cases, the optimal uMPS ground state approximation artificially breaks a $SU(2)$ symmetry and develops antiferromagnetic order, breaking translation invariance. We therefore choose a two-site unit cell in the case of the Hubbard model (e) and the Haldane-Shastry model (f). In the case of the Heisenberg antiferromagnet (d), translation invariance can be restored through a unitary transformation by rotating every second spin by π around the z -axis, transforming $H_{XXZ}(\Delta) \rightarrow -H_{XXZ}(-\Delta)$, and the artificially symmetry broken ground state becomes ferromagnetically ordered along the x and y directions. We can therefore choose a single-site unit cell for (d), and the staggered magnetization along z is thus zero. A similar approach could be chosen to restore translation invariance also for the gapped antiferromagnet (c) and the Hubbard model (e), but we do not choose to do so for demonstrative reasons. To summarize, we used the single-site Algorithm 2 for (a), (b) and (d), and the sequential Algorithm 3 with a two-site unit cell for (c), (e) and (f). For a comparison between the sequential and parallel approach see Sec. III B 5.



0.696198978154358	0.001665909341209
0.696198978154358	0.001487766860776
0.086098881485240	0.001487766860776
0.086098881485240	0.001487766860776
0.086098881485239	0.001487766860776
0.086098881485239	0.001487766860776
0.020013261627349	0.001487766860776
0.020013261627349	0.001106527294364
0.005770050481551	0.001106527294364
0.005770050481551	0.000410363691742
0.005770050481551	0.000410363691742
0.005770050481551	0.000307343959372
0.001665909341209	0.000307343959372
0.001665909341209	0.000307343959372
0.001665909341209	0.000307343959372

Figure 2. *Top*: Schmidt spectrum of the $S = 1$ Heisenberg antiferromagnet for $D = 120$, converged to gradient norm $\|B\| < 10^{-15}$. *Bottom*: The table shows the first 30 Schmidt values in descending order. The degeneracies are reproduced to 15 digits of precision, without exploiting any symmetries.

1. General Convergence

Above all we observe that VUMPS shows unprecedented fast convergence, both in the energy density e and the norm of its gradient $\|B\|$, and excellent accuracy of the final ground state approximation. Observe in Fig. 1 that in all cases the energy is already well converged to machine precision after $\mathcal{O}(10 - 50)$ iterations, while the state is still quite some distance from the variational optimum, according to the gradient norm $\|B\|$. Further optimizing the state, this quantity can also be converged to essentially machine precision (even in the presence of small Schmidt values), while the energy virtually doesn't change anymore. The resulting final state then corresponds to the variationally optimal state for the given bond dimension. This is very useful in the case where the quantum state itself is required to be accurate to high precision, e.g. when used as a starting state for real time evolution or as a starting point to compute excited states and scattering thereof.^{44–46}

The Schmidt spectrum of the ground state of the $S = 1$ Heisenberg antiferromagnet at $D = 120$, converged to

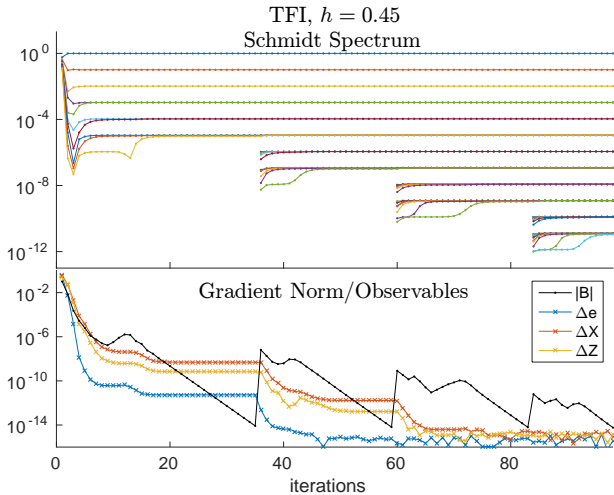


Figure 3. Evolution of the Schmidt spectrum (top) and the gradient norm $\|B\|$ and various observables (bottom) with iteration number for the TFI model at $h = 0.45$. Here we defined the deviation of an observable O from its exact value as $\Delta O = |\langle O \rangle - \langle O \rangle_{\text{exact}}|$, similar to (32) for the energy. We used bond dimensions $D = [9, 19, 33, 55]$, i.e. we increased the bond dimension three times during the optimization process as soon as $\|B\|$ dropped below 10^{-14} (at iterations 36, 60 and 84). It is apparent that while high lying Schmidt values converge quite quickly, the better part of the final iterations goes into converging low lying Schmidt values. Moreover one can see that there is quite some rearrangement of exactly these low lying Schmidt values every time $\|B\|$ reaches a local maximum (e.g. around iterations 12, 40, 70, 88 and 92) during a non-monotonous phase of gradient evolution (see also Sec. III B 2).

gradient norm $\|B\| < 10^{-15}$, is depicted in Fig. 2. It can be seen that the degeneracies are reproduced perfectly to the same precision, without explicitly exploiting any symmetries in the implementation of the algorithm.

In cases where the final desired bond dimension D_{final} is not known beforehand, one can successively enlarge a state of some small initial bond dimension every few iterations until the state fulfills the desired criteria, e.g. current bond dimension above some threshold, truncation error (see below) or smallest Schmidt value below some threshold, etc. This strategy is particularly useful when using an implementation exploiting physical symmetries of the system, such as e.g. conservation of magnetization or particle number, as the correct number and size of the required symmetry sectors in the MPS tensors is generally not known beforehand.⁴⁷ On the other hand, if D_{final} is known beforehand, it generally appears to be more efficient to immediately start from an initial state with $D = D_{\text{final}}$. The gain in computational time due to the cheaper initial iterations with small bond dimension is usually outweighed by a considerable number of required additional iterations. On the other hand, for some hard problems (e.g. the Hubbard model) stability and convergence speed can profit from a strategy of sequentially

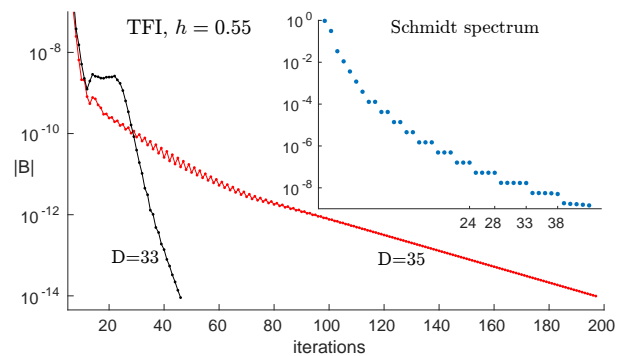


Figure 4. Comparison of the convergence rate of the gradient norm $\|B\|$ for the TFI model at $h = 0.55$, with $D = 33$ and $D = 35$. Convergence is roughly 4 times faster for $D = 33$ as compared to $D = 35$. The inset shows the Schmidt spectrum of the ground state (up to $D = 43$). For $D = 35$ the smallest Schmidt values form an incomplete degenerate multiplet, whereas for $D = 33$ the multiplet is complete.

increasing the bond dimension from some small initial value.

To conclude the discussion of general convergence, we plot the evolution of the Schmidt spectrum, as well as the gradient norm $\|B\|$ and various observables vs. iteration number during a ground state optimization for the TFI model (28) in the ferromagnetic phase at $h = 0.45$ in Fig. 3. During the simulation we used a sequence of bond dimensions $D = [9, 19, 33, 55]$, where we started with an initial random state with $D = 9$ and increased the bond dimension to the next value as soon as $\|B\|$ dropped below 10^{-14} . We chose this set of bond dimension in order to not cut any degenerate multiplets of Schmidt values (see also Sec. III B 3). It can be seen that the high lying Schmidt values converge quite quickly, while most of the computational time goes into converging the low lying Schmidt values. Moreover, there is quite some rearrangement of the small Schmidt values every time the gradient norm $\|B\|$ reaches a local maximum during a phase of non-monotonous evolution (see also next subsection). Lastly, from the evolution of the errors of the local observables $\langle X \rangle$ and $\langle Z \rangle$ it is apparent that they require a substantially higher bond dimension of $D = 55$ to reach the same accuracy as the energy, which is already correct to machine precision at $D = 19$.

2. Different Regimes of Gradient Norm Convergence

Depending on the complexity of the model, the gradient norm shows a period of irregular non-monotonous behavior before entering a regime of monotonous convergence. This can be understood as the (random) initial state having to adapt its initial structure (e.g. the Schmidt spectrum) to the requirements of the best variational ground state approximation. Monitoring the Schmidt values during this period nicely shows how

groups of Schmidt values slightly rearrange to the correct structure (see e.g. Fig. 3). This period is usually more dominant in critical systems – as can be seen in Fig. 1 (d) - (f), where it takes $\mathcal{O}(50 - 100)$ iterations – and of course strongly depends on the chosen initial state. One could argue, that the jumps in parameter space caused by the algorithm during this period are too big for the state to find the correct structure quickly, hindering a fast crossover to the regime of monotonous convergence. However, an approach of preconverging the state using smaller steps through parameter space – e.g. by means of imaginary time evolution with moderate time steps – has proven to be even slower in all cases tried. Thus, the best choice is still to use VUMPS during the entire optimization process. We want to emphasize here, that we have never observed a stagnation of the algorithm during this initial regime; the algorithm always reached the monotonous regime eventually in all cases, and instances where the algorithm remains in the irregular regime for an unusually long time are rare and only occur in the case of particularly hard problems.

As soon as the gradient norm reaches the monotonous regime, it always converges exponentially fast. Surprisingly, this is true even for critical systems, where one would in principle expect algebraic convergence. This can be qualitatively understood from the theory of finite entanglement scaling,^{48–50} which states that the MPS approximation itself introduces a small perturbation away from criticality, and thus a finite gap. However, as VUMPS improves convergence speed over existing methods (see also Sec. III C), it is ideally suited to study critical systems via the theory of finite entanglement scaling, which still requires that one finds the optimal MPS approximation in the first place.

3. Degenerate Schmidt Values

In the presence of multiplets of degenerate Schmidt values, the convergence rate is severely affected if the smallest few Schmidt values are part of an incomplete multiplet, i.e. if the last multiplet is “cut”. In that case the algorithm still shows stable convergence, albeit at a greatly reduced rate. For an example in the TFI model see Fig. 4. This issue can be easily circumvented by ensuring that the smallest few Schmidt values are part of a complete multiplet when dynamically increasing the bond dimension, or by choosing a viable (or reducing from some) fixed initial bond dimension.

4. Energy Convergence with Bond Dimension

In a careful MPS study, variational energies obtained for different bond dimension D are compared in order to extrapolate to the exact $D \rightarrow \infty$ limit. This can be done by plotting the energy $e(D)$ as a function of bond dimension against the inverse of the bond dimension $1/D$. The

infinite D limit is then obtained by fitting with a power law form and extrapolating to $1/D \rightarrow 0$. In DMRG, another popular measure for the quality of an MPS approximation is given by the truncation error or discarded weight ϵ_ρ , defined in Eq. (A28). The variational energy is found to scale linearly with ϵ_ρ ^{51,52} and an extrapolation to $\epsilon_\rho \rightarrow 0$ is thus generally easier and more stable. For further details on assessing the quality of the ground state approximation we refer to Appendix A 4.

We show an example for both extrapolation schemes for the isotropic $S = 1/2$ Heisenberg antiferromagnet in Fig. 5. The exact ground state energy is given by $e_{\text{exact}} = \frac{1}{4} - \log(2)$, or as numerical value by $e = -0.4431471805599453$ to 16 digits of precision. On the left we plot the energy vs. truncation error and obtain an estimate $e_T = -0.443147178(1)$ with 9 digits of precision from a linear fit $e(\epsilon_\rho) = e + a\epsilon_\rho$. On the right we plot the energy vs. inverse bond dimension and obtain an estimate $e_D = -0.4431471797(1)$ with 10 digits of precision from a power law fit $e(1/D) = e + a(1/D)^b$. Comparing to e_{exact} we observe that e_T has an error of $\Delta e_T \approx 3 \times 10^{-9}$, while e_D has an error $\Delta e_D \approx 8 \times 10^{-10}$.

5. Multi Site Unit Cell Implementations

Lastly we discuss and compare the performance of the *sequential* and *parallel* algorithms for multi-site unit cells presented in Sec. II E. As the two methods differ in their convergence with the number of iterations, as well as in the computational effort for each iteration of updating the entire unit cell, we compare the rate of convergence with absolute computing time t in seconds. To that end we only time operations that are absolutely necessary for each algorithm, i.e. we do not time measurements, data storage etc. We further start from the same (random) initial state and keep the bond dimension fixed throughout the entire simulation for both methods to make the simulations as comparable as possible. All calculations are performed using a non-parallelized MATLAB implementation on a single core of a standard laptop CPU.

We find that in general for gapped systems, the sequential approach outperforms the parallel approach, while in critical systems no definite statement about better performance can be made. There are instances where one algorithm takes substantially longer than the other to reach the regime of monotonous convergence, but such cases are rare and strongly depend on the model and initial state. In Fig. 6 we show two extreme examples of such behavior for the critical Hubbard model (30) with $U = 5$ and $D = 65$ and a $N = 2$ site unit cell. In the left plot the sequential approach is clearly more efficient than the parallel approach, while in the right plot the opposite is the case. The only difference between those two cases is the chosen initial state. Overall both approaches however generally show comparable performance, with the sequential approach appearing to be slightly more stable and reliable in the cases we considered.

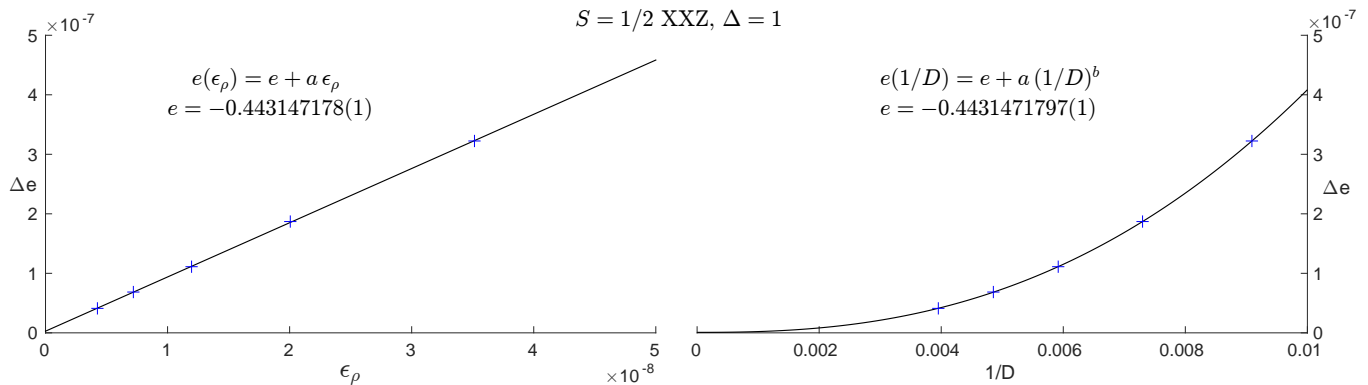


Figure 5. Scaling of the variational ground state energy e with truncations error ϵ_ρ and bond dimension D for the isotropic $S = 1/2$ Heisenberg antiferromagnet. We plot the energy e vs. truncation error ϵ_ρ on the left, and vs. inverse bond dimension $1/D$ on the right. The exact ground state energy is given by $e_{\text{exact}} = -0.4431471805599453$ to 16 digits of precision. We obtain estimates from a linear fit $e(\epsilon_\rho) = e + a\epsilon_\rho$ (left) and a power law fit $e(1/D) = e + a(1/D)^b$ (right), where the estimate on the right has one more digit of precision and is roughly 4 times more accurate than the estimate on the left.

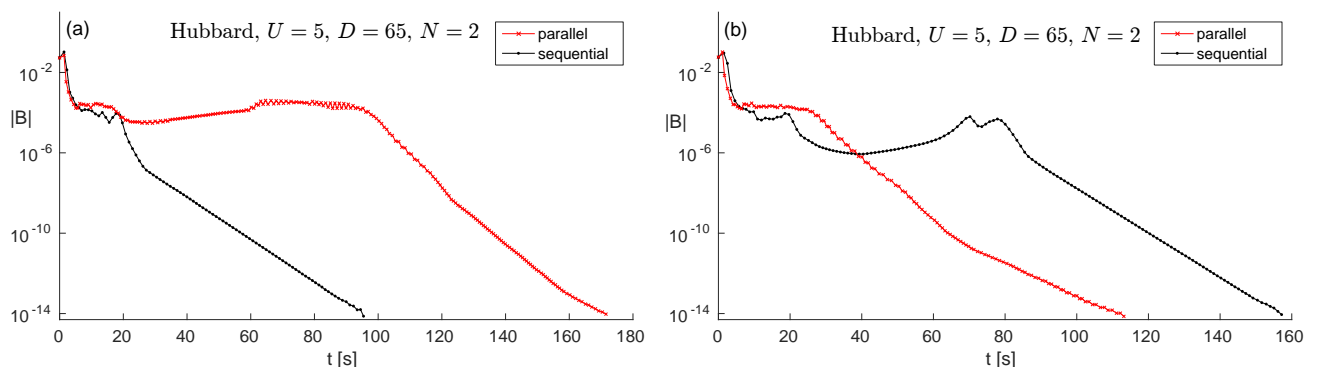


Figure 6. Example performance comparisons between the sequential and parallel algorithm presented in Sec. II E. We show two extreme examples where performance differs greatly between the two approaches in the critical Hubbard model at $U = 5$ and $D = 65$, on a two-site unit cell, starting from two different (random) initial states. In the left example the sequential approach is clearly faster than the parallel approach, while in the right example the opposite is the case. We want to emphasize here that such examples are the exception and overall both approaches generally show comparable performance.

Once both algorithms are in the monotonous regime, convergence speed in terms of absolute computing time is also similar, with the sequential approach generally taking longer for each iteration, but the parallel approach generally requiring more iterations to reach convergence.

C. Comparison with IDMRG and ITEBD

We further benchmark the performance of VUMPS against a standard two-site IDMRG implementation,^{5,6,24} and a standard two-site ITEBD implementation,^{26,36} and compare the rate of convergence of the energy error Δe and the norm of the energy gradient $\|B\|$ between the three methods. For VUMPS, we solve the effective eigenvalue problems in each iteration to precision $\epsilon_H = \epsilon_{\text{prec}}/100$ with ϵ_{prec} the current precision according to (24). For IDMRG we solve the effective two-site eigenvalue problem in each iteration to

precision $\epsilon_H = (1 - F)/100$, with F the current orthogonality fidelity F (see Sec. III.A in Ref. 24). For ITEBD we employ a fourth order Suzuki Trotter decomposition and measure every 10 time steps. We use a sequence of time steps $\delta t \in [10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}]$ where we decrease the time step as soon as the change in Schmidt values per unit of imaginary time drops below a certain threshold. Naturally, the strategy of time step reduction should be optimized carefully for each model under consideration, however we choose the same strategy for all example cases to maintain comparability.

We also explicitly calculate all necessary quantities for obtaining a truly variational energy (e.g. by reorthogonalizing the unit cell) and for measuring the energy gradient, even if these quantities are not necessary for the respective algorithm itself.

As the three methods differ quite substantially in the number of iterations required for convergence, as well as the computational effort for each iteration, we again

compare convergence against absolute computing time t in seconds, where we only time operations that are absolutely necessary for each algorithm and we do not time measurements, data storage, reorthogonalizing, etc. We further start from the same (random) initial state and keep the bond dimension fixed throughout the entire simulation for all three methods to make the simulations as comparable as possible. Again, all calculations were performed using a non-parallelized MATLAB implementation on a single core of a standard laptop CPU.

We show example comparisons for two gapped and two critical models in Fig. 7, similar to the cases studied in the previous section. Specifically, we show results for (a) the gapped isotropic $S = 1$ Heisenberg antiferromagnet, i.e. the $S = 1$ XXZ model (29) at $\Delta = 1$, (b) the $S = 1/2$ XXZ model (29) in the gapped symmetry broken antiferromagnetic phase at $\Delta = 2$, (c) the critical isotropic $S = 1/2$ Heisenberg antiferromagnet, i.e. the $S = 1/2$ XXZ model (29) at $\Delta = 1$, and (d) the critical Fermi Hubbard model (30) at $U = 5$ and half filling. We plot the energy error Δe on the left and the gradient norm $\|B\|$ on the right, vs. absolute computing time t in seconds. For VUMPS we used a single-site unit cell for (a) and (c), and a two-site unit cell for (b) and (d).

Above all we observe that VUMPS clearly outperforms both IDMRG and ITEBD by far, both in convergence speed and accuracy of the final state, especially for critical systems. In all shown cases, the final energy error Δe of all three algorithms only differ by a few percent; VUMPS however always yields the best variational energy, often already after a few seconds, and thus converges in energy much faster than IDMRG or ITEBD – in the case of critical systems even by orders of magnitude. Observe that especially for the two critical systems (c) and (d) a large part of the computational time of IDMRG and ITEBD goes in converging the last few digits of the energy (see also insets in Fig. 7). For (d) in particular, the final energy error obtained by IDMRG is still almost 10% higher than the value obtained by VUMPS.

In terms of convergence of the energy gradient $\|B\|$, we observe that IDMRG and ITEBD perform quite poorly. Surprisingly, IDMRG usually stagnates at some value $\|B\| > 10^{-7}$. ITEBD on the other hand would be in principle capable of converging $\|B\|$ essentially also to machine precision, albeit at prohibitively long simulation times, as the limiting factor appears to be the Trotter error, requiring very small time steps; we therefore also only reach values of $\|B\| \gtrsim 10^{-10}$ with ITEBD within reasonable simulation times.⁵³ VUMPS on the other hand is always capable to converge $\|B\|$ essentially to machine precision, and does so – contrary to other methods – exponentially fast and with unprecedented speed. For instance, in the case of the Hubbard model in example (d), ITEBD only reached a gradient norm of $\|B\| = 1.3 \times 10^{-7}$ after ≈ 60 hours of absolute computing time, while VUMPS already reached this value after only ≈ 30 seconds, and converged further to $\|B\| < 10^{-14}$ in

≈ 90 seconds. IDMRG on the other hand stagnated at a quite high value of $\|B\| \approx 2 \times 10^{-4}$.

1. Observables

We also measure and compare the regular and staggered (averaged) magnetizations m_r and m_s of the final state after convergence for the Hubbard model in example (d), as in this case all three methods use a two-site unit cell. The exact ground state is $SU(2)$ symmetric and thus has zero magnetization; a finite D ground state approximation however artificially breaks this symmetry. The final values for the regular magnetization m_r are zero to machine precision for both VUMPS and IDMRG, but $m_r = 8 \times 10^{-12}$ for ITEBD. The staggered magnetization is $m_s = 0.011162$ for VUMPS, $m_s = 0.080768$ for IDMRG, and $m_s = 0.034797$ for ITEBD. Both the regular and staggered magnetizations are thus smallest for the final state obtained from VUMPS. For IDMRG the staggered magnetization is highest, but the regular magnetization is zero, which in turn is finite for ITEBD. This result is not surprising, as VUMPS yields the best variational state out of the three methods.

2. Truncation Error

As a last figure of merit, popular in DMRG studies as a measure of the quality of the MPS ground state approximation and used for extrapolations to the exact infinite D limit, we also calculate the truncation error or discarded weight ϵ_ρ of the final state, defined in Eq. (A28). In the case of the Hubbard model in example (d), we obtain a truncation error $\epsilon_\rho = 2.54438 \times 10^{-6}$ from IDMRG, and a slightly lower $\epsilon_\rho = 2.45138 \times 10^{-6}$ from VUMPS.

IV. CONCLUSION AND OUTLOOK

We have introduced a novel algorithm for calculating MPS ground state approximations of strongly correlated one dimensional quantum lattices models with nearest neighbor or long range interactions, in the thermodynamic limit. It combines ideas from conventional DMRG and tangent space methods by variationally optimizing a uniform MPS by successive solutions of effective eigenvalue problems. The algorithm can easily be implemented by extending an existing single-site (I)DMRG implementation with routines for i) calculating effective Hamiltonian contributions from infinite environments and ii) solving an effective “zero site” eigenvalue problem in addition to the usual single-site problem. The new algorithm is free of any ill-conditioned inverses and therefore does not suffer from small Schmidt values, contrary to other tangent space methods such as TDVP. Additionally, as it does not rely on imaginary time evo-

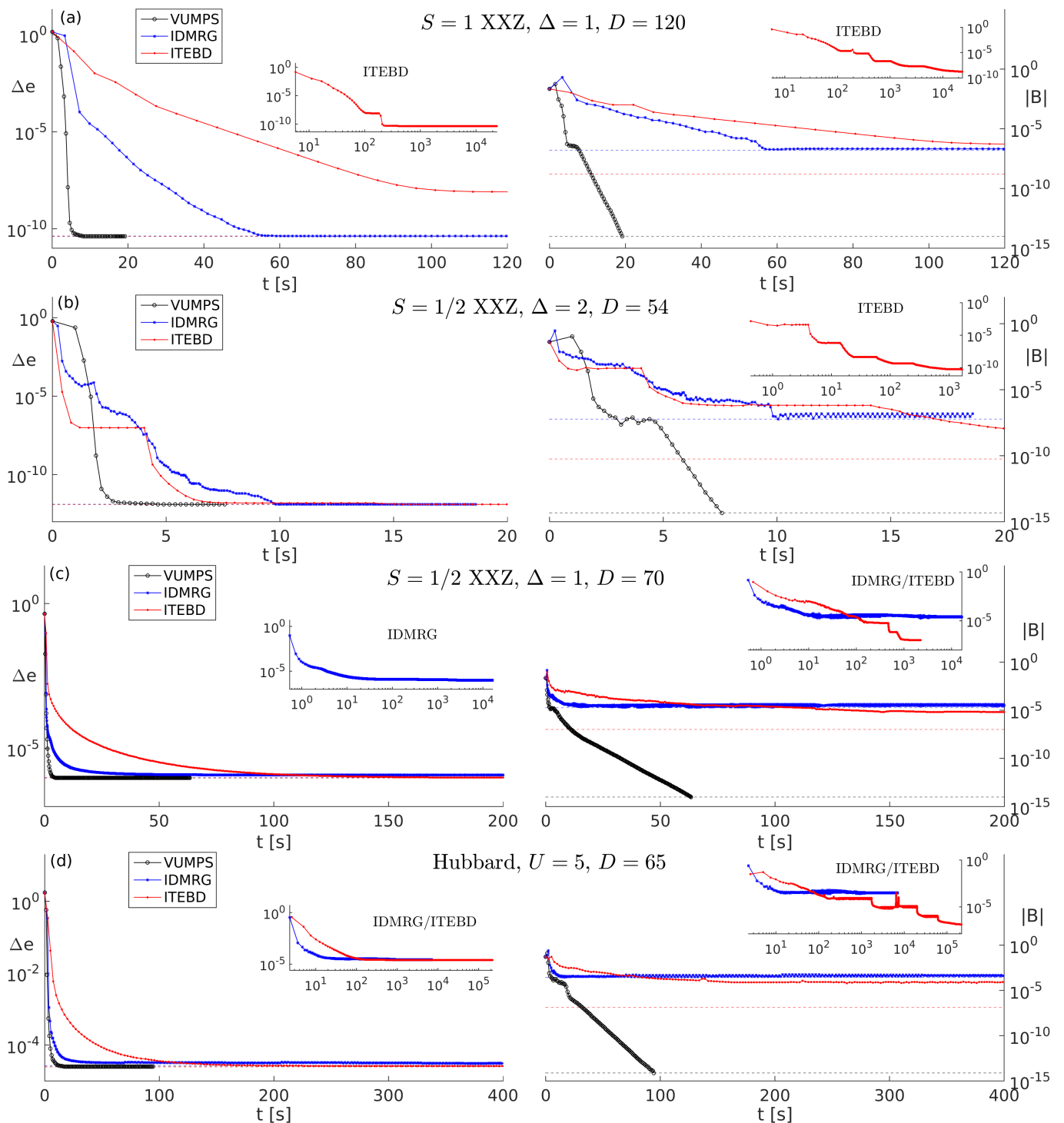


Figure 7. Comparative benchmark plots for VUMPS, IDMRG and ITEBD. We plot the error Δe on the left, and the gradient norm $\|B\|$ on the right, vs. total computing time t in seconds for (a) the gapped isotropic $S = 1$ XXZ antiferromagnet at $D = 120$, (b) the gapped $S = 1/2$ XXZ antiferromagnet at $\Delta = 2$ and $D = 54$, (c) the critical isotropic $S = 1/2$ XXZ antiferromagnet at $D = 70$, and (d) the critical Hubbard model at $U = 5$ and $D = 65$. The dashed lines are a guide to the eye and denote the minimum values of Δe and $\|B\|$ obtained by the respective algorithm. The insets show a plot of the entire ITEBD and/or IDMRG simulation with logarithmic time scale. It is obvious that VUMPS reaches convergence orders of magnitude faster than IDMRG or ITEBD, especially for critical systems. Notice also that, while Δe differs only by a few percent between the different algorithms (up to $\approx 10\%$ for (d)), VUMPS always manages to also converge $\|B\|$ essentially to machine precision, whereas IDMRG and ITEBD stagnate at some substantially higher values, remaining quite far from the variational optimum.

lution, it is especially fit for studying systems with long range interactions.

We described and benchmarked implementations for uniform MPS with both single-site and multi-site unit cells. We observed that the new algorithm clearly outperforms existing methods such as IDMRG and ITEBD, both in convergence speed and accuracy of the final state at convergence. The energy converges with unprecedented speed after $\mathcal{O}(10 - 50)$ iterations, even in critical systems (where this is orders of magnitude faster than conventional methods). The algorithm further proceeds to converge the state to the variational optimum by minimizing the energy gradient essentially to machine precision; it does so exponentially fast, even for critical systems, contrary to other methods. The new algorithm is thus the perfect choice for studying critical systems. Additionally, a state converged to the variational optimum is particularly useful in cases where the quantum state itself is required to be accurate to high precision, e.g. when used as a starting state for real time evolution or for a variational calculation of elementary excitations.^{44–46}

It is straightforward to include physical symmetries that come with good quantum numbers (such as e.g. conserved magnetization or particle number) after a proper definition of a symmetric uniform MPS unit cell, where absolute (diverging) values of these quantum numbers are replaced by densities. All steps of the algorithm then immediately also apply to MPS tensors with good quantum numbers. Symmetric ground states obtained this way are an excellent starting point for obtaining elementary excitations with well defined quantum numbers following Ref. 44, which for instance enables to target elementary excitation that lie within a multi-particle continuum.⁴⁷

Within the same framework it is also very natural to

recover real or imaginary time evolution by replacing the effective Hamiltonian ground state problems (23a) and (23b) by small finite time evolution steps, which yields the thermodynamic limit version of the time evolution algorithm presented in Ref. 28. This enables e.g. to efficiently study real time evolution of quantum states on systems with long range interactions in the thermodynamic limit.^{54,55}

We believe that the ideas presented in this paper should be relevant for other classes of tensor-network states as well. Specifically in the case of projected entangled pair states (PEPS), designed to capture ground states on two-dimensional quantum lattices, the further development of efficient variational algorithms – as an alternative to ITEBD inspired approaches – is still much desired.^{56,57} In particular, it motivates the search for (approximate) canonical forms for PEPS, which would enable a translation of the VUMPS algorithm to the two-dimensional setting.

ACKNOWLEDGMENTS

The authors acknowledge inspiring and insightful discussions with M. Ganahl, M. Mariën and I.P. McCulloch. We thank Boye Buyens for providing templates for drawing tensor network diagrams. This project has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant agreement No 647905). The authors gratefully acknowledge support from the Austrian Science Fund (FWF): F4104 SFB ViCoM and F4014 SFB FoQuS (V.Z.-S. and F.V.), and GRW 1-N36 (M.F.). J.H. and L.V. are supported by the Research Foundation Flanders (FWO).

-
- ¹ K. Wilson, *Rev. Mod. Phys.* **47**, 773 (1975).
² M. Fisher, *Rev. Mod. Phys.* **46**, 597 (1974).
³ M. Fisher, *Rev. Mod. Phys.* **70**, 653 (1998).
⁴ J. Zinn-Justin, *Quantum Field Theory and Critical Phenomena* (Oxford University Press, 1996).
⁵ S. White, *Phys. Rev. Lett.* **69**, 2863 (1992).
⁶ S. White, *Phys. Rev. B* **48**, 10345 (1993).
⁷ M. Fannes, B. Nachtergaele, and R. Werner, *Comm. Math. Phys.* **144**, 443 (1992).
⁸ A. Klümper, A. Schadschneider, and J. Zittartz, *Europhys. Lett.* **24**, 293 (1993).
⁹ S. Östlund and S. Rommer, *Phys. Rev. Lett.* **75**, 3537 (1995).
¹⁰ D. Perez-Garcia, F. Verstraete, M. Wolf, and J. Cirac, *Quant. Inf. Comput.* **75**, 401 (2007).
¹¹ F. Verstraete, V. Murg, and J. Cirac, *Adv. Phys.* **57**, 143 (2008).
¹² U. Schöllwöck, *Ann. Phys.* **326**, 96 (2011).
¹³ J. Haegeman, T. Osborne, and F. Verstraete, *Phys. Rev. B* **88**, 075133 (2013).
¹⁴ R. Orús, *Ann. Phys.* **349**, 117 (2014).
¹⁵ J. Bridgeman and C. Chubb, “Hand-waving and Interpretive Dance: An Introductory Course on Tensor Networks,” [arXiv:1603.03039](https://arxiv.org/abs/1603.03039).
¹⁶ M. Hastings, *J. Stat. Mech. Theor. Exp.* **2007**, P08024 (2007).
¹⁷ J. Eisert, M. Cramer, and M. Plenio, *Rev. Mod. Phys.* **82**, 277 (2010).
¹⁸ F. Verstraete and J. Cirac, *Phys. Rev. B* **73**, 094423 (2006).
¹⁹ M. Hastings, *Phys. Rev. B* **76**, 035114 (2007).
²⁰ J. Haegeman, S. Michalakis, B. Nachtergaele, T. Osborne, N. Schuch, and F. Verstraete, *Phys. Rev. Lett.* **111**, 080401 (2013).
²¹ F. Verstraete, D. Verstraete, and J. Cirac, *Phys. Rev. Lett.* **93**, 227205 (2004).
²² S. White, *Phys. Rev. B* **72**, 180403 (2005).
²³ I. McCulloch, *J. Stat. Mech.*, P10014 (2007).
²⁴ I. McCulloch, “Infinite size density matrix renormalization group, revisited,” [arXiv:0804.2509](https://arxiv.org/abs/0804.2509).
²⁵ G. Vidal, *Phys. Rev. Lett.* **91**, 147902 (2003).
²⁶ G. Vidal, *Phys. Rev. Lett.* **98**, 070201 (2007).
²⁷ J. Haegeman, J. Cirac, T. Osborne, I. Pižorn, H. Verschelde, and F. Verstraete, *Phys. Rev. Lett.* **107**, 070601

- (2011).
- ²⁸ J. Haegeman, C. Lubich, I. Oseledets, B. Vandereycken, and F. Verstraete, *Phys. Rev. B* **94**, 165116 (2016).
- ²⁹ In many cases $h_{j,j+1}$ is sparse and the number d_h of non-zero elements is usually of the order $\mathcal{O}(d^2)$. The first two terms can then be applied in $\mathcal{O}(d_h D^3)$ operations.
- ³⁰ R. Bhatia, *Matrix Analysis* (Springer-Verlag, 1997).
- ³¹ These values can be different and depend on the subtraction scheme for the divergent energy expectation value. If $\tilde{h} \rightarrow \tilde{h}$ is performed everywhere, we have $E_{AC} = E_C = 0$. If, in the case of nearest neighbor interactions, we only substitute $h \rightarrow \tilde{h}$ in the construction of H_L and H_R , but not in the local terms, we will have $E_{AC} = 2E_C = 2e$.
- ³² Further approximations to comparable accuracy can be made within the construction of the effective Hamiltonians, e.g. when determining H_L and H_R to precision ϵ_S . There, the approximations $\tilde{R} = CC^\dagger$ and $\tilde{L} = C^\dagger C$ for the true L and R needed for some of these operations are good enough, if ϵ_S is roughly of the same order of magnitude as ϵ_{prec} .
- ³³ H. Phien, I. McCulloch, and G. Vidal, *Phys. Rev. B* **91**, 115137 (2015).
- ³⁴ B. Pirvu, V. Murg, J. Cirac, and F. Verstraete, *New J. Phys.* **12**, 025012 (2010).
- ³⁵ M. Zaletel, R. Mong, C. Karrasch, J. Moore, and F. Pollmann, *Phys. Rev. B* **91**, 165112 (2015).
- ³⁶ M. Hastings, *J. Math. Phys.* **50**, 095207 (2011).
- ³⁷ S. Singh, R. Pfeifer, and G. Vidal, *Phys. Rev. B* **83**, 115125 (2011).
- ³⁸ P. Pfeuty, *Ann. Phys.* **57**, 79 (1970).
- ³⁹ M. Takahashi, *Thermodynamics of one-dimensional solvable models* (Cambridge University Press, 1999).
- ⁴⁰ E. Lieb and F. Wu, *Phys. Rev. Lett.* **20**, 1445 (1968).
- ⁴¹ F. Essler, H. Frahm, F. Göhmann, A. Klümper, and V. Korepin, *The One-Dimensional Hubbard Model* (Cambridge University Press, 2005).
- ⁴² F. Haldane, *Phys. Rev. Lett.* **60**, 635 (1988).
- ⁴³ B. Shastri, *Phys. Rev. Lett.* **60**, 639 (1988).
- ⁴⁴ J. Haegeman, B. Pirvu, D. J. Weir, J. I. Cirac, T. J. Osborne, H. Verschelde, and F. Verstraete, *Physical Review B* **85**, 100408 (2012).
- ⁴⁵ L. Vanderstraeten, J. Haegeman, T. J. Osborne, and F. Verstraete, *Physical review letters* **112**, 257202 (2014).
- ⁴⁶ L. Vanderstraeten, F. Verstraete, and J. Haegeman, *Physical Review B* **92**, 125136 (2015).
- ⁴⁷ V. Zauner-Stauber, I. McCulloch, and F. Verstraete, (unpublished).
- ⁴⁸ L. Tagliacozzo, T. R. de Oliveira, S. Iblisdir, and J. I. Latorre, *Phys. Rev. B* **78**, 024410 (2008).
- ⁴⁹ F. Pollmann, S. Mukerjee, A. M. Turner, and J. E. Moore, *Phys. Rev. Lett.* **102**, 255701 (2009).
- ⁵⁰ V. Stojevic, J. Haegeman, I. P. McCulloch, L. Tagliacozzo, and F. Verstraete, *Phys. Rev. B* **91**, 035120 (2015).
- ⁵¹ S. White and D. Huse, *Phys. Rev. B* **48**, 3844 (1993).
- ⁵² Ö. Legeza and G. FÁth, *Phys. Rev. B* **53**, 14349 (1996).
- ⁵³ The fact that ITEBD usually yields better variational states than IDMRG for the same bond dimension has already been observed in Ref. 24.
- ⁵⁴ J. Halimeh and V. Zauner-Stauber, “Enriching the dynamical phase diagram of spin chains with long-range interactions,” [arXiv:1610.02019](https://arxiv.org/abs/1610.02019).
- ⁵⁵ J. Halimeh, V. Zauner-Stauber, I. McCulloch, I. de Vega, U. Schollwöck, and M. Kastner, *Phys. Rev. B* **95**, 024302 (2017).
- ⁵⁶ P. Corboz, *Phys. Rev. B* **94**, 035133 (2016).
- ⁵⁷ L. Vanderstraeten, J. Haegeman, P. Corboz, and F. Verstraete, *Phys. Rev. B* **94**, 155123 (2016).
- ⁵⁸ J. Haegeman and F. Verstraete, “Diagonalizing transfer matrices and matrix product operators: a medley of exact and computational methods,” [arXiv:1611.08519](https://arxiv.org/abs/1611.08519).
- ⁵⁹ *Journal of Mathematical Physics* **55**, 021902 (2014).
- ⁶⁰ C. Hubig, I. McCulloch, U. Schollwöck, and F. Wolf, *Phys. Rev. B* **91**, 155115 (2015).
- ⁶¹ For fermionic Hamiltonians with long range interactions, we can employ a Jordan-Wigner transformation to spin operators, introducing a string operator counting the number of fermions between o_j and o_{j+n} ; geometric sums of transfer matrices in (C4) then turn into geometric sums of (string) operator transfer matrices.
- ⁶² D. Ruelle, *Comm. Math. Phys.* **9**, 267 (1968).
- ⁶³ F. Dyson, *Comm. Math. Phys.* **12**, 91 (1969).
- ⁶⁴ J. Cardy, *J. Phys. A: Math. Gen.* **144**, 1407 (1981).
- ⁶⁵ F. Verstraete, J. Garcia-Ripoll, and J. Cirac, *Phys. Rev. Lett.* **93**, 207204 (2004).
- ⁶⁶ G. Crosswhite and D. Bacon, *Phys. Rev. A* **78**, 012356 (2008).
- ⁶⁷ G. Crosswhite, A. Doherty, and G. Vidal, *Phys. Rev. B* **78**, 035116 (2008).
- ⁶⁸ F. Fröwis, V. Nebendahl, and W. Dür, *Phys. Rev. A* **81**, 062337 (2010).
- ⁶⁹ L. Michel and I. McCulloch, “Schur Forms of Matrix Product Operators in the Infinite Limit,” [arXiv:1008.4667](https://arxiv.org/abs/1008.4667).
- ⁷⁰ H. van der Vorst, *SIAM J. Sci. Stat. Comput.* **13**, 631 (1992).
- ⁷¹ Y. Saad and M. Schultz, *SIAM J. Sci. Stat. Comput.* **75**, 856 (1986).

Appendix A: Theoretical background

In this Appendix we reiterate definitions and concepts needed for the algorithm presented in Sec. II in more detail, and motivate the VUMPS algorithm from a variational perspective.

1. Variational principle on manifolds

The variational principle in quantum mechanics characterizes the ground state of a given Hamiltonian as the state $|\Psi\rangle$ which minimizes the normalized energy expectation value

$$E = \frac{\langle \Psi | H | \Psi \rangle}{\langle \Psi | \Psi \rangle}.$$

If (typically for computational reasons) we only have access to a subset of Hilbert space, the variational principle still gives a way to find an approximation to the true ground state, namely by solving the minimization problem within the restricted set. If this subset is a linear subspace spanned by a number of basis vectors $\{|i\rangle, i = 1, \dots, N\}$, we obtain a generalized eigenvalue problem

$$\langle i | H | j \rangle c_j = E \langle i | j \rangle c_j$$

for the expansion coefficients c_i in $|\Psi\rangle = \sum_{i=1}^N c_i |i\rangle$. This is known as the Rayleigh-Ritz method, and by orthonormalizing the basis it clearly amounts to projecting the full time-independent Schrödinger equation into the variational subspace.

If, more generally, we have a variational ansatz $|\Psi(\mathbf{A})\rangle$ which depends analytically on a number of complex parameters, as encoded in the complex vector \mathbf{A} , a variational minimum $|\Psi(\mathbf{A}^*)\rangle$ is characterized by a vanishing gradient of the energy expectation value, i.e.

$$\langle \partial_i \Psi(\bar{\mathbf{A}}^*) | H - E(\bar{\mathbf{A}}^*, \mathbf{A}^*) | \Psi(\mathbf{A}^*) \rangle = 0, \quad (\text{A1})$$

with $\bar{\mathbf{A}}$ the (formally independent) complex conjugate of \mathbf{A} , ∂_i and $\bar{\partial}_i$ the complex derivatives with respect to the i 'th component of \mathbf{A} and $\bar{\mathbf{A}}$ and

$$E(\bar{\mathbf{A}}, \mathbf{A}) = \frac{\langle \Psi(\bar{\mathbf{A}}) | H | \Psi(\mathbf{A}) \rangle}{\langle \Psi(\bar{\mathbf{A}}) | \Psi(\mathbf{A}) \rangle}.$$

Eq. (A1) can be interpreted as a Galerkin condition: It forces the residual $(H - E)|\Psi\rangle$ of the full Schrödinger eigenvalue equation – which does not have an exact solution in the variational subset – to be orthogonal to the space spanned by the states $|\partial_i \Psi(\mathbf{A}^*)\rangle$. If the variational subset is a manifold, these states can be interpreted as a basis for the tangent space of the manifold at the point of the variational optimum. Hence, geometrically, the residual has to be orthogonal to the manifold (and thus to its tangent space) at the point of the variational

optimum. Interpreting Eq. (A1) as a Galerkin condition on the ground state eigenvalue problem is useful because it can be generalized to other eigenvalue problems which do not necessarily have a variational characterization (and thus no gradient), as e.g. when the operator is non-Hermitian. Indeed, a similar approach as is developed here was described for finding fixed points of transfer matrices, encoded as matrix product operators, in Ref. 58.

However, before discussing Eq. (A1) in the context of MPS, let us conclude this section by relating it to the time-dependent variational principle (TDVP).²⁷ Geometrically, the TDVP also amounts to an orthogonal projection of the equation of motion (the time-dependent Schrödinger equation) onto the tangent space of the variational manifold. In the case of imaginary time evolution, it can be written as

$$g_{i,j}(\bar{\mathbf{A}}, \mathbf{A}) \frac{d}{dt} A^j = - \langle \partial_i \Psi(\bar{\mathbf{A}}) | H - E(\bar{\mathbf{A}}, \mathbf{A}) | \Psi(\mathbf{A}) \rangle \quad (\text{A2})$$

where

$$g_{i,j}(\bar{\mathbf{A}}, \mathbf{A}) = \langle \partial_i \Psi(\bar{\mathbf{A}}) | \partial_j \Psi(\mathbf{A}) \rangle$$

is the Gram matrix of the tangent vectors and thus the metric of the manifold. The right hand side of Eq. (A2) is again the gradient of the objective function, and the TDVP will thus converge when it reaches a variational optimum \mathbf{A}^* where Eq. (A1) is satisfied. However, the metric $g_{i,j}$ in the left hand side shows that the TDVP equation is not a normal gradient flow, but rather a proper covariant gradient flow that takes the geometry of the manifold and its embedding into the Hilbert space into account. We can thus also associate a quantum state with the gradient, which is given by

$$|\partial_i \Psi\rangle g^{i,\bar{j}} \langle \partial_j \Psi | (H - E) | \Psi \rangle = \mathcal{P}_{T_{|\Psi\rangle}, \mathcal{M}}(H - E) | \Psi \rangle \quad (\text{A3})$$

where we have omitted the arguments \mathbf{A} and $\bar{\mathbf{A}}$, $g^{i,\bar{j}}$ is the inverse of the metric and

$$\mathcal{P}_{T_{|\Psi\rangle}, \mathcal{M}} = |\partial_i \Psi\rangle g^{i,\bar{j}} \langle \partial_j \Psi | \quad (\text{A4})$$

is the projector onto the tangent space at the point $|\Psi\rangle$ in the variational manifold \mathcal{M} . This latter expression is only valid when the variational parameters are proper coordinates for the manifold (i.e. a bijective mapping). While this is not the case for MPS because of gauge freedom (see below), the geometrical interpretation for the Galerkin condition

$$\mathcal{P}_{T_{|\Psi\rangle}, \mathcal{M}}(H - E) | \Psi \rangle = 0 \quad (\text{A5})$$

remains valid; the correct expression for the MPS tangent-space projector will be discussed in more detail in the following section. Independent of whether a variational algorithm is based on a gradient flow, the Hilbert space norm of the gradient $\|\mathcal{P}_{T_{|\Psi\rangle}, \mathcal{M}} H | \Psi \rangle\|$ provides an objective measure for the convergence of the state towards the variational optimum. Note that this is different from the standard Euclidean norm of the naive gradient vector with components $\langle \partial_j \Psi | H - E | \Psi \rangle$.

the orthogonal projection of a general translation invariant state $|\Xi\rangle$ onto the tangent space can be readily found by solving the minimization problem

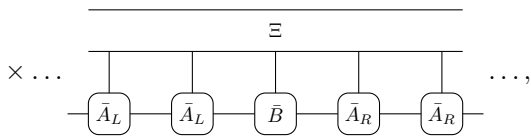
$$\min_B \|\Xi - |\Phi(B)\rangle\|^2,$$

or, equivalently,

$$\min_B (\langle \Phi(B) | \Phi(B) \rangle - \langle \Xi | \Phi(B) \rangle - \langle \Phi(B) | \Xi \rangle).$$

In order to use Eq. (A12) for the first term, we however need to impose the constraint in Eq. (A10) or (A11). In the former case, this will add a term $\text{Tr}[\Lambda \sum_s (A_L^s)^\dagger B^s] + \text{Tr}[\bar{\Lambda} \sum_s (B^s)^\dagger A_L^s]$ to the objective function, with Λ and $\bar{\Lambda}$ corresponding Lagrange multipliers. The solution is readily obtained by demanding $\partial_{\bar{B}}(\dots) = 0$, where Eq. (A12) simply results in $\partial_{\bar{B}} \langle \Phi(B) | \Phi(B) \rangle = |\mathbb{Z}|B$. The overlap between a tangent vector and $|\Xi\rangle$ is given by

$$\langle \Phi(B) | \Xi \rangle = |\mathbb{Z}|$$



By inserting the solution for B back into Eq. (A9), we can read off the tangent space projector. While the value of B depends on the gauge condition, the resulting projector is of course gauge independent and given by

$$\mathcal{P}_{|\Psi(A)\rangle} = \sum_{n \in \mathbb{Z}} \dots \begin{array}{c} \begin{array}{c} \text{---} \bar{A}_L \text{---} \bar{A}_L \text{---} \\ | \quad | \\ \bar{A}_L \text{---} \bar{A}_L \text{---} \\ | \quad | \\ s_{n-2} \quad s_{n-1} \end{array} \quad \Bigg| \quad \begin{array}{c} \begin{array}{c} \text{---} \bar{A}_R \text{---} \bar{A}_R \text{---} \\ | \quad | \\ \bar{A}_R \text{---} \bar{A}_R \text{---} \\ | \quad | \\ s_{n+1} \quad s_{n+2} \end{array} \quad \dots - \dots \\ \begin{array}{c} \begin{array}{c} \text{---} \bar{A}_L \text{---} \bar{A}_L \text{---} \bar{A}_L \text{---} \\ | \quad | \quad | \\ \bar{A}_L \text{---} \bar{A}_L \text{---} \bar{A}_L \text{---} \\ | \quad | \quad | \\ s_{n-2} \quad s_{n-1} \quad s_n \end{array} \quad \begin{array}{c} \begin{array}{c} \text{---} \bar{A}_R \text{---} \bar{A}_R \text{---} \\ | \quad | \\ \bar{A}_R \text{---} \bar{A}_R \text{---} \\ | \quad | \\ s_{n+1} \quad s_{n+2} \end{array} \quad \dots \end{array} \end{array}$$

We can represent the tangent space projector as

$$\mathcal{P}_{|\Psi(A)\rangle} = \sum_{n \in \mathbb{Z}} P_{AC}(n) - P_C(n), \quad (\text{A15})$$

by defining the partial projectors

$$P_{AC}(n) = P_L(n-1) \otimes \mathbb{1}_n \otimes P_R(n+1), \quad (\text{A16a})$$

$$P_C(n) = \mathcal{P}_L(n) \otimes \mathcal{P}_R(n+1), \quad (\text{A16b})$$

$$P_L(n) = \sum_{\alpha} |\Psi_L^{\alpha}(n)\rangle \langle \Psi_L^{\alpha}(n)|, \quad (\text{A16c})$$

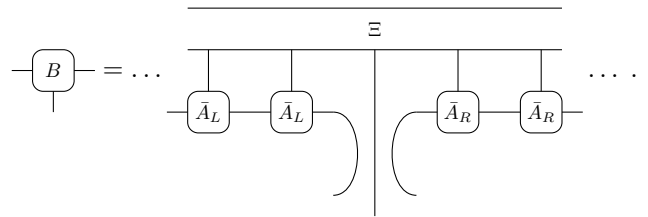
$$P_R(n) = \sum_{\alpha} |\Psi_R^{\alpha}(n)\rangle \langle \Psi_R^{\alpha}(n)|. \quad (\text{A16d})$$

We can verify that $\mathcal{P}_{|\Psi(A)\rangle}^2 = \mathcal{P}_{|\Psi(A)\rangle}$ by using

$$P_L(m) P_L(n) = P_L(\max(m, n)), \quad (\text{A17a})$$

$$P_R(m) P_R(n) = P_R(\min(m, n)). \quad (\text{A17b})$$

so that its derivative $\partial_{\bar{B}} \langle \Phi(B) | \Xi \rangle$ is easily obtained by omitting the tensor \bar{B} from the diagram and interpreting the open legs as defining the indices of a new tensor. Without the Lagrange multiplier, we would simply obtain



With the additional constraint, the solution is still straightforward. It can be easily verified that the correct value of the Lagrange multiplier is such that the additional term acts as a projection

$$B^s \rightarrow B^s - A_L^s \left[\sum_t A_L^t \dagger B^t \right] \quad (\text{A13})$$

or similarly

$$B^s \rightarrow B^s - \left[\sum_t B^t A_R^t \dagger \right] A_R^s \quad (\text{A14})$$

if we would have chosen the right gauge of Eq. (A11).

3. Gradient and Effective Hamiltonians

As discussed in the beginning of this section, a variational optimum can be characterized geometrically as

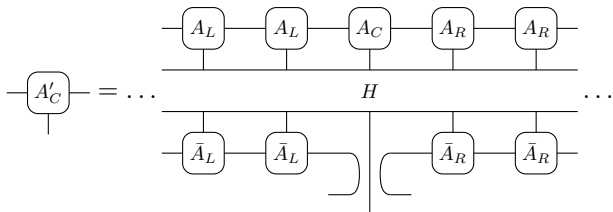
$$\mathcal{P}_{|\Psi(A)\rangle} (H - E) |\Psi(A)\rangle = 0 \quad (\text{A18})$$

where $E = \langle \Psi(A) | H | \Psi(A) \rangle$ (unit normalization is assumed). Since the Galerkin condition is automatically ensured in the direction of the MPS itself, the only non-trivial information of Eq. (A18) is thus contained in the part of the tangent space orthogonal to $|\Psi(A)\rangle$. This is convenient, as $\mathcal{P}_{|\Psi(A)\rangle}$ was actually constructed as the projector onto the part of tangent space orthogonal to $|\Psi(A)\rangle$ in the first place. While this implies that the E subtraction does not contribute, it is convenient to keep it around, as it ensures that the individual terms in the final expression are finite in the thermodynamic limit.

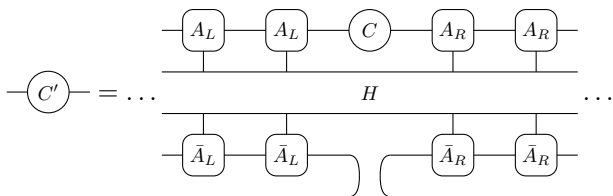
Applying the tangent space projection as in the previous section to the state $|\Xi\rangle = H|\Psi(A)\rangle$ gives rise to a tangent vector of the form Eq. (A9), with

$$B^s = A_C^s - C' A_R^s \quad \text{or} \quad B^s = A_C^s - A_L^s C' \quad (\text{A19})$$

where A_C' originates from applying $P_{A_C}(n)$ and C' from applying $P_C(n)$. By writing $|\Psi(A)\rangle$ itself in a compatible gauge for every individual term, we obtain the diagrammatic expressions



and



We can thus obtain A_C' and C' by acting with the effective Hamiltonians H_{A_C} and H_C introduced in (9) and (10) in the main text onto A_C and C

$$A_C' = H_{A_C} A_C', \quad C' = H_C C. \quad (\text{A20})$$

Even without subtracting the energy, the two choices of B (which are related by the additive gauge transform with $X = C'$) will be finite in the thermodynamic limit. However, the individual tensors A_C' and C' will have a divergent contribution proportional to A_C and C , respectively. Indeed, as discussed in the main text for the case of nearest neighbor interactions, the effective Hamiltonians H_{A_C} and H_C have a divergent contribution corresponding to the total energy times the identity operator. It is thus by subtracting $H \rightarrow \tilde{H} = H - E$ (or $h \rightarrow \tilde{h} = h - e$ for the local terms) that these divergences are canceled. Appendix C provides a detailed description of the construction of the effective Hamiltonian for other types of interactions and illustrates explicitly that the diverging contributions cancel exactly.

A variational extremum is characterized by $|\Phi(B)\rangle = 0$, which leads to $B = 0$ for either gauge choice, as these choices completely fix the gauge freedom. This gives rise to the following simultaneous conditions:

$$A_C^s = A_L^s C' = C' A_R^s \quad (\text{A21})$$

$$A_C^s = A_L^s C = C A_R^s. \quad (\text{A22})$$

However, because the gauge transformation that relates A_L and A_R is unique up to a factor (for injective MPS),

C and C' have to be proportional, and we have actually obtained the eigenvalue equations

$$A_C' = H_{A_C} A_C = E_{A_C} A_C \quad (\text{A23})$$

$$C' = H_C C = E_C C. \quad (\text{A24})$$

As we are looking for a variational minimum, the eigenvalues E_{A_C} and E_C should be the lowest eigenvalues of the effective Hamiltonians H_{A_C} and H_C . Depending on how we have regularized the divergent contributions, these eigenvalues might be different. If we have completely subtracted the energy expectation value from every term, we then have $E_C = E_{A_C} = 0$.

4. Convergence and error measures

While neither VUMPS nor IDMRG or ITEBD directly use the gradient itself, the Hilbert space norm of $|\Phi(B)\rangle$ can be used as an objective convergence measure to indicate how far the current state is from the variational optimum. For either choice of B , we obtain $\|\Phi(B)\| = \sqrt{N}\|B\|$, with N the diverging number of sites and $\|B\|$ the 2-norm of the tensor B . Its square is given by

$$\begin{aligned} \|B\|^2 &= \sum_{s,\alpha,\beta} |B_{\alpha,\beta}^s|^2 \\ &= \sum_s \|A_C^s - A_L^s C'\|^2 \\ &= \sum_s \|A_C^s - C' A_R^s\|^2 \\ &= \|A_C'\|^2 - \|C'\|^2 \end{aligned} \quad (\text{A25})$$

where the equalities follow from $C' = \sum_s (A_L^s)^\dagger A_C^s = \sum_s A_C^s (A_R^s)^\dagger$. Note that none of these expressions are well suited for numerically evaluating the norm close to convergence, as they involve subtracting quantities that are almost equal, especially when the state is close to convergence.

An alternative strategy for evaluating $\|B\|$ is by using the matrix notation for tensors (17) to write $\mathcal{B}^{[\ell]} = \mathcal{A}_C'^{[\ell]} - \mathcal{A}_L C'$. Since \mathcal{A}_L is an isometry, we can extend it to a $dD \times dD$ unitary matrix $U = [\mathcal{A}_L \ \mathcal{N}_L]$, where \mathcal{N}_L contains an orthonormal basis for the $(d-1)D$ -dimensional null space of \mathcal{A}_L^\dagger , i.e. $\mathcal{A}_L^\dagger \mathcal{N}_L = 0$. As the 2-norm is unitarily invariant, we can write

$$\|B\| = \|\mathcal{B}^{[\ell]}\| = \|U^\dagger \mathcal{B}^{[\ell]}\| = \|\mathcal{N}_L^\dagger \mathcal{B}^{[\ell]}\| = \|\mathcal{N}_L^\dagger \mathcal{A}_C'^{[\ell]}\|.$$

The second equality follows from $\mathcal{A}_L^\dagger \mathcal{B}^{[\ell]} = 0$ and the third from the null space property of \mathcal{N}_L . We then obtain $\|B\|$ as the Frobenius norm of a single matrix, which can be calculated accurately as a sum of strictly positive numbers. For further reference, we reshape \mathcal{N}_L into a $D \times d \times (d-1)D$ tensor N_L and similarly introduce a

$(d-1)D \times d \times D$ tensor N_R via the defining relations

$$\sum_s (N_L^s)^\dagger A_L^s = 0, \quad \sum_s (N_L^s)^\dagger N_L^s = \mathbb{1}, \quad (\text{A26})$$

$$\sum_s A_R^s (N_R^s)^\dagger = 0, \quad \sum_s N_R^s (N_R^s)^\dagger = \mathbb{1}. \quad (\text{A27})$$

While the norm of the gradient provides a measure for the quality of approaching the variational minimum, it does not provide any information about the quality of the (u)MPS approximation to the true ground state itself. In the context of (two-site) DMRG schemes, a popular measure is the truncation error, as it is naturally accessible throughout the algorithm (see e.g. Ref. 12, 51, and 52). But also within VUMPS we can compute this quantity by first writing the state $|\Psi(A)\rangle$ in a mixed canonical form with a two-site center block

$$|\Psi(A)\rangle = \sum_{n,\alpha,\beta,s_n,s_{n+1}} (A_{2C})_{\alpha,\beta}^{s_n s_{n+1}} |\Psi_L^\alpha\rangle |s_n\rangle |s_{n+1}\rangle |\Psi_R^\beta\rangle.$$

The two-site center tensor $A_{2C}^{ss'} = A_L^s A_C^s = A_C^s A_R^{s'}$ (known as the two-site wave function $\psi^{ss'}$ in standard DMRG) has an associated effective Hamiltonian $H_{A_{2C}}$. We can compute its lowest eigenvector \tilde{A}_{2C} and compute its singular value decomposition (by first reshaping it to a $dD \times dD$ matrix) $\tilde{A}_{2C}^{ss'} = U^s S V^{s'}$. The truncation error then corresponds to the discarded weight

$$\epsilon_\rho = \sum_{k=D+1}^{dD} S_k^2 \quad (\text{A28})$$

when truncating the inner bond dimension of this two-site tensor to its original value D .

A more generic measure for the error in the variational approximation is given by the energy variance $\langle \Psi(A) | (H - E)^2 | \Psi(A) \rangle = \|(H - E) |\Psi(A)\rangle\|^2$. This quantity is also used in the context of e.g. variational Monte Carlo and various other methods. We can systematically decompose $(H - E) |\Psi(A)\rangle$ into various parts: the projection onto $|\Psi(A)\rangle$ is automatically zero by the definition of E . The projection onto the tangent space is zero when we are at the variational minimum. Next, we can project $(H - E) |\Psi(A)\rangle$ onto the space of all 2-site variations, which is given by states of the form

$$|\Phi^{(2)}(B_2)\rangle = \sum_s \sum_{n \in \mathbb{Z}} (\dots A_L^{s_{n-1}} B_2^{s_n s_{n+1}} A_R^{s_{n+2}} \dots) |s\rangle \quad (\text{A29})$$

In the case of nearest neighbor Hamiltonians, this space captures $(H - E) |\Psi(A)\rangle$ completely, namely by choosing $B_2^{st} = \langle st | \tilde{h} | s't' \rangle A_{2C}^{s't'}$ with A_{2C} the two-site center tensor defined in the previous paragraph and \tilde{h} the local terms of the Hamiltonian, with the current expectation value subtracted, i.e. $H - E = \sum_n \tilde{h}_{n,n+1}$. However, there is again an additive representation redundancy (gauge freedom) $B_2^{st} \rightarrow B_2^{st} + A_L^s X^t - X^s A_R^t$ which enables us to choose representations B satisfying e.g. a left gauge condition

$\sum_s (A_L^s)^\dagger B_2^{st} = 0$ ($\forall t$). The advantage of this representation is again that it facilitates the calculation of the norm, as $\|\Phi^{(2)}(B_2)\|^2 = N \|B_2\|^2$ with N the diverging number of sites. The projection of $(H - E) |\Psi(A)\rangle$ onto this space can be worked out similarly as for the tangent space, and leads to the general result (for any Hamiltonian)

$$B_2^{st} = A_{2C}^{st} - A_L^s A_C^t \quad \text{or} \quad B_2^{st} = A_{2C}^{st} - A_C^s A_R^t \quad (\text{A30})$$

with $A_{2C}^t = H_{A_{2C}} A_{2C}^t$ a single application of the two-site effective Hamiltonian. Using $A_L^s (A_L^s)^\dagger + N_L^s (N_L^s)^\dagger = \delta_{s,t} \mathbb{1}$, we can rewrite the first form of B_2 as

$$B_2^{st} = N_L^s \sum_{s'} (N_L^{s'})^\dagger A_{2C}^{s't}.$$

We now also apply $(A_R^s)^\dagger A_R^t + (N_R^s)^\dagger N_R^t = \delta_{s,t} \mathbb{1}$ to the right hand side and recognize $A_C^s = A_{2C}^{st} (A_R^t)^\dagger$. But since $\sum_s (N_L^s)^\dagger A_C^s = 0$ at the variational minimum, we obtain at the variational minimum

$$B_2^{st} = N_L^s \left[\sum_{s't'} (N_L^{s'})^\dagger A_{2C}^{s't'} (N_R^{t'})^\dagger \right] N_R^t$$

and in particular

$$\|B_2\| = \left\| \sum_{s't'} (N_L^{s'})^\dagger A_{2C}^{s't'} (N_R^{t'})^\dagger \right\|. \quad (\text{A31})$$

We can also relate $\|B_2\|^2$ to the truncation error defined in the previous paragraph. For the truncation error arising in the context of two-site DMRG schemes, the lowest eigenvector \tilde{A}_{2C} of the two-site effective Hamiltonian is used. When the DMRG algorithm has converged, the rank D approximation of \tilde{A}_{2C} should again be the original two-site center tensor $A_{2C}^{st} = A_L^s A_C^t A_R^t$. But this means that we can construct N_L and N_R exactly from the singular vectors corresponding to the $(d-1)D$ singular values that were truncated away, and thus that the truncation error is given by $\epsilon_\rho = \|\sum_{st} (N_L^s)^\dagger \tilde{A}_{2C}^{st} (N_R^t)^\dagger\|^2$. This definition is close to $\|B_2\|^2$, except that the latter uses the tensor A_{2C}' arising from applying the two-site effective Hamiltonian once. As A_{2C} and \tilde{A}_{2C} are anyway close, we can think of A_{2C}' as providing the leading order correction from A_{2C} to \tilde{A}_{2C} in the sense of a Krylov scheme. Indeed, in the first iteration of the Lanczos method, the eigenvector \tilde{A}_{2C} would be approximated in the form $\alpha A_{2C} + \beta A_{2C}'$. Since the first term drops out when projecting onto N_L and N_R , the DMRG truncation error and $\|B_2\|^2$ will be of the same order of magnitude.

Note, however, that $\|B_2\|^2$ only captures the full energy variance (per site) for nearest neighbor Hamiltonians, whose action on $|\Psi(A)\rangle$ is completely contained within the space of two-site variations as noticed above. In that case, we can see that the only term that survives in A_{2C}' after projection onto N_L and N_R is the local term, where \tilde{h} acts on the two-site center tensor. We can thus

also write

$$\|B_2\| = \left\| \sum_{s't'st} \langle s't' | \tilde{h} | st \rangle (N_L^{s'})^\dagger A_{2C}^{st} (N_R^{t'})^\dagger \right\|. \quad (\text{A32})$$

We can also relate this to the truncation step in (I)TEBD, where we would apply $\exp(-\Delta t \tilde{h})$ to every two-site block of the state. The resulting truncation would lead to a discarded weight of the order $\Delta t^2 \|B_2\|^2$.

The considerations regarding the projection of $(H - E)|\Psi(A)\rangle$ onto the space of two-site variations can also be used to devise a scheme for expanding the bond dimension of the uMPS. This approach is presented in the next section.

Appendix B: Dynamic Control of the Bond Dimension

A characteristic feature of two-site implementations of conventional MPS methods – such as e.g. (I)TEBD or (I)DMRG – is that the bond dimension D of the MPS is automatically increased in every iteration and has to be *truncated* in order to remain at a finite maximum bond dimension. This truncation step lies at the basis

We have developed a similar subspace expansion technique that works for a uMPS in the thermodynamic limit. It is based on projecting the full action of the Hamiltonian $(H - E)|\Psi(A)\rangle$ onto the space of two-site variations, as developed in the previous section. There we have found the representation

$$\begin{aligned} B_2^{st} &= N_L^s \left[\sum_{s't'} (N_L^{s'})^\dagger A_C^{s't'} (A_R^{t'})^\dagger \right] A_R^t + N_L^s \left[\sum_{s't'} (N_L^{s'})^\dagger A_{2C}^{s't'} (N_R^{t'})^\dagger \right] N_R^t \\ &= A_C^s A_R^t - A_L^s C^t A_R^t + N_L^s \left[\sum_{s't'} (N_L^{s'})^\dagger A_{2C}^{s't'} (N_R^{t'})^\dagger \right] N_R^t \end{aligned}$$

Even when we have not yet reached the variational minimum, the first term (on line 1) or the first two terms (on line 2) are captured in the tangent space, and only the last term (on either line) contains a new search direction. To capture it completely, we would need to expand the bond dimension from value D to dD . If we want to expand to a new dimension $\tilde{D} = D + \Delta D$, we can use a singular value decomposition to compute the rank ΔD approximation of $\sum_{s't'} (N_L^{s'})^\dagger A_{2C}^{s't'} (N_R^{t'})^\dagger = USV$. By keeping only the largest ΔD singular values, U and V are left and right isometries of size $(d-1)D \times \Delta D$ and $\Delta D \times (d-1)D$ respectively. As remarked in the previous section, in the case of nearest neighbor interactions, the projection of A_{2C}^t onto N_L and N_R does not require the full two-site effective Hamiltonian but reduces to the local term.

We do not directly update the current MPS, but rather write it in an expanded basis in a mixed canonical form with matrices

$$\tilde{A}_L^s = \begin{bmatrix} A_L^s & N_L^s U \\ 0 & 0 \end{bmatrix}, \quad \tilde{A}_R^s = \begin{bmatrix} A_R^s & 0 \\ V^\dagger N_R^s & 0 \end{bmatrix} \tilde{C} = \begin{bmatrix} C & 0 \\ 0 & 0 \end{bmatrix}.$$

With these initial tensors, we can now start a new iteration of VUMPS. Note that we can straightforwardly update the environments used to construct the effective Hamiltonians into this expanded basis, which is necessary if we want to use them as initial guess.

of why such schemes for finding ground states will never truly converge to the variational minimum up to machine precision, as observed in the results. Indeed, even in finite size simulations, two-site DMRG is used to initialize the state and one-site DMRG to obtain final convergence. However, the truncation step in two-site methods has the advantage that the bond dimension can be dynamically increased (or decreased) according to some quality constraint, such as the magnitude of the smallest Schmidt-value or the discarded weight. Especially in the presence of symmetry, this is important to automatically obtain the correct symmetry sectors within the virtual MPS space.

The VUMPS algorithm presented in the main text is variational from the start and therefore works at fixed bond dimension D , i.e. it is a one-site scheme in DMRG terminology. Alternative subspace expansion strategies for dynamically increasing the bond dimension in such one-site schemes have been proposed.^{22,60} These methods use information from acting with the global Hamiltonian onto the current state to either add a tiny perturbation to the current MPS or to generate a larger basis in which the effective eigenvalue problem is solved.

Appendix C: Explicit Construction of Effective Hamiltonians

In this section we describe how to efficiently apply the effective Hamiltonians H_{A_C} and H_C onto the center site tensor A_C^s and bond matrix C and how the necessary individual terms are explicitly constructed. Such a procedure is needed for solving the effective eigenvalue problems (23a) and (23b) by means of an iterative eigensolver. The case of systems with nearest neighbor interaction has already been discussed in Sec. II B. In the following we

consider the cases of Hamiltonians with long range interactions in Sec. C1 and general Hamiltonians given in terms of Matrix Product Operators (MPOs) in Sec. C2.

1. Long Range Interactions

Consider Hamiltonians with long range interactions of the form $H = \sum_{j \in \mathbb{Z}} h_j$, where h_j is itself an infinite sum

$$h_j = \sum_{n>0} f(n) o_j o_{j+n} \quad (\text{C1})$$

and operators o_i act on a single-site i and commute when acting on different sites $[o_i, o_j] = 0$, $i \neq j$. Without loss of generality, we restrict to a single pair of (bounded) operators o , which commute when acting on different sites $[o_i, o_j] = 0$, $i \neq j$ ⁶¹. The generalization to Hamiltonians containing several terms of that form is straight forward. Furthermore, we assume distance functions $f(n)$ that are bounded in the sense of $\sum_{n>0} |f(n)| < \infty$, such that $\|h_j\| < \infty$, and that can be well approximated by a sum of K exponentials, i.e.

$$f(n) \approx \sum_{k=1}^K c_k \lambda_k^{n-1}, \quad (\text{C2})$$

with $|\lambda_k| < 1$ and $n > 0$. In practice, for an infinite system we fit $f(n)$ with a suitable number of K exponentials over a distance N large enough, such that $f(N)$ and the largest residuals are below some desired threshold.

Examples of Hamiltonians that fall in this class are the transverse field Ising (TFI) model or XXZ model with power-law interactions,⁶²⁻⁶⁴ as well as the famous Haldane-Shastry model,^{42,43} for which the ground state is exactly known.

Similar to the case of nearest neighbor interactions in Sec. IIB, the effective Hamiltonians factorize into a number of terms which can all be applied efficiently. For H_{AC} these are the five terms, out of which four are already familiar from the case of nearest neighbor interactions. Two of these are the left and right block Hamiltonians H_L and H_R with infinitely many local contributions from h_j acting on sites strictly left or right of the current center site, and the other two are the terms containing interactions between the center site and the left and right block respectively, i.e. where h_j partially acts on A_C . For long range interactions we have one additional term, containing infinitely many interaction terms between the left and the right block only without involving the center site, i.e. where o_j acts to the left of the current center site, and o_{j+n} acts to the right.

To construct all these terms we start by defining the operator transfer matrices

$$T_L^{[o]} = \sum_{st} o_{st} \bar{A}_L^s \otimes A_L^t \quad T_R^{[o]} = \sum_{st} o_{st} \bar{A}_R^s \otimes A_R^t. \quad (\text{C3})$$

The current energy density expectation value $e = \langle \Psi(A) | h | \Psi(A) \rangle$ can thus be written as

$$\begin{aligned} e &= (\mathbb{1} | T_L^{[o]} \left[\sum_{n>0} f(n) (T_L)^{n-1} \right] T_L^{[o]} | R) \\ &= (L | T_R^{[o]} \left[\sum_{n>0} f(n) (T_R)^{n-1} \right] T_R^{[o]} | \mathbb{1}), \end{aligned} \quad (\text{C4})$$

or using (C2)

$$\begin{aligned} e &= \sum_k c_k (\mathbb{1} | T_L^{[o]} \left[\sum_{n \geq 0} (\lambda_k T_L)^n \right] T_L^{[o]} | R) \\ &= \sum_k c_k (L | T_L^{[o]} \left[\sum_{n \geq 0} (\lambda_k T_R)^n \right] T_L^{[o]} | \mathbb{1}). \end{aligned} \quad (\text{C5})$$

Since $|\lambda_k| < 1$ the geometric series converge and we can perform them explicitly. We proceed by defining

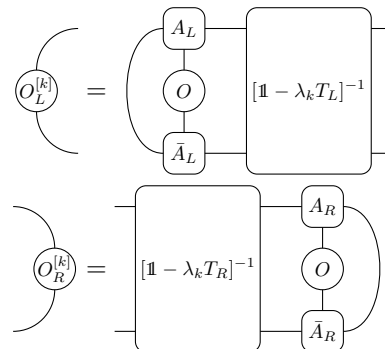
$$\begin{aligned} (O_L^{[k]} | &= (\mathbb{1} | T_L^{[o]} [\mathbb{1} - \lambda_k T_L]^{-1} \\ | O_R^{[k]}) &= [\mathbb{1} - \lambda_k T_R]^{-1} T_L^{[o]} | \mathbb{1}). \end{aligned} \quad (\text{C6})$$

These terms can again either be calculated recursively by explicitly evaluating the geometric sums term by term until convergence, or more efficiently by iteratively solving the following systems of linear equations

$$\begin{aligned} (O_L^{[k]} | [\mathbb{1} - \lambda_k T_L] &= (\mathbb{1} | T_L^{[o]} \\ [\mathbb{1} - \lambda_k T_R] | O_R^{[k]}) &= T_L^{[o]} | \mathbb{1} \end{aligned} \quad (\text{C7})$$

using iterative methods.

We represent these terms by the diagrams



and collect all such terms into single left and right environment contributions

$$(O_L | = \sum_k c_k (O_L^{[k]} | \quad | O_R) = \sum_k c_k | O_R^{[k]}) \quad (\text{C8})$$

and further

$$(h_L | = (O_L | T_L^{[o]} \quad | h_R) = T_R^{[o]} | O_R). \quad (\text{C9})$$

We can then write for the energy density

$$e = (h_L | R) = (L | h_R). \quad (\text{C10})$$

Comparing with (C4) we have thus defined

$$\begin{aligned} |h_L\rangle &= (\mathbb{1}|T_L^{[o]} \left[\sum_{n>0} f(n)(T_L)^{n-1} \right] T_L^{[o]} \\ |h_R\rangle &= T_R^{[o]} \left[\sum_{n>0} f(n)(T_R)^{n-1} \right] T_R^{[o]} |\mathbb{1}\rangle. \end{aligned} \quad (\text{C11})$$

With these definitions at hand we can write the left

We are now ready to formulate the action of H_{A_C} onto A_C^s as

$$A_C^s = H_L A_C^s + A_C^s H_R + O_L \left[\sum_t o_t^s A_C^t \right] + \left[\sum_t o_t^s A_C^t \right] O_R + \sum_k c_k \lambda_k O_L^{[k]} A_C^s O_R^{[k]} \quad (\text{C13})$$

The additional factor of λ_k in the sum in the last term arises due to A_C^s adding an additional site between the left and right operators o . Similarly, the action of H_C onto C becomes

$$C' = H_L C + C H_R + \sum_k c_k O_L^{[k]} C O_R^{[k]} \quad (\text{C14})$$

In (C13) the first two terms can be applied in $\mathcal{O}(dD^3)$, the second two in $\mathcal{O}(d^2D^2) + \mathcal{O}(dD^3)$ and the last term in $\mathcal{O}(KdD^3)$ operations, and in (C14) the first two terms in $\mathcal{O}(D^3)$ and the last term in $\mathcal{O}(KD^3)$ operations. In general we have to perform $2(K+1)$ iterative inversions involving $\mathcal{O}(D^3)$ operations and collect K terms to arrive at the necessary terms for (C13) and (C14), where the solutions from the previous iteration can be used as starting vectors to speed up convergence.

If there are additional simple single or nearest neighbor two-site terms present in the Hamiltonian, appropriate terms as described in Sec. IIB can be added. For a pseudocode summary for obtaining the necessary explicit terms of H_{A_C} and H_C for Hamiltonians with long range interactions, and their applications onto a state, required for solving the effective eigenvalue problems using an iterative eigensolver, see Table IV.

and right block Hamiltonians as

$$(H_L | = \langle h_L | \sum_{n=0}^{\infty} [T_L]^n \quad |H_R\rangle = \sum_{n=0}^{\infty} [T_L]^n |h_R\rangle. \quad (\text{C12})$$

These equations are exactly the same as Eq. (13) for the case of nearest neighbor interactions, but with different $\langle h_L |$ and $|h_R\rangle$. We can thus evaluate the geometric sums recursively or by solving a linear system iteratively, as explained in Sec. IIB. Note that we again start by applying an energy shift $\langle h_L | \rightarrow \langle \tilde{h}_L | = (\langle h_L | - e |R\rangle) (\mathbb{1} |$ and similar for $|h_R\rangle$, such that $\langle \tilde{h}_L | R\rangle = \langle L | \tilde{h}_R\rangle = 0$.

2. General Hamiltonians given in terms of MPOs

Consider the Hamiltonian H given in terms of an infinite Matrix Product Operator (MPO) ^{23,34,65-69} with 4-index MPO elements $W_{ss'}^{ab}$ with $a, b = 1, \dots, d_W$ and $s, s' = 1, \dots, d$ and we call d_W the MPO bond dimension. In terms of the operator valued matrices $\hat{W}^{ab} = \sum_{ss'} W_{ss'}^{ab} |s\rangle \langle s'|$ the Hamiltonian can then be written as

$$H = \hat{w}_L \left[\prod_{j \in \mathbb{Z}} \hat{W}_{[j]} \right] \hat{w}_R$$

Algorithm 5 Explicit terms of effective Hamiltonians with long range interactions and their application onto a state

Input: operator o defining (C1), parameters c_k and λ_k defining (C2), current uMPS tensors A_L, A_R in left and right gauge, left dominant eigenvector $|L\rangle$ of T_R , right dominant eigenvector $|R\rangle$ of T_L , desired precision ϵ_S for terms involving infinite geometric sums

Output: Explicit terms of effective Hamiltonians H_{AC} and H_C , updated A'_C and C'

```

1: function HEFFTERMS( $H = \{o, \{c_k\}, \{\lambda_k\}\}, A_L, A_R, L, R, \epsilon_S$ )           ▷ Calculates explicit terms of effective Hamiltonians
2:   Calculate  $O_L^{[k]}$  and  $O_R^{[k]}$  by iteratively solving (C7) for each  $\lambda_k$  to machine precision
3:   Calculate single environment contributions  $O_L$  and  $O_R$  from (C8) and  $h_L$  and  $h_R$  from (C9)
4:   Calculate  $H_L$  and  $H_R$  by iteratively solving (14) or (preferably) (15), to precision  $\epsilon_S$ 
5:    $H_{AC} \leftarrow \{o, \{c_k\}, \{\lambda_k\}, \{O_L^k\}, \{O_R^k\}, O_L, O_R, H_L, H_R\}$ 
6:    $H_C \leftarrow \{\{c_k\}, \{O_L^k\}, \{O_R^k\}, H_L, H_R\}$ 
7:   return  $H_{AC}, H_C$ 
8: end function
9: function APPLYHAC( $A_C, H_{AC}$ )                                           ▷ Terms of  $H_{AC}$  from HEFFTERMS( $H, A_L, A_R, L, R, \epsilon_S$ )
10:  Calculate updated  $A'_C$  from (C13)
11:  return  $A'_C$ 
12: end function
13: function APPLYHC( $C, H_C$ )                                             ▷ Terms of  $H_C$  from HEFFTERMS( $H, A_L, A_R, L, R, \epsilon_S$ )
14:  Calculate updated  $C'$  from (C14)
15:  return  $C'$ 
16: end function

```

Table IV. Pseudocode for obtaining the explicit terms of the effective Hamiltonians H_{AC} and H_C for systems with with long range interactions and their applications onto a state.

where $\hat{W}_{[j]}$ contains operators acting on site j only and \hat{w}_L and \hat{w}_R are operator valued boundary vectors.

An example for such an MPO decomposition for the Transverse Field Ising (TFI) Hamiltonian with exponentially decaying long range interaction

$$H_{\text{TFI}} = -J \sum_j \sum_{n>0} \lambda^{n-1} X_j X_{j+n} - h \sum_j Z_j$$

with $\lambda < 1$ is given by

$$\begin{aligned} \hat{W} &= \begin{bmatrix} \mathbb{1} & 0 & 0 \\ -JX & \lambda \mathbb{1} & 0 \\ -hZ & X & \mathbb{1} \end{bmatrix} \\ \hat{w}_L &= [-hZ \quad X \quad \mathbb{1}] \\ \hat{w}_R &= [\mathbb{1} \quad -JX \quad -hZ]^T, \end{aligned} \quad (\text{C15})$$

where X and Z are Pauli matrices. For the TFI Hamiltonian we thus have $d_W = 3$ and the limit $\lambda = 0$ corresponds to the nearest neighbor interaction case.

In order to efficiently apply the effective Hamiltonians H_{AC} and H_C , it is necessary to determine the left and right (quasi) fixed points $L_a^{[W]}$ and $R_a^{[W]}$ of the MPO transfer matrices

$$T_{L/R}^{[W]ab} = \sum_{ss'} W_{s's}^{ab} \bar{A}_{L/R}^{s'} \otimes A_{L/R}^s, \quad (\text{C16})$$

where – similar to MPS tensors – $L_a^{[W]}$ and $R_a^{[W]}$ are collections of d_W matrices of dimension $D \times D$, with $a = 1, \dots, d_W$. These two objects are in fact the thermodynamic limit versions of the objects defined in Eq. (190) and (191) in Ref. 12.

Typically, MPO representations \hat{W}^{ab} of (quasi)local Hamiltonians (such as e.g. Eq. (C15)) are of Schur form,⁶⁹ such that the MPO transfer matrix contains Jordan blocks and that the dominant eigenvalue is one and of twofold algebraic degeneracy. Such MPO transfer matrices therefore technically do not have well defined fixed points. We can however find quasi fixed points $L_a^{[W]}$ and $R_a^{[W]}$, that are fixed points up to a term contributing to the energy density expectation value in one of the d_W elements of $L_a^{[W]}$ and $R_a^{[W]}$. An application of $T_{L/R}^{[W]ab}$ onto both quasi fixed points will therefore accumulate an additional term contributing to the extensive global energy expectation value. Similar to the terms in Eq. (13) or Eq. (C12) in the previous cases involving infinite geometric sums, we can however safely discard these diverging contributions, which is equivalent to setting the energy expectation values of the semi-infinite left and right half of the system to zero (see below).

In the following we briefly reiterate the procedure of Ref. 69 to systematically determine $L_a^{[W]}$ and $R_a^{[W]}$ from given \hat{W}^{ab} and $A_{L/R}^s$. The obtained solutions will of course contain the results of Sec. IIB and Sec. C1 as special cases.

Without loss of generality we assume \hat{W}^{ab} to be of lower triangular form, i.e. $\hat{W}^{ab} = 0, \forall b > a$. Furthermore, we assume the typical case of any nonzero diagonal elements being proportional to the identity, i.e. $\hat{W}^{aa} = \lambda_a \mathbb{1}$, where $\lambda_a \leq 1$ and usually $\lambda_1 = \lambda_{d_W} = 1$, as is e.g. the case in (C15). By defining the result of the action of the

MPO transfer matrix as

$$(Y_{L_a}| = \sum_{b>a} (L_b^{[W]}|T_L^{[W]ba} \quad (\text{C17})$$

$$|Y_{R_a}) = \sum_{b<a} T_R^{[W]ab}|R_b^{[W]}), \quad (\text{C18})$$

the system of fixed point equations can be written as

$$(L_a^{[W]}| = (L_a^{[W]}|T_L^{[W]aa} + (Y_{L_a}| \quad (\text{C19})$$

$$|R_a^{[W]} = T_R^{[W]aa}|R_a^{[W]} + |Y_{R_a}). \quad (\text{C20})$$

Notice that due to the lower triangular structure of \hat{W}^{ab} , the terms $(Y_{L_a}|$ and $|Y_{R_a})$ only contain contributions from $(L_{b>a}^{[W]}|$ and $|R_{b<a}^{[W]})$ and we can solve (C19) and (C20) recursively, starting with $a = d_W$ for $L_a^{[W]}$ and with $a = 1$ for $R_a^{[W]}$, which initially amounts to $(L_{d_W}^{[W]}| = |\mathbb{1}|$ and $|R_1^{[W]} = |\mathbb{1}|$. Terms with $T_{L/R}^{[W]aa} = 0$ are particularly simple and simply reduce to the identification $(L_a^{[W]}| = (Y_{L_a}|$ and $|R_a^{[W]} = |Y_{R_a})$.

Terms with $T_{L/R}^{[W]aa} = \lambda_a T_{L/R}$ where $\lambda_a < 1$, now result in solutions of the form

$$(L_a^{[W]}| = (Y_{L_a}|[\mathbb{1} - \lambda_a T_L]^{-1} \quad (\text{C21})$$

$$|R_a^{[W]} = [\mathbb{1} - \lambda_a T_R]^{-1}|Y_{R_a}), \quad (\text{C22})$$

equivalent to terms such as (C6) stemming from infinite geometric sums of (weighted) MPS transfer matrices.

Equivalently, terms with $T_{L/R}^{[W]aa} = T_{L/R}$ then result in relations of the form

$$(L_a^{[W]}|[\mathbb{1} - T_L] = (Y_{L_a}| \quad (\text{C23})$$

$$[\mathbb{1} - T_R]|R_a^{[W]} = |Y_{R_a}), \quad (\text{C24})$$

which in general do not have a formal solution, since the left hand sides of these equations live in the subspace orthogonal to the dominant eigenspaces of $T_{L/R}$, while the right hand sides generally do have contributions in the dominant eigenspace. We can however discard these contributions by projecting onto the complementary subspace, and then obtain $(L_a^{[W]}|$ and $|R_a^{[W]})$ by solving the systems of equations (see also Appendix D)

$$(L_a^{[W]}|[\mathbb{1} - T_L + |R)(\mathbb{1}] = (Y_{L_a}| - (Y_{L_a}|R)(\mathbb{1}) \quad (\text{C25a})$$

$$[\mathbb{1} - T_R + |\mathbb{1})(L)|R_a^{[W]} = |Y_{R_a}) - |\mathbb{1})(L|Y_{R_a}) \quad (\text{C25b})$$

We have encountered exactly the same type of equations in (15) when evaluating infinite geometric sums of transfer matrices, after a constant shift in energy to remove diverging terms. The MPO formalism thus automatically yields these contributions in a form where the sums have already been explicitly performed.

Such a situation typically occurs only for the final terms in the recursive solution of the fixed point equations, i.e. for $(L_1^{[W]}|$ and $|R_{d_W}^{[W]})$. A concrete evaluation

(see below) of the discarded terms in these cases shows that they correspond to contributions to the energy density expectation value, i.e. discarding these terms is equivalent to a constant shift in energy, such that the energy density is zero and we have $(L_1^{[W]}|R) = (L|R_{d_W}^{[W]}) = 0$. After applying $T_{L/R}^{[W]}$ once onto the quasi fixed points we thus have for the first element of $L^{[W]}$ and the last element of $R^{[W]}$

$$(Y_{L_1}| = (L_1^{[W]}| + (Y_{L_1}|R)(\mathbb{1}) \quad (\text{C26})$$

$$|Y_{R_{d_W}} = |R_{d_W}^{[W]} + |\mathbb{1})(L|Y_{R_{d_W}}),$$

i.e. the fixed point relations only hold up to an additive diagonal correction for these elements. These corrections correspond to the energy density expectation value

$$e = (Y_{L_1}|R) = (L|Y_{R_{d_W}}) \quad (\text{C27})$$

and they can in fact be used for its evaluation.

As a concrete example, for the long range TFI Hamiltonian given by MPO (C15) we obtain

$$(L_1^{[W]}|[\mathbb{1} - T_L] = -h(\mathbb{1}|T_L^Z - J(\mathbb{1}|T_L^X[\mathbb{1} - \lambda T_L]^{-1}T_L^X$$

$$(L_2^{[W]}| = (\mathbb{1}|T_L^X[\mathbb{1} - \lambda T_L]^{-1}$$

$$(L_3^{[W]}| = (\mathbb{1}|$$

and

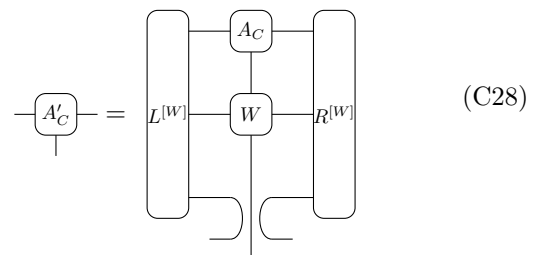
$$|R_1^{[W]} = |\mathbb{1})$$

$$|R_2^{[W]} = -J[\mathbb{1} - \lambda T_R]^{-1}T_R^X|\mathbb{1})$$

$$[\mathbb{1} - T_R]|R_3^{[W]} = -hT_R^Z|\mathbb{1}) - JT_R^X[\mathbb{1} - \lambda T_R]^{-1}T_R^X|\mathbb{1}).$$

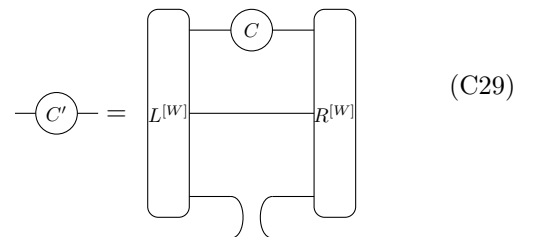
Having determined the left and right quasi fixed points of the MPO transfer matrices, it is now particularly easy to calculate the action of the effective Hamiltonians H_{A_C} onto A_C^s as

$$A_C^s = \sum_{abt} W_{st}^{ab} L_a^{[W]} A_C^t R_b^{[W]}$$



and equivalently the action of H_C onto C as

$$C' = \sum_a L_a^{[W]} C R_a^{[W]}$$



Algorithm 6 Explicit terms of effective Hamiltonians in MPO form and their application onto a state

Input: MPO \hat{W} defining the Hamiltonian, current uMPS tensors A_L , A_R in left and right gauge, left dominant eigenvector $|L\rangle$ of T_R , right dominant eigenvector $|R\rangle$ of T_L , desired precision ϵ_S for iterative solution of linear system of equations

Output: Explicit terms of effective Hamiltonians H_{A_C} and H_C , updated A'_C and C'

```

1: function HEFFTERMS( $H = \hat{W}, A_L, A_R, L, R, \epsilon_S$ ) ▷ Calculates explicit terms of effective Hamiltonians
2:    $L^{[W]} \leftarrow \text{CALCLW}(\hat{W}, A_L, R, \epsilon_S)$ 
3:    $R^{[W]} \leftarrow \text{CALCRW}(\hat{W}, A_R, L, \epsilon_S)$ 
4:    $H_{A_C} \leftarrow \{\hat{W}, L^{[W]}, R^{[W]}\}$ 
5:    $H_C \leftarrow \{L^{[W]}, R^{[W]}\}$ 
6:   return  $H_{A_C}, H_C$ 
7: end function
8: function CALCLW( $\hat{W}, A_L, R, \epsilon_S$ ) ▷ Calculates left quasi fixed point of MPO transfer matrix  $T_L^{[W]}$ 
9:    $(L_{d_w}^{[W]} | \leftarrow |\mathbb{1}\rangle$ 
10:  for  $a = d_w - 1, \dots, 1$  do
11:    Calculate  $(Y_{La} |$  from (C17)
12:    if  $T_L^{[W]aa} == \lambda_a T_L$  then
13:      Calculate  $(L_a^{[W]} |$  by iteratively solving (C21) to machine precision
14:    else if  $T_L^{[W]aa} == T_L$  then
15:      Calculate  $(L_a^{[W]} |$  by iteratively solving (C25a) to precision  $\epsilon_S$ 
16:    else if  $T_L^{[W]aa} == 0$  then
17:       $(L_a^{[W]} | \leftarrow (Y_{La} |$ 
18:    end if
19:  end for
20:  return  $L^{[W]}$ .
21: end function
22: function CALCRW( $\hat{W}, A_R, L, \epsilon_S$ ) ▷ Calculate right quasi fixed point of MPO transfer matrix  $T_R^{[W]}$ 
23:    $|R_1^{[W]} \rangle \leftarrow |\mathbb{1}\rangle$ 
24:   for  $a = 2, \dots, d_w$  do
25:     Calculate  $|Y_{Ra}\rangle$  from (C18)
26:     if  $T_R^{[W]aa} == \lambda_a T_R$  then
27:       Calculate  $|R_a^{[W]}\rangle$  by iteratively solving (C22) to machine precision
28:     else if  $T_R^{[W]aa} == T_R$  then
29:       Calculate  $|R_a^{[W]}\rangle$  by iteratively solving (C25b) to precision  $\epsilon_S$ 
30:     else if  $T_R^{[W]aa} == 0$  then
31:        $|R_a^{[W]}\rangle \leftarrow |Y_{Ra}\rangle$ 
32:     end if
33:   end for
34:   return  $R^{[W]}$ .
35: end function
36: function APPLYHAC( $A_C, H_{A_C}$ ) ▷ Terms of  $H_{A_C}$  from HEFFTERMS( $H, A_L, A_R, L, R, \epsilon_S$ )
37:   Calculate updated  $A'_C$  from (C28)
38:   return  $A'_C$ 
39: end function
40: function APPLYHC( $C, H_C$ ) ▷ Terms of  $H_C$  from HEFFTERMS( $H, A_L, A_R, L, R, \epsilon_S$ )
41:   Calculate updated  $C'$  from (C29)
42:   return  $C'$ 
43: end function

```

Table V. Pseudocode for obtaining the explicit terms of the effective Hamiltonians H_{A_C} and H_C for general Hamiltonians in MPO form and their applications onto a state.

which can be performed in $\mathcal{O}(dd_W D^3) + \mathcal{O}(d^2 d_W^2 D^2)$, respective $\mathcal{O}(d_W D^3)$ operations. In total we also have to perform an iterative inversion for each diagonal element of \hat{W} .

This framework is very flexible, general and powerful, once a routine for determining quasi fixed points of general MPO transfer matrices has been implemented. The effective Hamiltonians of Sec. II B and Sec. C 1 are contained within as special cases. A pseudocode summary for obtaining the necessary explicit terms of H_{AC} and H_C for Hamiltonians given in terms of an MPO, and their applications onto a state, required for solving the effective eigenvalue problems using an iterative eigensolver, is presented in Table V.

Appendix D: Geometric Sums of Transfer Matrices

We wish to evaluate terms involving infinite geometric sums of the form

$$(y| = (x| \sum_{n=0}^{\infty} T^n \quad |y) = \sum_{n=0}^{\infty} T^n |x). \quad (\text{D1})$$

Such expressions typically arise in situations where one sums up contributions of successive applications of T onto some fixed virtual boundary vector x , with the initial contribution being the boundary vector x itself. This is reflected in the above expression by summing from $n = 0$ and using the definition $T^0 = \mathbb{1}$.

We assume a spectral decomposition of the transfer matrix given by

$$T = \sum_{j=0}^{D^2-1} \lambda_j |j)(j|, \quad (\text{D2})$$

where the left and right eigenvectors are mutually orthonormal, i.e. $(j|k) = \delta_{jk}$. Note that T is in general not hermitian and thus $(j| \neq |j)^\dagger$.

For a generic injective normalized state, T has a unique eigenvalue of largest magnitude given by $\lambda_0 = 1$, whereas all other eigenvalues are contained in the unit circle ($|\lambda_{j>0}| < 1$).

We divide into dominant and complementary subspaces and get for powers of T

$$T^n = |0)(0| + \sum_{j=1}^{D^2-1} \lambda_j^n |j)(j|. \quad (\text{D3})$$

We can safely perform the geometric sum for all eigenvalues $|\lambda_{j>0}| < 1$, while $\lambda_0 = 1$ contributes a formally

diverging term

$$\sum_{n=0}^{\infty} T^n = \sum_{n=0}^{\infty} |0)(0| + \sum_{j=1}^{D^2-1} \sum_{n=0}^{\infty} \lambda_j^n |j)(j| \quad (\text{D4})$$

$$= |\mathbb{N}| |0)(0| + \sum_{j=1}^{D^2-1} (1 - \lambda_j)^{-1} |j)(j|. \quad (\text{D5})$$

The interpretation of this diverging contribution depends on the situation. By using the projectors

$$P = |0)(0| \quad Q = \mathbb{1} - |0)(0| \quad (\text{D6})$$

onto the dominant and complementary subspaces we define the projected transfer matrix

$$\mathcal{T} = \sum_{j=1}^{D^2-1} \lambda_j |j)(j| = QT = TQ = T - P. \quad (\text{D7})$$

We realize that the spectral decomposition of $(\mathbb{1} - \mathcal{T})^{-1}$ has a component of $|0)(0|$

$$(\mathbb{1} - \mathcal{T})^{-1} = |0)(0| + \sum_{j=1}^{D^2-1} (1 - \lambda_j)^{-1} |j)(j| \quad (\text{D8})$$

and therefore identify the second term in (D5) as

$$\sum_{j=1}^{D^2-1} (1 - \lambda_j)^{-1} |j)(j| = Q(\mathbb{1} - \mathcal{T})^{-1}Q. \quad (\text{D9})$$

For the geometric sum we then obtain

$$\sum_{n=0}^{\infty} T^n = |\mathbb{N}| |0)(0| + Q(\mathbb{1} - \mathcal{T})^{-1}Q \quad (\text{D10})$$

with a diverging contribution from P . Plugging into (D1) we finally get

$$\begin{aligned} (y| &= |\mathbb{N}| (x|0)(0| + (x|Q(\mathbb{1} - \mathcal{T})^{-1} \\ |y) &= |\mathbb{N}| |0)(0|x) + (\mathbb{1} - \mathcal{T})^{-1}Q|x). \end{aligned} \quad (\text{D11})$$

Usually it is not necessary to calculate the full matrix expression of $\sum_n T^n$, but to just act with it onto some $(x|$ or $|x)$. The diverging contributions can typically be safely discarded, as they correspond to a constant (albeit infinite) offset of some extensive observable (e.g. the Hamiltonian). The action of the finite remaining part can be calculated efficiently by iteratively solving the linear system of equations of the type $A\vec{y} = \vec{x}$ or $\vec{y}^\dagger A = \vec{x}^\dagger$

$$\begin{aligned} (y|(\mathbb{1} - \mathcal{T}) &= (x|Q \\ (\mathbb{1} - \mathcal{T})|y) &= Q|x) \end{aligned} \quad (\text{D12})$$

with inhomogeneities $\vec{x} = Q|x)$ and $\vec{x}^\dagger = (x|Q$. One can then efficiently compute $|y)$ and $(y|$ by employing an iterative Krylov subspace method such as `bicgstab`⁷⁰ or `gmres`.⁷¹ For such methods only the implementation of a (left or right) action of $(\mathbb{1} - \mathcal{T})$ onto a vector is necessary, which can be done efficiently with $\mathcal{O}(dD^3)$ operations. If the transfer matrix is in left or right canonical form, we recover the linear systems in Eqs. (15) and (C25)

$$\begin{aligned} (y|[\mathbb{1} - T_L + |R)(\mathbb{1}] &= (x| - (x|R)(\mathbb{1}] \\ [\mathbb{1} - T_R + |\mathbb{1})(L]|y) &= |x) - |\mathbb{1})(L|x). \end{aligned} \quad (\text{D13})$$