# Topological quantum order: stability under local perturbations

Sergey Bravyi<sup>\*</sup>, Matthew Hastings<sup>†</sup>, and Spyridon Michalakis<sup>‡</sup>

January 3, 2010

#### Abstract

We study zero-temperature stability of topological phases of matter under weak timeindependent perturbations. Our results apply to quantum spin Hamiltonians that can be written as a sum of geometrically local commuting projectors on a D-dimensional lattice with certain topological order conditions. Given such a Hamiltonian  $H_0$  we prove that there exists a constant threshold  $\epsilon > 0$  such that for any perturbation V representable as a sum of short-range bounded-norm interactions the perturbed Hamiltonian  $H = H_0 + \epsilon V$  has well-defined spectral bands originating from O(1) smallest eigenvalues of  $H_0$ . These bands are separated from the rest of the spectrum and from each other by a constant gap. The band originating from the smallest eigenvalue of  $H_0$  has exponentially small width (as a function of the lattice size).

Our proof exploits a discrete version of Hamiltonian flow equations, the theory of relatively bounded operators, and the Lieb-Robinson bound.

<sup>\*</sup>IBM Watson Research Center, Yorktown Heights NY 10594 (USA); sbravyi@us.ibm.com

<sup>&</sup>lt;sup>†</sup>Microsoft Research Station Q, CNSI Building, University of California, Santa Barbara, CA, 93106 (USA); mahastin@microsoft.com

<sup>&</sup>lt;sup>‡</sup>T-4 and CNLS, LANL - Los Alamos, NM, 87544 (USA); spiros@lanl.gov

# Contents

1	Introduction	3
	1.1 Summary of results	4
	1.2 Sketch of the stability proof	
<b>2</b>	Hamiltonians describing TQO	7
	2.1 Frustration-free commuting Hamiltonians	7
		8
	2.3 Verification of TQO conditions for stabilizer Hamiltonians	10
	2.4 Unstable version of the toric code model	11
3	Relatively bounded perturbations	<b>12</b>
	3.1 Definition and basic properties	12
	3.2 Stability of TQO under block-diagonal perturbations	13
4	Hamiltonian flow equations	16
	4.1 Outline of the method	16
	4.2 Local decompositions of Hamiltonians	18
	4.3 Proof of the main theorem	20
5	Linearized block-diagonalization problem	23
	5.1 Statement of the problem	23
	5.2 Finding the transformation $S$	23
	5.3 Local decomposition of the transformed Hamiltonian	28
6	Lieb-Robinson bounds	29
7	Adiabatic continuation of logical operators	34
	7.1 Dressed Operators	34

## 1 Introduction

The traditional classification of different phases of matter due to Landau rests on symmetry breaking. Given a pair of gapped Hamiltonians  $H_1$ ,  $H_2$  with some symmetry group G, the ground states of  $H_1$  and  $H_2$  were considered to be in different phases if their symmetry breaking patterns are different. The discovery of topologically ordered phases, however, changes this paradigm. Models such as Kitaev's toric code [1] have "topologically non-trivial" ground states despite lacking any symmetry breaking. Such states cannot be changed into a "topologically trivial" state such as a product state by any unitary locality-preserving operator [2].

One possible approach to classifying topological phases is to call a pair of gapped Hamiltonians  $H_1, H_2$  topologically equivalent iff it is possible to connect  $H_1$  and  $H_2$  by a continuous path in the space of local gapped Hamiltonians. Using the idea of quasi-adiabatic continuation [3], one can describe the evolution of the ground state subspace along such a path by a unitary locality-preserving operator. In particular ground state degeneracy and the geometry of "logical operators" acting on the ground subspace is the same for  $H_1$  and  $H_2$ .

Most of the Hamiltonians describing TQO models such as Kitaev's quantum double model [1] or Levin-Wen string-net model [4] are not quite physical since they involve interactions affecting more than two spins at a time. One may hope however that such models emerge as low-energy effective Hamiltonians describing some simpler high-energy theories [5, 6, 7]. For example, the toric code model with four-spin interactions can be "implemented" as the fourth-order effective Hamiltonian describing low-energy limit of the honeycomb model [8] which involves only two-spin interactions. The higher-order corrections to the effective Hamiltonian must be regarded as a perturbation. Thus in order to show that the honeycomb model is topologically equivalent to the toric code (in the  $J_z \gg J_x$ ,  $J_y$  phase) one has to prove that the spectral gap in the toric code Hamiltonian does not close in a presence of weak perturbations V that can be represented as a sum of bounded-norm short-range (exponentially decaying) interactions.

Even if one leaves aside the question of how multi-spin interactions can be implemented in a lab, one has to worry about precision up to which an ideal model Hamiltonian can be approximated in a real life. If the presence or absence of the gap depends on tiny variations of the Hamitonian parameters that are beyond experimentalist's control, the distinction between gapped and gapless Hamiltonians is meaningless. The best we can hope for is to approximate individual interactions of the ideal model with some constant precision  $\epsilon$  independent of the system size N. Accordingly, the ideal Hamiltonian can be approximated only up to an extensive error  $O(\epsilon N)$ . Proving stability of topological phases thus reduces to proving that the spectral gap of the ideal TQO models does not close in the presence of such extensive perturbations.

Currently, the tools for proving lower bounds on the spectral gap are fairly limited. For example, one of the outstanding problems in mathematical physics is to prove the existence of a spectral gap for the spin-1 Heisenberg chain, making rigorous the arguments of Haldane [9]. Some progress toward this was obtained by Yarotsky [10], who showed the stability of the gap near the AKLT point [11]. Yarotsky's tools however are limited to perturbations of Hamiltonians which are topologically trivial. Thus, new methods are needed to analyze topologically ordered phases. Some partial results were recently obtained by Trebst et al [12] and Klich [13] who proved gap stability for the toric code under a special type of perturbations diagonal in the

z-basis as well as for anyon lattices on a sphere.

In the present paper we succeed in proving gap stability under generic local perturbations. Our results are valid not just for the toric code, but more generally for any Hamiltonian which can be written as a sum of geometrically local commuting projectors on a *D*-dimensional lattice with certain topological order conditions that we define later. This includes models such as Kitaev's quantum double model [1] and the Levin-Wen string-net model [4]. Furthermore, we prove stability of the spectral gaps separating sufficiently low-lying eigenvalues of the unperturbed Hamiltonian. In the case of 2D models with anyonic excitations it allows us to define string-like operators that create particle excitations for the perturbed Hamiltonian and prove stability of invariants describing the braiding statistics of excitations. We explain how this may be used to adiabatically control a perturbed topological model to perform braiding operations to manipulate topologically protected quantum information.

## 1.1 Summary of results

Consider a system composed of finite-dimensional quantum particles (qudits) occupying sites of a D-dimensional lattice  $\Lambda$  of linear size L. The corresponding Hilbert space is a tensor product of the local Hilbert spaces,  $\mathcal{H} = \bigotimes_{u \in \Lambda} \mathcal{H}_u$ , dim  $(\mathcal{H}_u) = O(1)$ . Suppose the unperturbed Hamiltonian  $H_0$  can be written as a sum of geometrically local pairwise commuting projectors,

$$H_0 = \sum_{A \subseteq \Lambda} Q_A,$$

where the sum runs over all subsets of the lattice of diameter O(1) and  $Q_A$  is a projector acting non-trivially only on sites of A (one may have  $Q_A = 0$  for some subsets A). The commutativity assumption implies that all projectors  $Q_A$  can be diagonalized in the same basis. Accordingly, all eigenvalues of  $H_0$  are non-negative integers. We assume that the smallest eigenvalue of  $H_0$  is zero, that is, ground states of  $H_0$  are annihilated by every projector  $Q_A$ . Such states span the ground subspace P,

$$P = \{ |\psi\rangle \in \mathcal{H} : Q_A |\psi\rangle = 0 \text{ for all } A \subseteq \Lambda \}.$$

For any subset  $B \subseteq \Lambda$  we shall also define a local ground subspace as

$$P_B = \{ |\psi\rangle \in \mathcal{H} : Q_A |\psi\rangle = 0 \text{ for all } A \subseteq B \}.$$

We shall use the notations P and  $P_B$  both for linear subspaces and for the corresponding projectors. Note that the projector  $P_B$  acts non-trivially only on the subset B.

We shall impose two extra conditions on  $H_0$  and the ground subspace P that guarantee the gap stability. Let us first state these conditions informally (see Section 2.2 for formal definitions):

**TQO-1:** The ground subspace P is a quantum code with a macroscopic distance<sup>1</sup>,

TQO-2: Local ground subspaces are consistent with the global one

<sup>&</sup>lt;sup>1</sup>For our purposes it suffices that the distance grows as a positive power of the lattice size L

Condition TQO-1 is the traditional definition of TQO. It guarantees that a local operator cannot induce transitions between orthogonal ground states or distinguish a pair of orthogonal ground states from each other. Thus a local perturbation can lift the ground state degeneracy only in the n-th order of perturbation theory, where n can be made arbitrarily large by increasing the lattice size, see [1]. Surprising, condition TQO-1 by itself is not sufficient for stability, see a simple counter-example in Section 2.4.

Condition TQO-2 demands that a local ground subspace  $P_B$  and the global ground subspace P must be consistent, namely, the projectors  $P_B$  and P must have the same reductions on any subset  $A \subset B$  which is "sufficiently far" from the boundary of B. We need to impose TQO-2 only for regions with trivial topology such a cube or a ball. The consistency between the global and the local ground subspaces may be violated for regions with non-trivial topology. For example, if B has a hole, the local ground subspace  $P_B$  may include sectors with a non-trivial topological charge inside the hole as opposed to the global ground subspace. Condition TQO-2 by itself is also not sufficient for stability, see a counter-example in Section 2.4.

Let us emphasize that all our results apply also to the special case when  $H_0$  has non-degenerate ground state. In this case TQO-1 is automatically satisfied since P is a one-dimensional subspace and thus condition TQO-2 alone guarantees the gap stability.

We consider a perturbation V that can be written as a sum of bounded-norm interactions

$$V = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} V_{r,A},$$

where S(r) is a set of cubes of linear size r and  $V_{r,A}$  is an operator acting on sites of A. We assume that the magnitude of interactions decays exponentially for large r,

$$\max_{A \in \mathcal{S}(r)} \|V_{r,A}\| \le J e^{-\mu r},$$

where  $J, \mu > 0$  are some constants independent of L. Our main result is the following.

**Theorem 1.** Suppose  $H_0$  obeys TQO-1,2. Then there exist constants  $J_0, c_1, c_2 > 0$  depending only on  $\mu$  and the spatial dimension D such that for all  $J \leq J_0$  the spectrum of  $H_0 + V$  is contained (up to an overall energy shift) in the union of intervals  $\bigcup_{k\geq 0} I_k$ , where k runs over the spectrum of  $H_0$  and

$$I_k = \{ \lambda \in \mathbb{R} : k(1 - c_1 J) - \delta \le \lambda \le k(1 + c_1 J) + \delta \},$$

and

$$\delta = poly(L) \exp\left(-c_2 L^{3/8}\right).$$

In other words, the perturbation V changes positive eigenvalues of  $H_0$  at most by a constant factor  $1 \pm c_1 J$  (neglecting the exponentially small correction  $\delta$ ) while the smallest eigenvalue k = 0 is transformed into a band  $I_0$  of exponentially small width  $2\delta$ , see Fig. 1. One can easily

check that for any fixed k the band  $I_k$  is separated from all other bands  $I_m$ ,  $m \neq k$ , by a gap at least 1/2 provided that  $J < J_k$ , where

$$J_k = \frac{1}{c_1(4k+2)}.$$

Thus the bands originating from eigenvalues  $0, 1, \ldots, k$  of  $H_0$  are separated from each other and from the rest of the spectrum by a gap at least 1/2 provided that  $J < J_k$ .

In the case when excitations of  $H_0$  are anyons, one can infer all topological invariants such as S, R, and F-matrices by evaluating fusion and braiding diagrams with only a few particles (for example 4 particles suffice to compute all F matrices). Accordingly, any matrix element of, say, F-matrix, can be represented as an expectation value  $\langle \psi_0 | O_m \dots O_2 O_1 | \psi_0 \rangle$  where  $|\psi_0\rangle$  is the ground state and  $O_i$  are operators creating pairs of excitations from the ground state, moving, fusing, and annihilating them. Our stability result for excited states with O(1) excitations allows us to construct quasi-adiabatic continuation of operators  $O_i$  thus explicitly demonstrating that the perturbed Hamiltonian has the same S, R, and F matrices as the ideal one, see Section 7.

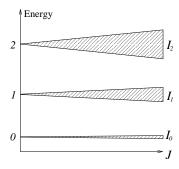


Figure 1: Energy bands  $I_k$  describing the spectrum of a perturbed Hamiltonian  $H_0 + V$ .

## 1.2 Sketch of the stability proof

Let us sketch the main steps of the proof of Theorem 1. We start from proving the theorem for a special class of perturbation V such that all individual interactions  $V_{r,A}$  preserve the ground subspace P, that is,  $[P, V_{r,A}] = 0$ . We call such perturbations block-diagonal. In Section 3 we prove that block-diagonal perturbations are relatively bounded by  $H_0$ , that is,  $||V\psi|| \le b||H_0\psi||$  for any state  $|\psi\rangle \in \mathcal{H}$  and for some coefficient b = O(J). Here for simplicity we ignore some exponentially small corrections. A nice feature of relatively bounded perturbations is that the spectrum of a perturbed Hamiltonian  $H_0 + V$  is contained in the union of intervals  $I_k$  where k runs over the spectrum of  $H_0$  and  $I_k = (k(1-b), k(1+b))$ , see Section 3.1. The proof of the relative boundness is rather elementary and uses certain decomposition of the Hilbert space in terms of syndrome subspaces which is a standard tool in the theory of quantum error correcting codes. In order to get a strong enough bound on the coefficient b we use a novel technique of "coarse-graining" the syndrome subspaces, see Section 3.2 for details.

In the second part of the proof we reduce generic perturbations V to block-diagonal perturbations. Specifically, we construct a unitary operator U such that  $U(H_0 + V)U^{\dagger} \approx H_0 + W$ , where W is a block-diagonal perturbation. Since U does not change eigenvalues, we can use the techniques described above to analyze the spectrum of  $H_0 + V$ . The operator U is constructed using a discrete version of Hamiltonian flow equations developed by Glazek, Wilson, and Wegner [14]. Specifically, we define a hierarchy of Hamiltonians  $H(n) = H_0 + V(n) + W(n)$  labeled by an integer level  $n \geq 0$ , such that W(n) is a block-diagonal perturbation while V(n) is a generic perturbation. We start at the level n = 0 with the perturbed  $H_0 + V$ , that is, V(0) = V and W(0) = 0. As we go to higher levels, the Hamiltonian H(n) becomes more close to a block-diagonal form. The transformation from H(n) to H(n+1) is described by a unitary operator U(n) that block-diagonalizes H(n) up to errors of order  $V(n)^2$ . These errors are dealt with at the next level of the hierarchy, see Section 4 for details. We construct U(n) by solving a linearized block-diagonalization problem, see Section 5. The solution can be easily constructed in terms of the series while convergence of the series follows from the fact that W(n) is relatively bounded by  $H_0$ .

We prove that the strength of V(n) decays doubly-exponentially as a function of n, while W(n) does not change essentially after the first few levels. We then choose the desired unitary operator U as  $U = U(n_f) \cdots U(1)U(0)$  where the highest level  $n_f \sim \log(L)$  is chosen to make the norm of  $V(n_f)$  exponentially small (as a function of L). The most technical part of the proof is to show that the unitary operators U(n) are locality preserving such that all Hamiltonians V(n) and W(n) remain sufficiently local. To this end we first prove that U(n) can be generated by a quasi-local Hamiltonian, see Section 5, and then employ the Lieb-Robinson bound, see Section 6.

# 2 Hamiltonians describing TQO

## 2.1 Frustration-free commuting Hamiltonians

To simplify notations we shall restrict ourselves to the spatial dimension D=2. A generalization to an arbitrary D is straightforward. Let  $\Lambda=\mathbb{Z}_L\times\mathbb{Z}_L$  be a two-dimensional square lattice of linear size L with periodic boundary conditions. We assume that every site  $u\in\Lambda$  is occupied by a finite-dimensional quantum particle (qudit) such that the Hilbert space describing  $\Lambda$  is a tensor product

$$\mathcal{H} = \bigotimes_{u \in \Lambda} \mathcal{H}_u, \quad \dim \mathcal{H}_u = O(1).$$

Let S(r) be a set of all square blocks  $A \subseteq \Lambda$  of size  $r \times r$ , where r is a positive integer. Note that S(r) contains  $L^2$  translations of some elementary square of size  $r \times r$  for all r < L,  $S(L) = \Lambda$ , and  $S(r) = \emptyset$  for r > L. We can always assume that the unperturbed Hamiltonian  $H_0$  involves only  $2 \times 2$  interactions (otherwise consider a coarse-grained lattice):

$$H_0 = \sum_{A \in \mathcal{S}(2)} G_A. \tag{2.1}$$

There will be three essential restrictions on the form of interactions  $G_A$ . Firstly, we require that  $G_A$  are pairwise commuting operators, that is,

$$G_A G_B = G_B G_A$$
 for all  $A, B \in \mathcal{S}(2)$ .

Thus all interactions  $G_A$  can be diagonalized in the same basis. Secondly, we require that  $H_0$  is a frustration free Hamiltonian, that is, the ground state of  $H_0$  minimizes energy of every individual term  $G_A$ . Performing an overall energy shift we can always assume that all  $G_A$  are positive-semidefinite operators,

$$G_A > 0$$
.

Then the condition of being frustration-free demands that ground states of  $H_0$  are common zero eigenvectors of every term  $G_A$ . Thus the ground subspace of  $H_0$  is

$$P = \{ |\psi\rangle \in \mathcal{H} : G_A |\psi\rangle = 0 \quad \text{for all } A \in \mathcal{S}(2) \}. \tag{2.2}$$

Thirdly, we shall assume that every operator  $G_A$  has a constant spectral gap, that is, the smallest positive eigenvalue of  $G_A$  is bounded from below by a constant independent of the lattice size L. We can always normalize the Hamiltonian  $H_0$  such that the spectral gap of any  $G_A$  is at least 1. This is equivalent to a condition

$$G_A^2 \ge G_A$$
.

Let  $P_A$  be the projector onto the zero subspace of  $G_A$  and  $Q_A = I - P_A$ . Note that all the projectors  $P_A$ ,  $Q_A$  are pairwise commuting. For any square  $B \in \mathcal{S}(r)$ ,  $r \geq 2$  define a projector onto the local ground subspace

$$P_B = \prod_{\substack{A \in \mathcal{S}(2) \\ A \subseteq B}} P_A \tag{2.3}$$

and  $Q_B = I - P_B$ . Note that  $P_B$  and  $Q_B$  have support on B. We shall often use the same notation for a subspace and for the corresponding projector.

## 2.2 Formal definition of TQO

We shall need two extra property of  $H_0$  and the ground subspace P that guarantee the gap stability and robustness of the ground state degeneracy. We shall assume that there exists a constant c > 0 such that the following conditions hold for some integer  $L^* \geq cL$  for all sufficiently large L:

**TQO-1:** Let  $A \in \mathcal{S}(r)$  be any square of size  $r \leq L^*$ . Let  $O_A$  be any operator acting on A. Then

$$PO_AP = cP$$

for some complex number c.

**TQO-2:** Let  $A \in \mathcal{S}(r)$  be any square of size  $r \leq L^*$  and let  $B \in \mathcal{S}(r+2)$  be the square that contains A and all nearest neighbors of A. Define reduced density matrices  $\rho_A = \operatorname{Tr}_{A^c}(P)$  and  $\rho_A^{(B)} = \operatorname{Tr}_{A^c}(P_B)$ . Then the kernel of  $\rho_A$  coincides with the kernel of  $\rho_A^{(B)}$ .

Remark 1. Using the language of quantum error correcting codes one can define the minimum distance of P as the smallest integer d such that erasure of any subset of d particles can be corrected for any encoded state  $|\psi\rangle \in P$ , see [15] for details. Note that TQO-1 holds for  $L^* = \lfloor \sqrt{d} \rfloor$  since the reduced state of any square  $A \in \mathcal{S}(L^*)$  does not depend on the encoded state. What is less trivial, TQO-1 holds also for  $L^* = \Omega(d)$ , see [15]. Thus  $L^*$  coincides with the distance of the code P up to a constant coefficient (as far as condition TQO-1 is concerned).

Remark 2. Condition TQO-2 can be easily 'proved' if the excitations of  $H_0$  are anyons (since the latter assumption lacks a rigorous formulation, the argument given below is not completely rigorous either). Indeed, in this case we can choose a complete basis of the excited subspace Q such that the basis vectors correspond to various configurations of anyons. For non-abelian theories one may have several basis vectors for a fixed configuration of anyons that describe different fusion channels, see [8]. Note that any state  $|\psi\rangle \in P_B$  is a superposition of configurations with no anyons inside B. Since A is a topological trivial region, any such configuration can be prepared from the vacuum P by some unitary operator  $U_{A^c}$  acting on complementary region  $A^c = \Lambda \backslash A$ . Thus  $|\psi\rangle = U_{A^c}|\psi_0\rangle$  for some ground state  $|\psi_0\rangle \in P$ . Since all ground states  $|\psi_0\rangle$  have the same reduced matrix on A, it means that  $|\psi\rangle$  and P have the same reduced matrix on A. This implies TQO-2. The above arguments suggest that TQO-2 holds for all 2D models of TQO that can be described by commuting frustration-free such as quantum double models [1] and Levin-Wen string-net models [4].

Remark 3. As was already mentioned, the consistency between the global and the local ground subspaces may be violated for regions with non-trivial topology. For example, if A has a hole, the local ground subspace  $P_A$  may include sectors with a non-trivial topological charge inside the hole as opposed to the global ground subspace.

We shall need the following corollary of TQO-2.

**Corollary 2.1.** Let  $A \in \mathcal{S}(r)$  be any square of size  $r \leq L^*$  and  $O_A$  be any operator acting on A such that  $O_AP = 0$ . Let  $B \in \mathcal{S}(r+2)$  be the square that contains A and all nearest neighbors of A. Then  $O_AP_B = 0$ .

*Proof.* Let  $\rho_A = \operatorname{Tr}_{A^c}(P)$ . The assumption  $O_A P O_A^{\dagger} = 0$  implies that  $O_A \rho_A O_A^{\dagger} = 0$ , that is,  $O_A$  annihilates any state in the range of  $\rho_A$ . From TQO-2, the range of  $\operatorname{Tr}_{A^c}(P_B)$  coincides with the range of  $\rho_A$ , and thus  $\operatorname{Tr}(O_A P_B O_A^{\dagger}) = 0$ . It implies  $O_A P_B = 0$ .

## 2.3 Verification of TQO conditions for stabilizer Hamiltonians

Conditions TQO-1,2 can be easily checked for those models of TQO that can be described using the stabilizer formalism such as the toric code model [1] or topological color codes [16]. For such models each site of the lattice  $\Lambda$  represents one or several qubits, while the ground state subspace P is a stabilizer code, i.e., the invariant subspace of some abelian stabilizer group  $\mathcal{G} \subseteq \operatorname{Pauli}(\Lambda)$ . Here  $\operatorname{Pauli}(\Lambda)$  is a group generated by single-qubit Pauli operators  $\sigma_i^x, \sigma_i^y, \sigma_i^z$ . The stabilizer group must have a set of geometrically local generators, that is,  $\mathcal{G} = \langle S_1, \ldots, S_M \rangle$  where any generator  $S_a \in \operatorname{Pauli}(\Lambda)$  acts non-trivially only on O(1) qubits located within distance O(1) from each other. Note that the generators need not to be independent. We choose the corresponding stabilizer Hamiltonian  $H_0$  as

$$H_0 = \sum_a (I - S_a)/2$$

such that states invariant under action of stabilizers have zero energy. The minimal distance of the code is the smallest integer d such that there exists a Pauli operator O that commutes with all elements of  $\mathcal{G}$  but does not belong to  $\mathcal{G}$ . Such an operator O can be regarded as a logical Pauli operator acting on encoded states. It follows from results of [15] that condition TQO-1 holds if we choose  $L^* = \Omega(d)$ .

Assume that the set of qubits is coarse-grained into sites of the lattice  $\Lambda$  such that the support of any generator  $S_a$  is contained in at least one  $2 \times 2$  square. One can bring this Hamiltonian into the form Eq. (2.1) by distributing the generators over  $2 \times 2$  squares in an arbitrary way. For any square  $B \in \mathcal{S}(r)$  one can define two subgroups of  $\mathcal{G}$ : (i) a subgroup  $\mathcal{G}_B$  generated by generators  $S_a$  whose support is contained in B, and (ii) a subgroup  $\mathcal{G}(B)$  that includes all stabilizers  $S \in \mathcal{G}$  whose support is contained in B. By definition,  $\mathcal{G}_B \subseteq \mathcal{G}(B)$ , but in general  $\mathcal{G}_B \neq \mathcal{G}(B)$ .

**Lemma 2.1.** The stabilizer Hamiltonian  $H_0$  satisfies condition TQO-2 iff for any square  $A \in \mathcal{S}(r)$ ,  $r \leq L^*$ , one has  $\mathcal{G}(A) \subseteq \mathcal{G}_B$ , where  $B = b_1(A)$ .

Thus TQO-2 demands that any element of the stabilizer group whose support is contained in a square A can be written as a product of generators whose support is contained in A and a small neighborhood of A. We leave verification of this condition for the toric code model as an exercise for the reader.

*Proof.* Indeed, the reduced density matrix  $\rho_A$  computed using the global ground subspace P is proportional to the projector onto the codespace of the stabilizer code  $\mathcal{G}(A)$ . The reduced density matrix  $\rho_A$  computed using the local ground subspace  $P_B$  is proportional onto the codespace of the stabilizer code  $\mathcal{G}_B(A)$ , where  $\mathcal{G}_B(A)$  includes all elements of  $\mathcal{G}_B$  whose support is contained in A. Thus TQO-2 holds iff  $\mathcal{G}(A) = \mathcal{G}_B(A)$ . This is equivalent to the condition of the lemma.  $\square$ 

## 2.4 Unstable version of the toric code model

In this section we demonstrate that condition TQO-1 alone is not sufficient for stability. Let us start from the standard toric code model [1],

$$H_{tc} = -\sum_{p} B_{p} - \sum_{s} A_{s},$$

where qubits live on edges of a square 2D lattice, p and s labels plaquettes and sites of the lattice,  $B_p$  is a product of  $\sigma^z$  over the four boundary edges of p, and  $A_s$  is a product of four  $\sigma^x$  over the four edges incident to s. We shall refer to  $B_p$  and  $A_s$  as plaquette and star operators. The ground subspace P is defined by eigenvalue equations  $B_p = 1$  for all p and  $A_s = 1$  for all s. It is well known that P is a quantum code with the minimal distance d = L - 1. Hence P obeys TQO-1,2 with  $L^* = L - 1$ .

Consider now a modified toric code model

$$H'_{tc} = -\sum_{(p,q)} B_p B_q - \sum_s A_s - B_{p^*}$$

where (p,q) labels pairs of adjacent plaquettes and  $p^*$  is some selected plaquette. We assume that the total number of plaquettes  $N_p$  is even. Note that  $H'_{tc}$  is a frustration-free commuting Hamiltonian. In addition,  $H'_{tc}$  and  $H_{tc}$  have the same ground subspace P corresponding to  $B_p = 1$ ,  $A_s = 1$  for all s and p. Hence  $H'_{tc}$  obeys TQO-1. We claim that  $H'_{tc}$  violates TQO-2. Indeed, choose any square A located sufficiently far from the selected plaquette  $p^*$ . Then the local ground subspace  $P_A$  has equal contributions from sectors  $B_p = 1$ ,  $A_s = 1$  and  $B_p = -1$ ,  $A_s = 1$  thus being inconsistent with the global ground state.

Let us now argue that the spectral gap of  $H'_{tc}$  closes in a presence of local perturbations with a strength of order  $1/N_p$ . This instability has the same origin as the instability of the classical 2D Ising model under external magnetic field. Note that  $H'_{tc}$  has spectral gap  $\Delta = 2$  and the second smallest eigenvalue belongs to the sector  $A_s = 1$ ,  $B_p = -1$  for all s and p. Consider a perturbation describing an "external magnetic field",

$$V = h \sum_{p} B_{p}, \quad h > 0.$$

For sufficiently large h, say,  $h = 4/N_p$ , the ground state of  $H'_{tc} + V$  moves from the sector  $A_s = 1$ ,  $B_p = 1$  to the sector  $A_s = 1$ ,  $B_p = -1$ . Hence the gap above the ground state closes for some intermediate value of h.

Needless to say, condition TQO-2 alone is also not sufficient for stability. The simplest counter-example is the 2D classical Ising model in which the gap is unstable under external magnetic field.

# 3 Relatively bounded perturbations

## 3.1 Definition and basic properties

In this section we introduce necessary facts from the theory of relatively bounded perturbations. It mostly follows Chapter IV of [17] although our definitions and proofs are much simpler since we are interested only in finite-dimensional Hilbert spaces.

Let  $H_0$  and W be any Hamiltonians acting on some Hilbert space  $\mathcal{H}$ . We shall say that W is relatively bounded by  $H_0$  iff there exist  $0 \le b < 1$  such that

$$||W\psi|| \le b \, ||H_0\psi|| \quad \text{for all } |\psi\rangle \in \mathcal{H}.$$
 (3.1)

The notion of a relatively bounded perturbation allows one to define a "weak perturbation" and rigorously justify application of perturbative expansions even when the norm of W is much larger than the spectral gap of  $H_0$ . We shall be mostly interested in the case when b is a constant independent of the lattice size L. Note that the condition Eq. (3.1) is equivalent to an operator inequality  $W^2 \leq b^2 H_0^2$ .

The following lemma asserts that a relatively bounded perturbation can change eigenvalues of  $H_0$  at most by a factor  $1 \pm b$ .

**Lemma 3.1.** Suppose W is relatively bounded by  $H_0$ . Then the spectrum of  $H_0+W$  is contained in the union of intervals  $[\lambda_0(1-b), \lambda_0(1+b)]$  where  $\lambda_0$  runs over the spectrum of  $H_0$ .

*Proof.* Indeed, suppose  $(H_0 + W) |\psi\rangle = \lambda |\psi\rangle$ , that is,

$$(H_0 - \lambda I) |\psi\rangle = -W |\psi\rangle. \tag{3.2}$$

The relative boundness then implies  $\|(H_0 - \lambda I)\psi\| \le b\|H_0\psi\|$ , that is,

$$\langle \psi | (H_0 - \lambda I)^2 | \psi \rangle \le b^2 \langle \psi | H_0^2 | \psi \rangle.$$
 (3.3)

Let  $H_0 = \sum_{\lambda_0} \lambda_0 P_{\lambda_0}$  be the spectral decomposition of  $H_0$ . Here the sum runs over the spectrum of  $H_0$  and  $P_{\lambda_0}$  is a projector onto the eigenspace with an eigenvalue  $\lambda_0$ . Define a probability distribution  $p(\lambda_0) = \langle \psi | P_{\lambda_0} | \psi \rangle$ . Substituting it into Eq. (3.3) one gets

$$\sum_{\lambda_0} (\lambda_0 - \lambda)^2 p(\lambda_0) \le \sum_{\lambda_0} b^2 \lambda_0^2 p(\lambda_0). \tag{3.4}$$

Therefore there exists at least one eigenvalue  $\lambda_0$  such that

$$(\lambda_0 - \lambda)^2 \le b^2 \lambda_0^2. \tag{3.5}$$

This is equivalent to  $\lambda_0(1-b) \leq \lambda \leq \lambda_0(1+b)$ .

## 3.2 Stability of TQO under block-diagonal perturbations

In this section we shall consider perturbations

$$W = \sum_{A \in \mathcal{S}(q)} W_A$$

such that all local terms  $W_A$  are block-diagonal,

$$[W_A, P] = 0$$
 for all  $A \in \mathcal{S}(q)$ .

We shall assume that  $q \leq L^*$ , so condition TQO-1 implies that the restriction of  $W_A$  onto the P-subspace is a multiple of the identity. Since we are not interested in the overall shift in energy, we can assume that

$$W_A P = 0$$
 for all  $A \in \mathcal{S}(q)$ . (3.6)

The interaction strength of W will be measured by a parameter

$$w = \max_{A \in \mathcal{S}(q)} \|W_A\|. \tag{3.7}$$

**Lemma 3.2.** Let W be a perturbation satisfying Eqs. (3.6,3.7). Then W is relatively bounded by  $H_0$  with a constant

$$b = O(wq^2).$$

*Proof.* Let us start from introducing some notations. A syndrome  $s: \mathcal{S}(2) \to \{0,1\}$  is a function that assigns an eigenvalue  $s_A \in \{0,1\}$  to every projector  $Q_A$ ,  $A \in \mathcal{S}(2)$ . Given a syndrome s and a square  $A \in \mathcal{S}(2)$  we shall say that A is a defect iff  $s_A = 1$ . Thus one can consider s as a configuration of defects. For any syndrome s define a projector

$$R_s = \prod_{A \in S(2)} [s_A Q_A + (1 - s_A)(I - Q_A)]$$

projecting onto a subspace spanned by states with a syndrome s. Clearly the family of projectors  $R_s$  defines an orthogonal decomposition of the Hilbert space, that is,  $\sum_s R_s = I$ .

Let us fix some partition of the lattice into contiguous  $q \times q$  squares  $B_1, \ldots, B_M \in \mathcal{S}(q)$  (if L is not a multiple of q, the squares  $B_i$  may have size  $q \pm O(1)$ ). We shall refer to a set of  $2 \times 2$  squares contained in a particular square  $B_i$  as a box. We shall need the following properties:

- 1. Every square  $A \in \mathcal{S}(2)$  is contained in exactly one box  $B_i$ ,
- 2. Each box  $B_i$  overlaps with  $O(q^2)$  squares  $C \in \mathcal{S}(q)$ ,
- 3. Each square  $C \in \mathcal{S}(q)$  overlaps with O(1) boxes  $B_i$ .

Given a syndrome s and a box  $B_i$  we shall say that  $B_i$  is occupied if  $B_i$  contains at least one defect, that is, there is a  $2 \times 2$  square  $A \subset B$  such that  $s_A = 1$ . Otherwise we shall say that the box  $B_i$  is empty.

Given a syndrome s let  $b(s) \subseteq [M]$  be the subset of occupied boxes. For any subset of boxes  $\mathcal{Y} \subseteq [M]$  define a projector

$$R_{\mathcal{Y}} = \sum_{s:b(s)=\mathcal{Y}} R_s.$$

It projects onto the subspace in which all boxes in  $\mathcal{Y}$  are occupied and the remaining boxes are empty. Clearly, the family of projectors  $R_{\mathcal{Y}}$  defines an orthogonal decomposition, that is,  $\sum_{\mathcal{Y}\subseteq[M]}R_{\mathcal{Y}}=I$ . We claim that any operator  $W_A$  acting on a square  $A\in\mathcal{S}(q)$  and satisfying Eq. (3.6) has only a few off-diagonal blocks with respect to this decomposition. Specifically, Corollary 2.1 implies that

$$R_{\mathcal{V}}W_{A}R_{\mathcal{Z}} \neq 0 \tag{3.8}$$

only if A has distance O(1) from some occupied box in  $\mathcal{Y}$  and A has distance O(1) from some occupied box in  $\mathcal{Z}$ , and the configurations  $\mathcal{Y}, \mathcal{Z}$  differ only at those boxes that overlap with A. Clearly, for any fixed  $\mathcal{Y} \subseteq [M]$  such that  $\mathcal{Y}$  has k occupied boxes the number of pairs  $(A \in \mathcal{S}(q), \mathcal{Z} \subseteq [M])$  that could satisfy Eq. (3.8) is at most  $O(kq^2)$ . Thus for any state  $|\psi\rangle$  we get

$$\langle \psi | W^{2} | \psi \rangle = \sum_{\mathcal{Y}, \mathcal{Z}, \mathcal{V} \subseteq [M]} \langle \psi | R_{\mathcal{Y}} W R_{\mathcal{Z}} W R_{\mathcal{V}} | \psi \rangle$$

$$\leq \sum_{\mathcal{Y}, \mathcal{Z}, \mathcal{V} \subseteq [M]} \| R_{\mathcal{Y}} W R_{\mathcal{Z}} \| \cdot \| R_{\mathcal{Z}} W R_{\mathcal{V}} \| \cdot \| R_{\mathcal{Y}} \psi \| \cdot \| R_{\mathcal{V}} \psi \|$$

$$\leq \sum_{\mathcal{Y}, \mathcal{Z}, \mathcal{V} \subseteq [M]} \| R_{\mathcal{Y}} W R_{\mathcal{Z}} \| \cdot \| R_{\mathcal{Z}} W R_{\mathcal{V}} \| \cdot \frac{1}{2} (\langle \psi | R_{\mathcal{Y}} | \psi \rangle + \langle \psi | R_{\mathcal{V}} | \psi \rangle)$$

$$= \sum_{\mathcal{Y}, \mathcal{Z}, \mathcal{V} \subseteq [M]} \| R_{\mathcal{Y}} W R_{\mathcal{Z}} \| \cdot \| R_{\mathcal{Z}} W R_{\mathcal{V}} \| \cdot \langle \psi | R_{\mathcal{Y}} | \psi \rangle$$

$$\leq \sum_{k \geq 0} \sum_{\mathcal{Y}: |\mathcal{Y}| = k} O(k^{2} q^{4} w^{2}) \langle \psi | R_{\mathcal{Y}} | \psi \rangle = O(w^{2} q^{4}) \langle \psi | G | \psi \rangle, \tag{3.9}$$

where

$$G = \sum_{k \ge 0} \sum_{\mathcal{Y}: |\mathcal{Y}| = k} k^2 R_{\mathcal{Y}}.$$
 (3.10)

The inequality Eq. (3.9) follows from the fact that  $\mathcal{Y}$  and  $\mathcal{Z}$  differ at at most O(1) boxes and an obvious bound  $k(k+O(1))=O(k^2)$ . Finally, note that  $G \leq H_0^2$  since any configuration of defects with k occupied boxes must have at least k defects and since creating a defect costs at least a unit of energy. We arrive at

$$\langle \psi | W^2 | \psi \rangle \le b^2 \langle \psi | H_0^2 | \psi \rangle, \quad b = O(wq^2).$$
 (3.11)

It completes the proof.

We shall also need a local version of Lemma 3.2. For any region  $C \subseteq \Lambda$  define a local version of the Hamiltonian  $H_0$ ,

$$H_0(C) = \sum_{\substack{A \in \mathcal{S}(2) \\ A \subseteq C}} G_A. \tag{3.12}$$

Let  $P_C$  and  $Q_C$  be the local versions of the projectors P and Q defined in Eq. (2.3). The following is a straightforward corollary of Lemma 3.2.

**Corollary 3.1.** Let  $C \subseteq \Lambda$  be any square of size smaller than  $L^*$ . Let  $W = \sum_{A \in \mathcal{S}(q)} W_A$  be a perturbation satisfying Eqs. (3.6,3.7). Suppose also that  $W_A = 0$  unless C contains both A and the nearest neighbors of A. Then W is relatively bounded by  $H_0(C)$  with a constant  $b = O(wq^2)$ .

*Proof.* Combining Eq. (3.6) and TQO-2 we conclude that  $W_A P_C = 0$  for all A. Thus we can apply all steps in the proof of Lemma 3.2 to the square C considered as the entire lattice  $\Lambda$ .  $\square$ 

We shall need another technical lemma that provides a bound on the norm of a commutator involving a block-diagonal Hamiltonian. Let us start from the simplest scenario. Let W be any operator such that W is relatively bounded by  $H_0$  with a constant  $0 \le b < 1$ . Then for any operator S we have

$$\begin{split} \|Q[S,W]P\| &= \|QWSP\| = \|QWH_0^{-1}QH_0SP\| \\ &= \|QWH_0^{-1}Q[H_0,S]P\| \leq \|WH_0^{-1}Q\| \cdot \|Q[H_0,S]P\|. \end{split}$$

Here the first equality follows from WP=0 and the third equality uses  $H_0P=0$ . Let  $|\psi\rangle\in Q$  be a normalized state such that  $||WH_0^{-1}Q||=||WH_0^{-1}\psi||$ . Using the relative boundness assumption we get

$$||WH_0^{-1}Q|| = ||WH_0^{-1}\psi|| \le b||H_0H_0^{-1}\psi|| = b.$$

To conclude, we have proved that

$$||Q[S, W]P|| \le b ||Q[S, H_0]P||.$$

Applying the same arguments as above with P and Q replaced by their local versions  $P_C$  and  $Q_C$ , see Eq. (2.3), we arrive at the following lemma.

**Lemma 3.3.** Let  $C \subseteq \Lambda$  be any region. Let W be any Hamiltonian such that W is relatively bounded by  $H_0(C)$  with a constant  $0 \le b < 1$ . Then for any operator S one has

$$||Q_C[S, W]P_C|| \le b ||Q_C[S, H_0(C)]P_C||. \tag{3.13}$$

Combining this lemma and Corollary 3.1 we get a simple upper bound on  $||Q_C[S, W]P_C||$ .

Corollary 3.2. Let  $C \subseteq \Lambda$  be any square of size smaller than  $L^*$ . Let  $W = \sum_{A \in \mathcal{S}(q)} W_A$  be a perturbation satisfying Eqs. (3.6,3.7). Suppose also that  $W_A = 0$  unless C contains both A and the nearest neighbors of A. Then for any operator S one has

$$||Q_C[S, W]P_C|| \le O(wq^2) ||Q_C[S, H_0(C)]P_C||.$$
 (3.14)

# 4 Hamiltonian flow equations

#### 4.1 Outline of the method

We are interested in the low-energy spectrum of a perturbed Hamiltonian  $H_0 + V$ , where V is a local perturbation specified by a list of local interactions,

$$V = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} V_{r,A}. \tag{4.1}$$

Here  $V_{r,A}$  is some interaction supported on a square  $A \in \mathcal{S}(r)$ . Our strategy will be to reduce the case of a generic perturbation to the special case of a block-diagonal perturbation which we can analyze using techniques of Section 3. To this end we shall define a hierarchy of Hamiltonians unitarily equivalent to H,

$$H(n) = H_0 + \sum_{k=1}^{n} W(k) + V(n) + E(n) + \lambda(n)I, \quad n = 0, 1, 2, \dots,$$
(4.2)

such that  $H(0) = H_0 + V$  and H(n+1) can be obtained from H(n) by a unitary transformation,

$$H(n+1) = U(n)H(n)U(n)^{\dagger}, \quad U(n)U(n)^{\dagger} = I.$$
 (4.3)

Accordingly, the spectrum of H(n) is the same for all  $n \ge 0$ . The purpose of the transformation U(n) is to make the Hamiltonian more close to the block-diagonal form. We shall refer to H(n) as a level-n Hamiltonian.

Let us describe the purpose of various terms in Eq. (4.2). The Hamiltonian W(k) represents a block-diagonal contribution to the total Hamiltonian H(n) that has been created at the level k. The Hamiltonian V(n) represents the part of the total Hamiltonian H(n) that has to be block-diagonalized at the level n+1. The Hamiltonians V(n) and W(n) will be represented by a sum of local interactions supported in squares of size  $r \leq L^*$ , and such that all interactions involved in W(n) are individually block-diagonal,

$$V(n) = \sum_{1 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} V_{r,A}(n),$$

and

$$W(n) = \sum_{1 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} W_{r,A}(n), \text{ where } QW_{r,A}(n)P = 0.$$

All contributions from squares of size  $r > L^*$  will be collected into the third Hamiltonian E(n) which can be regarded as an error Hamiltonian. The norm of E(n) will be exponentially small in L for all n. The norms of W(n) and V(n) will decay roughly as doubly-exponential functions of n. Thus at level  $n \sim \log L$  the total Hamiltonian H(n) will be block-diagonal up to corrections of order  $\exp(-poly(L))$  resulting from V(n) and E(n). Finally,  $\lambda(n)$  is an overall energy shift that we shall often ignore.

We start at the level n = 0 with initial conditions

$$W_{r,A}(0) = 0$$
,  $V_{r,A}(0) = V_{r,A}$  for  $r \le L^*$ ,  $E(0) = \sum_{r > L^*} \sum_{A \in \mathcal{S}(r)} V_{r,A}$ .

Accordingly,  $H(0) = H_0 + V$  is the Hamiltonian we are interested in. Suppose we have already defined the Hamiltonians  $W(0), \ldots, W(n), V \equiv V(n), E \equiv E(n)$  for some level n. Let

$$W = \sum_{k=1}^{n} W(k)$$

be the overall block-diagonal part of H(n). Let

$$H \equiv H(n) = H_0 + W + V + E.$$

We shall define the operator U(n) in Eq. (4.3) as  $U(n) = \exp(S)$ , where S is the solution of a linearized block-diagonalization problem

$$Q([S, H_0 + W] + V)P = 0, \quad S^{\dagger} = -S.$$
 (4.4)

The meaning of this equation can be easily understood if one treats V as a perturbation and  $H_0 + W$  as an unperturbed Hamiltonian. Expanding the transformed Hamiltonian  $e^S H e^{-S}$  in powers of S we get

$$e^{S}He^{-S} = H_0 + W + ([S, H_0 + W] + V) + O(S^2) + O(SV) + O(E).$$

Thus Eq. (4.4) says that the transformed Hamiltonian must be block-diagonal up to terms  $O(V^2)$  and O(E). The solution of Eq. (4.4) is constructed in Section 5, see Lemma 5.1. We start from defining raw versions of W(n+1) and V(n+1) which we shall denote  $\tilde{W}$  and  $\tilde{V}$  respectively, namely

$$\tilde{W} = [S, H_0 + W] + V,$$
(4.5)

and

$$\tilde{V} = e^{S}(H_0 + W + V)e^{-S} - (H_0 + W + V + [S, H_0 + W]). \tag{4.6}$$

A simple algebra shows that

$$e^{S}He^{-S} = H_0 + W + \tilde{W} + \tilde{V} + e^{S}Ee^{-S}.$$

Note also that  $\tilde{W}$  is block-diagonal due to Eq. (4.4). We shall construct a decomposition of  $\tilde{W}$  into a sum of local interactions that are individually block-diagonal using techniques of Section 5, see Corollary 5.1 of Lemma 5.1. It will yield

$$\tilde{W} = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} \tilde{W}_{r,A}, \text{ where } Q\tilde{W}_{r,A}P = 0.$$

We shall construct a decomposition of  $\tilde{V}$  into a sum of local interactions using Lemma 4.1 and Lemma 6.1 arriving at

$$\tilde{V} = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} \tilde{V}_{r,A}.$$

Next we use  $\tilde{W}$  and  $\tilde{V}$  to define W(n+1) and V(n+1) by taking out all contributions from squares of size  $r > L^*$  and adding these contributions to the error Hamiltonian, that is,

$$W(n+1) = \sum_{1 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} \tilde{W}_{r,A},$$

$$V(n+1) = \sum_{1 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} \tilde{V}_{r,A},$$

and

$$E(n+1) = e^{S} E e^{-S} + \sum_{r>L^*} \sum_{A \in \mathcal{S}(r)} \tilde{W}_{r,A} + \tilde{V}_{r,A}.$$

For the detailed analysis of these flow equations see Section 4.3

## 4.2 Local decompositions of Hamiltonians

In order to analyze convergence of the flow equations we shall need to set up some notations and terminology. Recall that S(r) is a set of all  $r \times r$  squares.

**Definition 4.1.** Let V be any operator acting on  $\mathcal{H}$ . A local decomposition of V is a list of operators  $\{V_{r,A}\}_{r,A}$  where  $r \geq 1$  and  $A \in \mathcal{S}(r)$  such that  $V_{r,A}$  has support on a square A and

$$V = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} V_{r,A} \tag{4.7}$$

Note that a local decomposition of an operator is not unique since the squares involved in the decomposition overlap with each other. Nevertheless, we shall often identify an operator and its local decomposition unless it may lead to confusions.

**Definition 4.2.** A local decomposition of an operator V is  $(J, \mu, \alpha)$ -decaying iff

$$\max_{r \ge 1} \max_{A \in \mathcal{S}(r)} \|V_{r,A}\| \, r^{\alpha} e^{\mu r} \le J. \tag{4.8}$$

Here we mean that  $J, \mu, \alpha$  are some constants independent of L. We shall often use a term  $(J, \mu, \alpha)$ -decaying operator meaning that this operator has a local decomposition which is  $(J, \mu, \alpha)$ -decaying. To simplify notations we shall often use an abbreviation  $(J, \mu)$ -decay for  $(J, \mu, 0)$ -decay.

In the rest of this section we derive several auxiliary technical results that can be skipped at the first reading.

We shall often need to construct a local decomposition for a commutator [S, V] given the local decompositions of S and V.

**Lemma 4.1.** Suppose S is  $(K, \mu, \alpha)$ -decaying and V is  $(J, \mu, \beta)$ -decaying for some  $\alpha, \beta \geq 4$ . Then [S, V] has a local decomposition which is  $(cKJ, \mu)$ -decaying for some constant c.

*Proof.* Consider local decompositions of S and V,

$$S = \sum_{p \ge 1} \sum_{A \in \mathcal{S}(p)} S_{p,A}, \quad ||S_{p,A}|| \le K p^{-\alpha} e^{-\mu p}, \tag{4.9}$$

$$V = \sum_{q \ge 1} \sum_{B \in \mathcal{S}(q)} V_{q,B}, \quad ||V_{q,B}|| \le J q^{-\beta} e^{-\mu q}. \tag{4.10}$$

If  $A \in \mathcal{S}(p)$  and  $B \in \mathcal{S}(q)$  is a pair of non-overlapping squares,  $A \cap B \neq \emptyset$ , then clearly  $A \cup B$  can be covered by some square  $C \in \mathcal{S}(p+q)$ . Thus we can choose a local decomposition of [S, V] as

$$[S, V] = \sum_{r>2} \sum_{C \in \mathcal{S}(r)} D_{r,C},$$
 (4.11)

where

$$D_{r,C} = \sum_{p+q=r} \sum_{\substack{A \in \mathcal{S}(p) \\ A \subseteq C}} \sum_{\substack{B \in \mathcal{S}(q) \\ B \subseteq C}} \left[ S_{p,A}, V_{q,B} \right] \chi(A, B, C) \tag{4.12}$$

where  $\chi(A, B, C) = 0, 1$  is some function that 'distributes' the commutators  $[S_{p,A}, V_{q,B}]$  over different terms  $D_{r,C}$ . A specific form of this function is not important for us. For fixed p, q such that p + q = r and a fixed  $C \in \mathcal{S}(r)$  we can bound the number of squares  $A, B \subseteq C$  as

$$\#\{A \in \mathcal{S}(p) : A \subseteq C\} = (r-p)^2 = q^2$$

and

$$\#\{B \in \mathcal{S}(q) : B \subseteq C\} = (r-q)^2 = p^2.$$

It yields

$$||D_{r,C}|| \le 2KJe^{-\mu r} \sum_{p+q=r} p^{2-\alpha} q^{2-\beta} \le 2KJe^{-\mu r} \sum_{p,q\ge 1} p^{-2} q^{-2} = cKJe^{-\mu r}$$
(4.13)

for some constant c provided that  $\alpha, \beta \geq 4$ .

We shall also need the following simple lemma that will allow us to amplify the degree  $\alpha$  by any constant by "borrowing" some decay from the exponential function. It makes the decay rate  $\mu$  a bit smaller and the amplitude J a bit larger.

**Lemma 4.2** (Degree Reset). Suppose V is  $(J, \mu, \beta)$ -decaying. Let  $0 < \epsilon < 1$  and  $\alpha > 0$  be any constants. Then V is also  $(J', \mu', \alpha + \beta)$ -decaying where

$$J' = cJ^{1-\epsilon} \quad and \quad \mu' = \mu - J^{\frac{\epsilon}{\alpha}}. \tag{4.14}$$

Here c is a constant depending on  $\epsilon$  and  $\alpha$ .

*Proof.* Indeed, let  $\mu' = \mu - \delta$  where  $\delta$  will be chosen later. Then

$$||V_{r,A}|| \le Jr^{-\beta}e^{-\mu r} \le Jr^{-\alpha-\beta}e^{-\mu' r} \max_{q \ge 0} q^{\alpha}e^{-\delta q} \le c \,\delta^{-\alpha}Jr^{-\alpha-\beta}e^{-\mu' r},\tag{4.15}$$

where  $c = \max_{x \geq 0} x^{\alpha} e^{-x} = O(1)$  is a constant. Choosing  $\delta = J^{\epsilon/\alpha}$  we achieve the desired scaling.

Finally, we shall use the following trivial observation.

**Lemma 4.3.** Suppose V is  $(J, \mu, \alpha)$ -decaying. Then V also has a local decomposition which is  $(cJe^{2\mu}, \mu, \alpha)$ -decaying with an extra property that  $V_{r,A}$  acts trivially on all sites that are adjacent to the boundary of A. Here c = O(1) is some constant.

*Proof.* Indeed, replace each square  $A \in \mathcal{S}(r)$  in the decomposition of V by a square  $A' \in \mathcal{S}(r+2)$  that contains A and all nearest neighbors of A. Let  $V'_{r+2,A'} = V_{r,A}$ . Then  $V = \sum_{r>1} \sum_{A \in \mathcal{S}(r)} V'_{r,A}$  is the desired decomposition.

#### 4.3 Proof of the main theorem

In this section we prove Theorem 1.

Proof. We shall derive simplified flow equations for a triple of parameters J(n),  $J_d(n)$ , and  $\mu(n)$  such that V(n) is  $(J(n), \mu(n), \beta)$ -decaying and W(n) is  $(J_d(n), \mu(n), \alpha)$ -decaying for all  $n \geq 0$ . Here  $\alpha$ ,  $\beta$  are sufficiently large constants that will be chosen later. Our manipulations with local decompositions will typically decrease the constants  $\alpha$  and  $\beta$ , so after each step of the flow equations we shall need to reset these constants back to their original values using Lemma 4.2.

Since W(0) = 0 we can choose initial conditions

$$J(0) = J, \quad J_d(0) = 0, \quad \text{and} \quad \mu(0) = \mu.$$
 (4.16)

Let us prove that for any constant  $0 < \epsilon < 1$  there exist constants  $c_1, c_2, c_3 > 0$  such that for all  $k \ge 0$  one has

$$J(k+1) \le c_1 J(k)^{2(1-\epsilon)}, \tag{4.17}$$

$$J_d(k+1) \leq c_2 J(k)^{1-\epsilon}, \tag{4.18}$$

$$\mu(k+1) = \frac{1}{2}\mu(k) - c_3 J(k+1)^{\frac{\epsilon}{10}}, \tag{4.19}$$

$$||E(k+1)|| \le ||E(k)|| + O(L^3)J(k)e^{-c_3L\mu(k)}$$
 (4.20)

provided that J(0) is below some constant threshold value. Note that although  $\epsilon$  can be chosen arbitrarily close to 0, Eq. (4.19) does not permit one to choose  $\epsilon = 0$  since otherwise the decay rate  $\mu(n)$  becomes negative after O(1) iterations.

Supposed we have already proved Eqs. (4.17-4.20) for k = 0, 1, ..., n. Let us denote  $V \equiv V(n), E \equiv E(n),$ 

$$W \equiv \sum_{k=0}^{n} W(k)$$

and

$$J_d \equiv \sum_{k=0}^n J_d(k).$$

Since  $\mu(k)$  is monotone decreasing for k = 0, ..., n, we can safely assume that W is  $(J_d, \mu(n), \alpha)$ -decaying. Also, since J(k) decreases doubly exponentially for k = 0, ..., n we can assume (for k > 0) that

$$J_d = O(J_d(1)) = O(J^{1-\epsilon}).$$
 (4.21)

Let S be the solution of the linearized block-diagonalization problem Eq. (4.4) constructed in Lemma 5.1. By construction S is  $(K, \mu(n), \beta)$ -decaying, where

$$K = \frac{c_1 J(n)}{1 - c_2 J_d} = O(J(n)). \tag{4.22}$$

Here we assumed that J is sufficiently small, so that  $c_2J_d \leq 1/2$ . Note that the assumptions of Lemma 5.1 require  $\beta \geq 2$  and  $\alpha \geq \beta + 4$ .

Recall that  $\tilde{W}$  and  $\tilde{V}$  are raw versions of W(n+1) and V(n+1) defined in Eq. (4.5) and Eq. (4.6). From Corollary 5.1, Section 5.3, we infer that the operator  $\tilde{W}$  has a local decomposition with block-diagonal terms which is  $(cJ(n), \mu(n))$ -decaying. Note that the assumptions of Corollary 5.1 require  $\beta \geq 2$  and  $\alpha \geq \beta + 4$ . The operator  $\tilde{V}$  can be rewritten as

$$\tilde{V} = [S, V] + \omega(H), \tag{4.23}$$

where  $H \equiv H_0 + W + V$  and  $\omega(H) \equiv e^S H e^{-S} - H - [S, H]$ . We can assume that H is  $(c, \mu(n), \beta)$ -decaying where c = O(1) and we assumed that  $\beta \leq \alpha$ . Applying Lemma 4.1 from Section 4.2 we get a local decomposition for [S, V] which is  $(c_1 J(n)^2, \mu(n))$ -decaying for some constant  $c_1$  provided that  $\beta \geq 4$ . Applying Lemma 6.1 from Section 6 we get a local decomposition for  $\omega(H)$  which is  $(c_2 J(n)^2, \mu(n)/2, -1)$ -decaying for some constant  $c_2$  provided that  $\beta \geq 6$ . Summarizing,  $\tilde{V}$  is  $(c_3 J(n)^2, \mu(n)/2, -1)$ -decaying where c is a constant. It will be convenient to keep the decay rates of  $\tilde{W}$  and  $\tilde{V}$  the same. Thus we shall assume that  $\tilde{W}$  is  $(c_3 J(n), \mu(n)/2)$ -decaying (which is a weaker version of what we proved above). One can easily check that a choice

$$\alpha = 10 \quad \text{and} \quad \beta = 6 \tag{4.24}$$

satisfies conditions of all lemmas used above.

Recall that W(n+1) and V(n+1) are defined by taking the local decompositions of  $\tilde{W}$  and  $\tilde{V}$  and removing all terms associated with squares of size larger than  $L^*$ . Therefore W(n+1) and V(n+1) have the same decay parameters as  $\tilde{W}$  and  $\tilde{V}$ , that is, we get

$$J(n+1) \le c_1 J(n)^2$$
,  $J_d(n+1) \le c_2 J(n)$ ,  $\mu(n+1) = \frac{1}{2}\mu(n)$  (4.25)

for some constants  $c_1, c_2$ . Note that we have not reset  $\alpha, \beta$  to their original values yet. The total number of squares of size  $r > L^*$  is at most  $L^3$ . Thus the contribution to the error Hamiltonian E(n+1) can be estimated as

$$||E(n+1)|| \le ||E(n)|| + O(L^3)J(n)e^{-\mu(n)L^*/2} = ||E(n)|| + O(L^3)J(n)e^{-c_3\mu(n)L}$$
(4.26)

for some constant  $c_3$ . Resetting the constants  $\alpha, \beta$  using Lemma 4.2 we arrive at the desired flow equations Eqs. (4.17-4.20).

Solving Eq. (4.17) we get

$$J(n) \le J\left(\frac{J}{J_0}\right)^{\theta^n}, \quad \theta = 2(1 - \epsilon)$$
 (4.27)

for some constant  $J_0 > 0$ . Note that we are free to choose the constant  $\epsilon$  as small as possible. To simplify the formulas let us set  $\theta = 2$ . Also note that for small enough  $J_0$  we can assume that  $\mu(n)$  decays exponentially with exponent arbitrarily close to 1/2. To simplify the formulas let us assume that

$$J(n) \le J\left(\frac{J}{J_0}\right)^{2^n}$$
 and  $\mu(n) = \mu 2^{-n}$ . (4.28)

Simple algebra shows that choosing the number of steps  $n = c \log L$  for some constant c one can achieve the bounds

$$||V(n)||, ||E(n)|| \le poly(L) \exp(-c\sqrt{L}).$$
 (4.29)

Here the exponent  $-c\sqrt{L}$  is determined by a tradeoff between the doubly exponential decay of J(n) and the exponential decay of  $\mu(n)$ . Neglecting these exponentially small errors we can assume for simplicity that the Hamiltonian H(n) contains only block-diagonal contributions, that is,

$$H(n) = H_0 + W, \quad W = \sum_{k=0}^{n} W(k).$$
 (4.30)

Here the local decomposition of W is  $(J_d, 0, \alpha)$ -decaying with  $J_d = O(J^{1-\epsilon})$ ,  $\alpha = 10$ , and all terms in this decomposition are individually block-diagonal. In addition, by definition of W(k), this decomposition contains only squares of size  $r \leq L^*$ . By performing an overall energy shift and using TQO-1 we can guarantee that every local term in the decomposition of W(k) has zero restriction on the P subspace (note that it increases the strength  $J_d$  at most by a factor of two). It allows us to apply the machinery of Section 3. In particular, Lemma 3.2 says that W is relatively bounded by  $H_0$  with a constant

$$b \le O(J_d) \sum_{r \ge 1} r^{2-\alpha} = O(J_d).$$
 (4.31)

Assuming that b < 1, Lemma 3.1 implies that the spectrum of  $H_0 + W$  is contained in the union of intervals  $I_k = (k(1-b), k(1+b))$ , where  $k = 0, 1, 2, \ldots$  The effect of V(n) and E(n) can now be taken into account using the standard perturbation theory by considering  $H_0 + W$  as an unperturbed Hamiltonian and using Eq. (4.29).

Finally, notice that since the dominant contribution to W comes from the first level k = 1, we do not really need to perform the degree reset to estimate b (since the degree reset only changes our description of an operator but does not change the operator itself). Hence we can set  $J_d = O(J)$ . The theorem is proved.

# 5 Linearized block-diagonalization problem

## 5.1 Statement of the problem

Consider a pair of perturbations

$$V = \sum_{1 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} V_{r,A}, \tag{5.1}$$

$$W = \sum_{1 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} W_{r,A},\tag{5.2}$$

such that all terms  $W_{r,A}$  are block-diagonal,

$$QW_{r,A}P = 0 \quad \text{for all } r, A. \tag{5.3}$$

A linearized block-diagonalization problem can be divided into two parts. The first part is to find an anti-hermitian operator S such that

$$Q([S, H_0 + W] + V)P = 0, \quad S^{\dagger} = -S$$
 (5.4)

and construct a local decomposition of S. The second part is to construct a local decomposition for the transformed Hamiltonian

$$\tilde{W} = [S, H_0 + W] + V = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} \tilde{W}_{r,A}$$
 (5.5)

such that every term in the local decomposition is block-diagonal, that is,  $Q\tilde{W}_{r,A}P=0$  for all r,A. We solve the two parts of the problem in Lemma 5.1 and its Corollary 5.1 respectively. Throughout this section we assume that  $H_0$  is a Hamiltonian defined in Eq. (2.1) that obeys conditions TQO-1 and TQO-2.

## 5.2 Finding the transformation S

In this section we prove the following

**Lemma 5.1.** Let V and W be perturbations defined in Eqs. (5.1,5.2,5.3). Suppose that V is  $(J, \mu, \beta)$ -decaying and W is  $(J_d, \mu, \alpha)$ -decaying such that  $\beta \geq 2$  and  $\alpha \geq \beta + 4$ . Then there exist constants  $c_1, c_2 > 0$  such that Eq. (5.4) has a solution S which is  $(K, \mu, \beta)$ -decaying with

$$K = \frac{c_1 J}{1 - c_2 J_d},\tag{5.6}$$

*Proof.* Let us start from several simplifying assumptions. Without loss of generality we can assume that

$$W_{r,A}P = 0 \quad \text{for all } r, A. \tag{5.7}$$

Indeed, condition TQO-1 guarantees that  $W_{r,A}P$  is a multiple of P. Shifting  $W_{r,A}$  by the corresponding multiple of the identity we can satisfy Eq. (5.7). On the other hand, one can easily check that Eq. (5.4) in invariant under such a shift. Note also that the shift can increase the norm  $||W_{r,A}||$  at most by a factor of two. Indeed, if  $W_{r,A}P = cP$  then  $|c| \leq ||W_{r,A}||$  and thus  $||W_{r,A} - cI|| \leq ||W_{r,A}|| + |c| \leq 2||W_{r,A}||$ . Thus we can assume that W satisfies Eq. (5.7) provided that we change  $J_d$  to  $2J_d$  in the final answer. By the same reasons, we can assume that

$$PV_{r,A}P = 0 \quad \text{for all } r, A \tag{5.8}$$

provided that we change J to 2J in the final answer. In addition, we can assume that  $V_{r,A}$  commutes with  $G_B$ ,  $B \in \mathcal{S}(2)$ , whenever B is not contained in A,

$$[V_{r,A}, G_B] = 0$$
 for all  $B \in \mathcal{S}(2)$  such that  $B \cap A^c \neq \emptyset$ . (5.9)

Indeed, by adding an idle layer of sites to each square in the decomposition of V as explained in Lemma 4.3 we guarantee Eq. (5.9). The price we pay for this simplification is that J has to be changed to  $cJe^{2\mu} = O(J)$  in the final answer, see Lemma 4.3. Note also that now the local decomposition of V in Eq. (5.1) starts from squares of size r = 3,

$$V = \sum_{3 \le r \le L^*} \sum_{A \in \mathcal{S}(r)} V_{r,A}. \tag{5.10}$$

We shall construct a solution S as a series  $S = \sum_{i=1}^{\infty} S^{(i)}$ , where  $S^{(i)}$  is anti-hermitian for all i and

$$Q([S^{(1)}, H_0] + V)P = 0, (5.11)$$

$$Q([S^{(i+1)}, H_0] + [S^{(i)}, W])P = 0, \text{ for } i \ge 1.$$
(5.12)

For any region  $A \subseteq \Lambda$  define a restricted Hamiltonian

$$H_0(A) = \sum_{\substack{B \in \mathcal{S}(2) \\ B \subseteq A}} G_B. \tag{5.13}$$

Define also a super-operator  $\mathcal{E}_A$  that takes as argument an arbitrary operator O and returns an operator

$$\mathcal{E}_A(O) = Q_A H_0(A)^{-1} O P_A - P_A O H_0(A)^{-1} Q_A.$$
(5.14)

Note that the  $P_A$  is the zero-subspace of  $H_0(A)$ , so that  $Q_A H_0(A)^{-1}$  is well-defined.

**Proposition 5.1.** Let  $O_A$  be any operator acting on A such that  $PO_AP = 0$ . Suppose  $O_A$  commutes with  $G_B$ ,  $B \in \mathcal{S}(2)$ , whenever B is not contained in A. Then

$$Q([\mathcal{E}_A(O_A), H_0] + O_A) P = 0.$$
(5.15)

If  $O_A$  is hermitian then  $\mathcal{E}_A(O_A)$  is anti-hermitian.

*Proof.* Indeed, since all terms in  $H_0$  which are not contained in A commute with  $\mathcal{E}_A(O_A)$  while  $H_0(A)P_A=0$  we have

$$[\mathcal{E}_A(O_A), H_0] = -Q_A O_A P_A - P_A O_A Q_A$$

It yields

$$Q[\mathcal{E}_A(O_A), H_0]P = -QQ_AO_AP = -Q(Q_A + P_A)O_AP = -QO_AP,$$

where the first equality follows from  $P_AP = P$ ,  $Q_AP = 0$ , and the second equality uses identity  $P_AO_AP = PO_AP = 0$ . Thus we have proved Eq. (5.15). The last statement of the proposition is obvious.

Using the assumptions Eqs. (5.8,5.9) and the proposition we can choose  $S^{(1)}$  in Eq. (5.11) as

$$S^{(1)} = \sum_{r \ge 3} \sum_{A \in \mathcal{S}(r)} S_{r,A}^{(1)}, \quad S_{r,A}^{(1)} = \mathcal{E}_A(V_{r,A}). \tag{5.16}$$

Taking into account that

$$\|\mathcal{E}_A(O_A)\| \le \|O_A\| \quad \text{for any } O_A \tag{5.17}$$

we conclude that Eq. (5.16) is a local decomposition of  $S^{(1)}$  which is  $(K_1, \mu, \beta)$ -decaying where

$$K_1 = J. (5.18)$$

This decomposition has an extra property that  $S_{r,A}^{(1)}$  is block-off-diagonal with respect to  $P_A$ ,  $Q_A$ , and  $S_{r,A}^{(1)}$  commutes with  $G_B$ ,  $B \in \mathcal{S}(2)$ , whenever B is not contained in A.

Let us now solve Eq. (5.12). We shall assume as our induction hypothesis that  $S^{(i)}$  possesses a local decomposition

$$S^{(i)} = \sum_{p \ge 3} \sum_{A \in \mathcal{S}(p)} S_{p,A}^{(i)} \tag{5.19}$$

such that

- **I1**  $S_{p,A}^{(i)}$  is block-off-diagonal with respect to  $P_A$ ,  $Q_A$
- **12**  $S_{p,A}^{(i)}$  commutes with  $G_B$ ,  $B \in \mathcal{S}(2)$ , whenever B is not contained in A
- **13**  $\| [H_0(A), S_{p,A}^{(i)}] \| \le K_i p^{-\beta} e^{-\mu p}$

Here the coefficient  $K_i$  will be determined inductively in terms of  $K_{i-1}$ . Note that combining (I1), (I3) with the fact that  $Q_A H_0(A) \geq I$  one gets

$$||S_{p,A}^{(i)}|| = ||Q_A S_{p,A}^{(i)} P_A|| \le ||Q_A H_0(A) S_{p,A}^{(i)} P_A|| = ||[H_0(A), S_{p,A}^{(i)}]|| \le K_i p^{-\beta} e^{-\mu p},$$

that is,  $S^{(i)}$  is  $(K_i, \mu, \beta)$ -decaying.

The base of induction is i = 1 which we have already proved. Our first step will be choosing a local decomposition for the commutator  $[S^{(i)}, W]$  in Eq. (5.12). We shall need the following geometrical fact.

**Proposition 5.2.** Let  $A \in \mathcal{S}(p)$ ,  $p \geq 2$ , and  $B \in \mathcal{S}(q)$ , q < L, be any squares such that  $A \cap B \neq \emptyset$ . Then there exists a square  $C \in \mathcal{S}(r)$ ,  $r = \min(p + q, L)$ , such that  $A \cup B \subseteq C$  and C contains all nearest neighbors of B.

*Proof.* If r = L the statement is obvious, so assume r = p + q < L. Define a metric on the lattice using the  $l_{\infty}$ -norm, that is, if  $u = (u_x, u_y)$  and  $v = (v_x, v_y)$  is a pair of sites, then

$$D(u, v) = \max\{|u_x - v_x|, |u_y - v_y|\}.$$

For any region  $M \subseteq \Lambda$  let D(M) be the diameter of M, i.e., the largest distance between a pair of sites  $u, v \in M$ . Clearly D(A) = p - 1 and D(B) = q - 1. Since  $A \cap B \neq \emptyset$  we have  $D(A \cup B) \leq D(A) + D(B) = p + q - 2$ . Therefore,  $A \cup B$  can be covered by a square  $C' \in \mathcal{S}(p+q-1)$ . If C' = B one actually has  $C' \in \mathcal{S}(q)$ . Now one can choose C as an arbitrary square of size r that contains C' and all nearest neighbors of C'. Suppose now that  $C' \neq B$ . Then either C' contains all nearest neighbors of B, or C' shares an edge or a corner with B. In the latter case, we can extend the size of C' by one obtaining a square  $C \in \mathcal{S}(p+q)$  with the desired properties.

The proposition implies that

$$[S^{(i)}, W] = \sum_{r>3} \sum_{C \in \mathcal{S}(r)} D_{r,C}^{(i)}, \tag{5.20}$$

where

$$D_{r,C}^{(i)} = \sum_{p+q=r} \sum_{\substack{A \in \mathcal{S}(p) \\ A \subset C}} \sum_{\substack{B \in \mathcal{S}(q) \\ B \subset C}} [S_{p,A}^{(i)}, W_{q,B}] \chi(A, B, C)$$
 (5.21)

and  $\chi(A, B, C) = 0, 1$  is some function that 'distributes' the commutators over different terms  $D_{r,C}^{(i)}$ . A particular choice of such distribution is not important for us. Using Proposition 5.2 we can assume that  $\chi(A, B, C) = 0$  unless  $A \cup B \subseteq C$  and C contains all nearest neighbors of B. By construction,  $D_{r,C}^{(i)}$  has support on C, that is, Eqs. (5.20,5.21) define a local decomposition of the commutator  $[S^{(i)}, W]$ .

Using Proposition 5.1 we can choose a solution  $S^{(i+1)}$  of Eq. (5.12) as

$$S^{(i+1)} = \sum_{r>3} \sum_{C \in \mathcal{S}(r)} S_{r,C}^{(i+1)}, \quad S_{r,C}^{(i+1)} = \mathcal{E}_C(D_{r,C}^{(i)}). \tag{5.22}$$

This is a local decomposition of  $S^{(i+1)}$  that satisfies (I1) by definition of the map  $\mathcal{E}_C$ . Let us check that it satisfies (I2). Indeed, let  $G_F$ ,  $F \in \mathcal{S}(2)$ , be such that F is not contained in C. Then F is not contained in  $A \subseteq C$  and thus  $G_F$  commutes with all  $S_{p,A}^{(i)}$  in Eq. (5.21). Since for all terms  $W_{q,B}$  in Eq. (5.21) the square C contains both B and the nearest neighbors of B, we conclude that F does not overlap with B, that is,  $G_F$  commutes with  $W_{q,B}$ . Finally,  $G_F$  commutes with all Krauss operators involved in the map  $\mathcal{E}_C$ . Thus  $G_F$  commutes with  $S_{r,C}^{(i+1)}$  which proves (I2).

It remains to verify that the local decomposition Eq. (5.22) satisfies (I3). It will be convenient to introduce an auxiliary Hamiltonian

$$W_q(A,C) = \sum_{\substack{B \in \mathcal{S}(q) \\ B \subseteq C}} \chi(A,B,C) W_{q,B}. \tag{5.23}$$

It allows us to rewrite Eq. (5.21) as

$$D_{r,C}^{(i)} = \sum_{p+q=r} \sum_{\substack{A \in \mathcal{S}(p) \\ A \subseteq C}} [S_{p,A}^{(i)}, W_q(A, C)]. \tag{5.24}$$

Applying Corollary 3.2, using Eq. (5.7), and taking into account that W is  $(J_d, \mu, \alpha)$ -decaying we get

$$\| [H_{0}(C), S_{r,C}^{(i+1)}] \| = \| Q_{C}D_{r,C}^{(i)}P_{C}\| \le \sum_{p+q=r} \sum_{\substack{A \in \mathcal{S}(p) \\ A \subseteq C}} \| Q_{C}[S_{p,A}^{(i)}, W_{q}(A, C)]P_{C}\|$$

$$\le \sum_{p+q=r} \sum_{\substack{A \in \mathcal{S}(p) \\ A \subseteq C}} b_{q} \| Q_{C}[H_{0}(C), S_{p,A}^{(i)}]P_{C}\|,$$

$$(5.25)$$

where

$$b_q \le cJ_d q^{2-\alpha} e^{-\mu q} \tag{5.26}$$

for some constant c. From (I1) we infer that  $S_{p,A}^{(i)}$  is block-off-diagonal with respect to  $P_A$ ,  $Q_A$  while (I2) implies  $[H_0(C), S_{p,A}^{(i)}] = [H_0(A), S_{p,A}^{(i)}]$ . Taking into account that  $P_C = P_A P_C$  we get a bound

$$\|Q_C[H_0(C), S_{p,A}^{(i)}]P_C\| \le \|Q_A[H_0(A), S_{p,A}^{(i)}]P_A\| = \|[H_0(A), S_{p,A}^{(i)}]\| \le K_i p^{-\beta} e^{-\mu p},$$
 (5.27)

where the last inequality follows from (I3). Combining Eqs. (5.25,5.27) and noting that the number of squares  $A \in \mathcal{S}(p)$  such that  $A \subseteq C$  is equal to  $(r-p)^2 = q^2$  we get

$$\| [H_0(C), S_{r,C}^{(i+1)}] \| \le K_i \sum_{q=1}^{r-1} q^2 b_q (r-q)^{-\beta} e^{-\mu(r-q)}$$

$$\le c J_d K_i e^{-\mu r} \sum_{q=1}^{r-1} q^{4-\alpha} (r-q)^{-\beta}.$$
(5.28)

It is convenient to split the sum over q into two parts:

$$\sum_{1 \le q \le r/2} q^{4-\alpha} (r-q)^{-\beta} \le cr^{-\beta} \sum_{q \ge 1} q^{4-\alpha} = c'r^{-\beta}$$
 (5.29)

for some constants c, c' since we assumed  $\alpha \geq \beta + 4 \geq 6$ . As for the other part, we have

$$\sum_{r/2 \le q \le r-1} q^{4-\alpha} (r-q)^{-\beta} \le cr^{4-\alpha} \sum_{q \ge 1} q^{-\beta} = c'r^{4-\alpha} \le c'r^{-\beta}$$
 (5.30)

for some constants c, c' since we assumed that  $\beta \geq 2$  and  $\alpha \geq \beta + 4$ . Therefore we arrive at

$$\|[H_0(C), S_{r,C}^{(i+1)}]\| \le cJ_dK_i r^{-\beta} e^{-\mu r}$$
 (5.31)

for some constant c which proves (I3) for

$$K_{i+1} = cJ_d K_i. (5.32)$$

Thus all induction assumptions are proved for  $S^{(i+1)}$ . By obvious reasons  $S = \sum_{i \geq 1} S_i$  is  $(K, \mu, \beta)$ -decaying with

$$K \le \sum_{i>1} K_i = \frac{c_1 J}{1 - c_2 J_d}.$$
 (5.33)

## 5.3 Local decomposition of the transformed Hamiltonian

Lemma 5.1 has the following corollary.

Corollary 5.1. Let V, W, and S be as in Lemma 5.1. Suppose V is  $(J, \mu, \beta)$ -decaying and W is  $(J_d, \mu, \alpha)$ -decaying for some  $\beta \geq 2$  and  $\alpha \geq \beta + 4$ . Then a transformed Hamiltonian  $\tilde{W} = [S, H_0 + W] + V$  has a local decomposition

$$\tilde{W} = \sum_{r \ge 1} \sum_{A \in \mathcal{S}(r)} \tilde{W}_{r,A} \tag{5.34}$$

such that  $Q\tilde{W}_{r,A}P = 0$  for all r,A. This decomposition is  $(\tilde{J}_d,\mu)$ -decaying, where

$$\tilde{J}_d \le \frac{cJ}{1 - cJ_d} \tag{5.35}$$

for some constant c.

*Proof.* We shall use notations and techniques introduced in the proof of Lemma 5.1. By definition of S we have

$$\tilde{W} = \sum_{i \ge 1} W^{(i)},\tag{5.36}$$

where

$$W^{(1)} = [S^{(1)}, H_0] + V$$
, and  $W^{(i)} = [S^{(i)}, H_0] + [S^{(i-1)}, W]$  for  $i \ge 2$ . (5.37)

Let us choose the local decomposition of  $\mathcal{W}^{(1)}$  as

$$W^{(1)} = \sum_{r \ge 3} \sum_{A \in \mathcal{S}(r)} W_{r,A}^{(1)}, \tag{5.38}$$

where

$$W_{r,A}^{(1)} = \left[\mathcal{E}_A(V_{r,A}), H_0\right] + V_{r,A} = P_A V_{r,A} P_A + Q_A V_{r,A} Q_A. \tag{5.39}$$

Here the last equality uses Proposition 5.1 (see the first equation in the proof of the proposition). Obviously,  $||W_{r,A}^{(1)}|| \le ||V_{r,A}||$  and thus  $W^{(1)}$  is  $(J, \mu, \beta)$ -decaying. As was noticed in the proof of Lemma 5.1 we can assume that  $V_{r,A}$  commutes with all operators  $G_B$ ,  $B \in \mathcal{S}(2)$  for which B is not contained in A. It means that

$$PW_{r,A}^{(1)} = PV_{r,A}P_A = PV_{r,A}P = W_{r,A}^{(1)}P,$$
(5.40)

that is,  $W_{r,A}^{(1)}$  is block-diagonal.

Recall that we have a decomposition

$$[S^{(i)}, W] = \sum_{r \ge 3} \sum_{C \in \mathcal{S}(r)} D_{r,C}^{(i)}, \tag{5.41}$$

where  $D_{r,C}^{(i)}$  commutes with all operators  $G_B$ ,  $B \in \mathcal{S}(2)$  for which B is not contained in C, see Eqs. (5.20,5.21) in the proof of Lemma 5.1. It means that we can choose the local decomposition of  $W^{(i)}$  with  $i \geq 2$  as

$$W^{(i)} = \sum_{r>3} \sum_{C \in \mathcal{S}(r)} W_{r,C}^{(i)}, \tag{5.42}$$

where

$$W_{r,C}^{(i)} = \left[\mathcal{E}_C(D_{r,C}^{(i-1)}), H_0\right] + D_{r,C}^{(i-1)} = P_C D_{r,C}^{(i-1)} P_C + Q_C D_{r,C}^{(i-1)} Q_C, \tag{5.43}$$

see Eq. (5.22). Obviously,  $||W_{r,C}^{(i)}|| \leq ||D_{r,C}^{(i-1)}||$  and  $W_{r,C}^{(i)}$  is block-diagonal. Using the fact that  $S^{(i)}$  is  $(K_i, \mu, \beta)$ -decaying where  $K_i$  is defined by Eqs. (5.18,5.32), and using Lemma 4.1 we conclude that  $W^{(i)}$  is  $(cJ_dK_{i-1}, \mu)$ -decaying. Using Eq. (5.36) we obtain a local decomposition of  $\tilde{W}$  with individually block-diagonal local terms which is  $(\tilde{J}_d, \mu)$ -decaying with

$$\tilde{J}_d = J + \sum_{i>2} cJ_d K_{i-1} \le J + \sum_{i>1} J(cJ_d)^i = \frac{J}{1 - cJ_d}.$$
 (5.44)

Finally we have to replace J and  $J_d$  by O(J) and  $O(J_d)$  to justify our assumptions Eqs. (5.7,5.8,5.9).

## 6 Lieb-Robinson bounds

Let S be some anti-hermitian operator and V be an arbitrary operator. We shall use notations

$$\tau(V) = e^{S}Ve^{-S},$$
  

$$\omega(V) = \tau(V) - V - [S, V].$$
(6.1)

Suppose we are given some local decompositions of S and V. In order to get a closed system of flow equations we need to construct a local decomposition for  $\omega(V)$ , see Section 4. The main result of this section is the following lemma.

**Lemma 6.1.** Suppose S is  $(K, \mu, \alpha)$ -decaying for some  $\alpha \geq 6$ . Suppose V is  $(J, \mu, \beta)$ -decaying for some  $\beta \geq 6$ . Then  $\omega(V)$  has a local decomposition which is  $(cJK^2, \mu/2, -1)$ -decaying for some constant c.

Thus the magnitude of interactions of range r in  $\omega(V)$  decays as  $cJK^2re^{-\mu r/2}$ . Our arguments will rely on the Lieb-Robinson bound, see [20]. More specifically, we shall exploit quasi-locality of dynamics in quantum spin systems with a fast decay of interactions as presented in [19]. The extra factor 1/2 in the decay rate of  $\omega(V)$  represents a simple geometrical fact that the diameter of the light-cone of any local region increases with a rate 2v, where v is the Lieb-Robinson velocity. This extra factor 1/2 is the price one has to pay for using the powerful machinery built on the Lieb-Robinson bound.

We shall start from solving a somewhat simpler problem. Let O be any operator with support on some square  $B \in \mathcal{S}(q)$ . Consider a local decomposition of S,

$$S = \sum_{p>2} \sum_{A \in \mathcal{S}(p)} S_{p,A}. \tag{6.2}$$

Let  $C \in \mathcal{S}(q+2j)$  be a square that contains B and all sites within distance j from B (with respect to the  $l_{\infty}$ -distance), that is,

$$C = \{ u \in \Lambda : D(u, B) \le j \}. \tag{6.3}$$

Let  $S_C$  be a localized version of S obtained by taking out all interactions whose support is not contained in C,

$$S_C = \sum_{p \ge 2} \sum_{\substack{A \in \mathcal{S}(p) \\ A \subseteq C}} S_{p,A}. \tag{6.4}$$

Define also a localized version of  $\omega(O)$ , that is,

$$\omega_C(O) = e^{S_C} O e^{-S_C} - O - [S_C, O]. \tag{6.5}$$

By definition,  $\omega_C(O)$  has support on C. We shall need a bound on the difference  $\|\omega(O) - \omega_C(O)\|$  that is proportional to  $\|O\| \cdot K^2 e^{-\mu j}$ .

**Lemma 6.2** (Quasi-Local Dynamics). Let S,  $S_C$  and O be the operators defined above. Suppose S is  $(K, \mu, \alpha)$ -decaying for some  $\alpha \geq 6$ . Then there exist constants  $c_0, c_1 > 0$  such that

$$\|\omega_C(O) - \omega(O)\| \le c_0(q+j)q^4K^2\|O\|e^{-\mu j}$$
 (6.6)

whenever  $K \leq c_1$ .

*Proof.* Define

$$\tau^{t}(O) = e^{St}Oe^{-St}, \quad \tau_{C}^{t}(O) = e^{S_{C}t}Oe^{-S_{C}t}.$$
 (6.7)

For any  $0 \le t \le 1$  define an operator

$$f(t) = \tau_C^t(O) - \tau^t(O) - [S_C - S, O]t.$$
(6.8)

Note that f(t) is an analytic function and

$$f(0) = \dot{f}(0) = 0, \quad f(1) = \omega_C(O) - \omega(O).$$
 (6.9)

Computing the derivatives over t we get

$$\dot{f}(t) = [S_C, \tau_C^t(O)] - [S, \tau^t(O)] - [S_C - S, O], \tag{6.10}$$

$$\ddot{f}(t) = [S_C, [S_C, \tau_C^t(O)]] - [S, [S, \tau^t(O)]]. \tag{6.11}$$

Let us extract from  $\ddot{f}(t)$  a norm preserving term  $[S,\dot{f}(t)]$ . After simple algebra we get

$$\ddot{f}(t) = [S, \dot{f}(t)] - [S - S_C, [S_C, \tau_C^t(O)]] - [S, [S - S_C, O]]. \tag{6.12}$$

It means that

$$\dot{f}(t) = \tau^t \left( \int_0^t \tau^{-s} \left( [[S_C, \tau_C^s(O)], \Delta S_C] + [[\Delta S_C, O], S] \right) ds \right), \tag{6.13}$$

where  $\Delta S_C \equiv S - S_C$ . Taking into account the initial conditions Eq. (6.9) we get

$$||f(1)|| \le \int_0^1 dt_1 \, ||\dot{f}(t_1)|| \le \frac{1}{2} ||[[\Delta S_C, O], S]|| + \int_0^1 dt_1 \int_0^{t_1} dt_2 \, ||[\tau_C^{t_2}([S_C, O]), \Delta S_C]||. \tag{6.14}$$

Let us start from bounding the time-independent term  $\|[[\Delta S_C, O], S]\|$ . Denote

$$\Gamma_{p_2,p_1} = \#\{(E_2,E_1) : E_2 \in \mathcal{S}(p_2), \quad E_1 \in \mathcal{S}(p_1), \quad E_2 \cap (E_1 \cup B) \neq \emptyset, \quad E_1 \cap B \neq \emptyset\}.$$

One can easily check that

$$\Gamma_{p_2,p_1} \le (q+p_1+p_2)^2(q+p_1)^2 \le 2(q+p_1)^4 + 2(q+p_1)^2p_2^2$$

The commutator  $[\Delta S_C, O]$  has contributions only from squares of size  $\geq j$  in the decomposition of S, which implies

$$\|[[\Delta S_C, O], S]\| \le cK^2 \|O\| \sum_{p_1 \ge j} \sum_{p_2 \ge 1} \Gamma_{p_2, p_1} p_1^{-\alpha} p_2^{-\alpha} e^{-\mu(p_1 + p_2)}.$$

Here and below c stands for a constant factor. Since  $\alpha \geq 6$  the sum over  $p_2$  is bounded by a constant and we arrive at

$$\|[[\Delta S_C, O], S]\| \le cK^2 \|O\| e^{-\mu j} \sum_{p_1 \ge j} (q + p_1)^4 p_1^{-\alpha} = cq^4 K^2 \|O\| e^{-\mu j}.$$
(6.15)

Let us now bound  $\|[\tau_C^{t_2}([S_C, O]), \Delta S_C]\|$ . Note that the commutator has contributions only from those terms in the decomposition of  $\Delta S_C$  that overlap with both C and its complement  $C^c$ . Define

$$\Omega_{p_1,p_2}(t) = \max_{\substack{E_1 \in \mathcal{S}(p_1) \\ E_1 \cap B \neq \emptyset}} \max_{\substack{E_2 \in \mathcal{S}(p_2) \\ E_2 \cap C^c \neq \emptyset}} \max_{O_1,O_2} \| [\tau_C^t([O_1,O]),O_2] \|,$$

where  $O_1, O_2$  are unit-norm operators acting on  $E_1, E_2$  respectively. Counting the number of squares contributing to the double commutator yields

$$\| \left[ \tau_C^t([S_C, O]), \Delta S_C \right] \| \le c(q+j)K^2 \sum_{p_1, p_2 \ge 1} (q+p_1)^2 p_1^{-\alpha} p_2^{2-\alpha} e^{-\mu(p_1+p_2)} \Omega_{p_1, p_2}(t).$$

As we show below, the unitary evolution under  $S_C$  can be characterized by a finite Lieb-Robinson velocity  $v_{LR} = O(K)$ . Using the Lieb-Robinson bound from [19] that governs unitary evolution under Hamiltonians with exponentially decaying interactions one gets

$$\Omega_{p_1,p_2}(t) \le c \|O\|(p_1+q)^2 p_2^2 \exp[v_{LR}t - \mu\theta(j-p_1-p_2)],$$

where  $\theta(x) = x$  for  $x \geq 0$  and  $\theta(x) = 0$  for x < 0. Note that  $\theta(j - p_1 - p_2)$  is a lower bound on the distance between supports of  $[O_1, O]$  and  $O_2$  in the definition of  $\Omega_{p_1, p_2}$ . Clearly,  $p_1 + p_2 + \theta(j - p_1 - p_2) \geq j$ . Since  $0 \leq t \leq 1$  we can assume that  $v_{LR}t = O(Kt) = O(1)$ . Hence

$$\| [\tau_C^t([S_C, O]), \Delta S_C] \| \le c(q+j) \|O\| K^2 e^{-\mu j} \sum_{p_1, p_2 \ge 1} (q+p_1)^4 p_1^{-\alpha} p_2^{4-\alpha} \le c(q+j) q^4 \|O\| K^2 e^{-\mu j}$$

$$(6.16)$$

for  $\alpha \geq 6$ . Combining Eqs. (6.15,6.16) and computing the integrals in Eq. (6.14) we arrive at

$$|f(1)| \le c(q+j)q^4||O||K^2e^{-\mu j}.$$

It remains to check that we have fast enough decay of interactions in S to use the Lieb-Robinson bound from Ref. [19]. Below, we use notations from Ref. [19]. Define a function

$$F_{\mu}(x) = \frac{\exp(-\mu x)}{1 + x^2}.$$
 (6.17)

As was shown in Ref. [19] the Lieb-Robinson velocity is bounded by a multiple of:

$$||S||_{\mu} := \sup_{u,v \in \Lambda} \sum_{r \ge 1} \sum_{\substack{A \in \mathcal{S}(r) \\ A \ni u,v}} \frac{||S_{r,A}||}{F_{\mu}(D(u,v))},\tag{6.18}$$

with the multiplying factor being  $2C_0(1)$ , where  $C_0(1)$  is a numerical constant (depending only on the dimensionality of the system; in our case the dimension is 2.) For any pair of sites  $u, v \in A$  with  $A \in \mathcal{S}(r)$  one has  $D(u,v) \leq r$  (here and below we use the  $l_{\infty}$ -distance). Since  $F_{\mu}(x)$  is monotone-decreasing we have  $F_{\mu}(D(u,v)) \geq F_{\mu}(r)$ . Taking into account that the number of squares  $A \in \mathcal{S}(r)$  such that  $A \ni u$  is at most  $r^2$  we get

$$||S||_{\mu} \le K \sum_{r>1} r^2 (1+r^2) r^{-\alpha} \le cK, \quad c = \sum_{r>1} (r^{-2} + r^{-4}) \le 4,$$
 (6.19)

provided that  $\alpha \geq 6$ . Thus, the Lieb-Robinson velocity is bounded by  $8C_0(1)K$ .

Let us use Lemma 6.2 to construct a local decomposition for  $\omega(O)$ . This local decomposition will involve a sequence of squares

$$B_0 = B \subset B_1 \subset B_2 \subset \ldots \subset \Lambda \tag{6.20}$$

such that  $B_j \in \mathcal{S}(q+2j)$  is the square obtained from B by adding a boundary region of thickness j on all sides of B, that is,

$$B_j = \{ u \in \Lambda : D(u, B) \le j \}.$$
 (6.21)

Note that  $B_j = \Lambda$  for large enough j. Then we can write  $\omega(O)$  as

$$\omega(O) = D_{B_0}(O) + \sum_{j \ge 1} D_{B_j}(O), \quad D_{B_0}(O) = \omega_B(O), \quad D_{B_j}(O) = \omega_{B_j}(O) - \omega_{B_{j-1}}(O). \quad (6.22)$$

Note that  $D_{B_j}(O)$  acts only on  $B_j$ , so Eq. (6.22) defines a local decomposition of  $\omega(O)$ . Using Lemma 6.2 we infer that

$$||D_{B_{j}}(O)|| \le ||\omega_{B_{j}}(O) - \omega(O)|| + ||\omega_{B_{j-1}}(O) - \omega(O)|| \le c(q+j)q^{4}K^{2} \cdot ||O|| e^{-\mu j}$$
(6.23)

for some constant c. In addition we have a standard bound (see for instance Ref. [7])

$$||D_{B_0}(O)|| = ||\omega_B(O)|| \le \frac{1}{2} ||[S_B, [S_B, O]]|| \le 2||S_B||^2 ||O||.$$
(6.24)

Taking into account that

$$||S_B|| \le \sum_{p=2}^q \sum_{\substack{A \in \mathcal{S}(p) \\ A \subseteq B}} ||S_{p,A}|| \le \sum_{p=2}^q (q-p)^2 K p^{-\alpha} e^{-\mu p} \le cq^2 K$$
(6.25)

for some constant c. Thus for all  $j \geq 0$  we have

$$||D_{B_j}(O)|| \le cq^4(q+j)K^2||O||e^{-\mu j}.$$
(6.26)

Having finished this warmup we can easily prove Lemma 6.1.

**Proof of Lemma 6.1.** Consider a local decomposition of V,

$$V = \sum_{q \ge 2} \sum_{B \in \mathcal{S}(q)} V_{q,B}, \quad ||V_{q,B}|| \le Jq^{-\beta} e^{-\mu q}.$$
(6.27)

We construct a local decomposition of  $\omega(V_{q,B})$  as in Eq. (6.22), where  $O \equiv V_{q,B}$ . We get

$$\omega(V) = \sum_{r \ge 2} \sum_{A \in \mathcal{S}(r)} \Omega_{r,A},\tag{6.28}$$

where

$$\Omega_{r,A} = \sum_{j=0}^{r/2-1} \sum_{\substack{B \in \mathcal{S}(r-2j)\\A=B:}} D_{B_j}(V_{r-2j,B}). \tag{6.29}$$

Note that for fixed j the sum over B contains a single square B such that  $A = B_j$ , see Eq. (6.21). Using Eq. (6.26) with q = r - 2j, we arrive at:

$$\|\Omega_{r,A}\| \leq \sum_{j=0}^{r/2-1} \sum_{\substack{B \in \mathcal{S}(r-2j)\\A=B_j}} K^2 \|V_{r-2j,B}\| (r-2j)^4 r e^{-\mu j}$$

$$\leq \sum_{j=0}^{r/2-1} K^2 J(r-2j)^{4-\beta} r e^{-\mu(r-j)} \leq c K^2 r J e^{-\frac{\mu r}{2}}.$$

for some constant c provided that  $\beta \geq 6$ . We have shown that  $\omega(V)$  is  $(cK^2J, \mu/2, -1)$ -decaying.

# 7 Adiabatic continuation of logical operators

## 7.1 Dressed Operators

Since the gap remains open up to a certain strength of perturbation, this means that the perturbed Hamiltonian  $H_0 + V$  is adiabatically connected to the original Hamiltonian,  $H_0$ . This adiabatic connection suggests that the perturbed Hamiltonian should have similar properties to the unperturbed Hamiltonian. For example, following [3], we can define string operators which have nontrivial action in the ground state subspace of the perturbed Hamiltonian and which have the correct commutation relations and expectation values. In fact, Theorem 1 will allow us to do even more, to define local operators that create defect excitations with well-defined energies.

To construct these operators, we use quasi-adiabatic continuation. We define a continuous family of Hamiltonians,

$$H_s = H_0 + sV, (7.1)$$

П

so that as s varies from 0 to 1,  $H_s$  continuously interpolates between  $H_0$  and the perturbed Hamiltonian.

We define a quasi-adiabatic continuation operator,  $\mathcal{D}_s$  by

$$\mathcal{D}_s \equiv i \int dt F(t) \exp(iH_s t) \left(\partial_s H_s\right) \exp(-iH_s t), \tag{7.2}$$

where the function F(t) is defined to have the following properties. First, the Fourier transform of F(t), which we denote  $\tilde{F}(\omega)$ , obeys

$$|\omega| \ge 1/2 \quad \to \quad \tilde{F}(\omega) = -1/\omega.$$
 (7.3)

Second,  $\tilde{F}(\omega)$  is infinitely differentiable, so that F(t) decays faster than any negative power of time for large |t|. Third, F(t) = -F(-t), so that  $\mathcal{D}_s$  is anti-Hermitian.

We define a unitary operator  $U_s$  by

$$U_s \equiv \mathcal{S}' \exp\left\{i \int_0^s \mathrm{d}s' \mathcal{D}_s\right\},\tag{7.4}$$

where the notation S' denotes that the above equation (7.4) is an s'-ordered exponential. The motivation for defining the above unitary operator is contained in the following lemmas:

**Lemma 7.1.** Let  $H_s$  be a differentiable family of Hamiltonians. Let  $|\Psi^i(s)\rangle$  denote eigenstates of  $H_s$  with energies  $E_i$ . Let  $E_{min}(s) < E_{max}(s)$  be continuous functions of s and

$$I(s) = \{ \lambda \in \mathbb{R} : E_{min}(s) \le \lambda \le E_{max}(s) \}.$$

Define a projector P(s) onto an eigenspace of  $H_s$  by

$$P(s) = \sum_{i: E_i \in I(s)} |\Psi^i(s)\rangle \langle \Psi^i(s)|.$$
 (7.5)

Assume that the space that P(s) projects onto is separated from the rest of the spectrum by a gap of at least 1/2 for all s with  $0 \le s \le 1$ . That is, any eigenvalue of  $H_s$  either belongs to I(s) or is separated from I(s) by a gap at least 1/2. Then, for all s with  $0 \le s \le 1$ , we have

$$P(s) = U_s P(0) U_s^{\dagger}. \tag{7.6}$$

*Proof.* By linear perturbation theory,

$$\partial_{s}P(s) = \sum_{i \in I(s)} \sum_{j \notin I(s)} \frac{1}{E_{i} - E_{j}} |\Psi^{j}(s)\rangle \Big( \langle \Psi^{j}(s) | \partial_{s}H(s) | \Psi^{i}(s)\rangle \Big) \langle \Psi^{i}(s) | + h.c$$

$$= -\sum_{i \in I(s)} \sum_{j \notin I(s)} |\Psi^{j}(s)\rangle \Big( \langle \Psi^{j}(s) | \int dt F(t) \exp(iH_{s}t) \partial_{s}H(s) \exp(-iH_{s}t) |\Psi^{i}(s)\rangle \Big) \langle \Psi^{i}(s) | + h.c$$

$$= i[\mathcal{D}_{s}, P_{s}].$$

$$(7.7)$$

The first equality in the above equation holds because

$$\langle \Psi^{j}(s)| \int dt F(t) \exp(iH_{s}t) \partial_{s} H(s) \exp(-iH_{s}t) |\Psi^{i}(s)\rangle$$

$$= \langle \Psi^{j}(s)| \int dt F(t) \exp[i(E_{j} - E_{i})t] \partial_{s} H(s) |\Psi^{i}(s)\rangle$$

$$= \tilde{F}(E_{j} - E_{i}) \langle \Psi^{j}(s)| \partial_{s} H(s) |\Psi^{i}(s)\rangle,$$
(7.8)

where we use Eq. (7.3) to show  $\tilde{F}(E_j - E_i) = -1/(E_j - E_i)$  using the assumption on the gap in the spectrum that  $|E_j - E_i| \ge 1/2$ .

Since 
$$\partial_s(U_sP(0)U_s^{\dagger}) = i[\mathcal{D}_s, U_sP(0)U_s^{\dagger}]$$
, and  $U_0 = I$ , Eq. (7.6) follows from Eq. (7.7).

This purpose for introducing this quasi-adiabatic continuation operator is to define certain "dressed" operators, following the idea introduced in [3]. Let  $O_1, O_2, ...$  be some operators that create defects when acting on the ground state of  $H_0$ . These operators may be defined to have certain commutation or anti-commutation requirements. For example, if we have certain operators  $O_i^E$  which create electric defects on a given neighboring pair of sites, and operators  $O_i^M$  which create magnetic defects on a given neighboring pair of plaquettes in a toric code

state, then these operators all commute with each other, except an electric and a magnetic operator anti-commute if the bond connecting the sites and the bond connecting the plaquettes on the dual lattice intersect. Then we define "dressed" operators  $O_i(s)$  by

$$O_i(s) = U_s O_i U_s^{\dagger}. \tag{7.9}$$

Since  $U_s$  is unitary, these operators  $O_i(s)$  obeys the same commutation or anti-commutation requirements:

$$[O_i, O_j] = 0 \to [O_i(s), O_j(s)] = 0,$$
  

$$\{O_i, O_j\} = 0 \to \{O_i(s), O_j(s)\} = 0.$$
(7.10)

Let  $P_n(s)$  project onto the eigenspace of H(s) with energy in the interval  $I_n$ , see Theorem 1. We assume that J is chosen sufficiently small (see Section 1.1) so that  $I_n$  is separated from the rest of the spectrum by a gap at least 1/2. So, Lemma 7.1 implies that if a product of operators  $O_1O_2...O_m$  acting any ground state of  $H_0$  creates an eigenstate of  $H_0$  with energy n, then the product of operators  $O_1(s)O_2(s)...O_m(s)$  acting on any ground state state of  $H_s$  creates a state with energy in the interval  $I_n$ . To see this, note that a ground state  $\Psi_0(s)$  of  $H_s$  can be written as  $U_s\Psi_0(0)$  for some ground state  $\Psi_0(0)$  of  $H_0$ . Then,

$$P_n(s)O_1(s)...O_m(s)|\Psi_0(s)\rangle = U_sP_n(0)O_1(0)...O_m(0)|\Psi_0\rangle$$

$$= U_s(0)O_1(0)...O_m(0)|\Psi_0\rangle = O_1(s)...O_m(s)|\Psi_0(s)\rangle.$$
(7.11)

Further, if some product of operators  $O_1O_2...$  has a given expectation value in the ground state of  $H_0$ , then the corresponding dressed operators have the same expectation value in the ground state of  $H_s$ . For example, since a product of  $\sigma^x$  around a contractable loop has expectation value unity in the ground state of the toric code, the same product of dressed operators will have the same expectation value in the ground state of  $H_s$ .

The next important property we want to show is that the operators  $O_i(s)$  are local. To do these, we need a Lieb-Robinson bound for quasi-adiabatic continuation. Before describing this, some comments. We have chosen to define the quasi-adiabatic continuation using "exact" expressions. That is, we have chosen a filter function F(t) such that its Fourier transform is exactly equal to  $1/\omega$  outside a certain interval. This is Osborne's [21] modification of the quasiadiabatic continuation of [3]. As a result, our filter function F(t) decays faster than any negative power of t, but does not decay exponentially in t (however, there do exist such filter functions with |F(t)| decaying exponentially in a polynomial of t [23]). This contrasts with the approach in [3], where an approximation was used that gave a filter function decaying exponentially in  $t^2$ . Each approach has certain advantages and disadvantages. The approach in [3] has the advantage that one can often get exponentially good approximations, rather than approximations which merely decay faster than any power. However, Osborne's modification which we use here has a few advantages also. First, we get an exact result that  $\Psi_0(s) = U_s \Psi_0(0)$  above, while the approach in [3] only gives approximate estimates. Second, it is much easier to derive locality estimates. In particular, the Lieb-Robinson bound for quasi-adiabatic continuation in the next lemma, which is the key step to prove locality of the dressed operators, is much easier than [3, 22]. The reason that this bound is so much simpler is the following: the Fourier transform

of the filter function here can be chosen to be some given bounded function of  $\omega$ . In [3, 22], we have a parameter-dependent family of filter functions. As this parameter is increased, the accuracy of the approximation improves, but also the upper bound on the Fourier transform of the filter function gets worse, and hence the bound on the norm of  $\mathcal{D}$  worsens.

The locality of the dressed operators depends on locality properties of  $H_0, V$ . We require that  $H_0, V$  are both sums of of local operators  $H_{0,Z}, V_Z$ , where both  $H_{0,Z}, V_Z$  obey a bound that, for all sites  $u \in \Lambda$ ,

$$\sum_{Z\ni u} ||H_{0,Z}|||Z| \exp(\mu \operatorname{diam}(Z)) = O(1),$$

$$\sum_{Z\ni u} ||V_Z|||Z| \exp(\mu \operatorname{diam}(Z)) \le J < \infty,$$
(7.12)

where |Z| denotes the cardinality of Z, for some positive constants  $\mu, J$ .

Given such locality property, for any finite dimensional system we have a Lieb-Robinson bound for  $H_0 + sV$  that, for any operator  $O_A$  supported on set A and any  $O_B$  supported on set B,

$$\|[\exp(iH_s t)O_A \exp(-iH_s t), O_B]\| \le \exp[-\mu(\operatorname{dist}(A, B) - v_{LR} t)]|A|\|O_A\|\|O_B\|, \tag{7.13}$$

where  $v_{LR}$  is some constant which depends on  $\mu, J$ . This bound is shown in [20].

**Lemma 7.2.** Let  $H_0$  and V obey Eq. (7.12). Then, if  $O_A$  is supported on set A and  $O_B$  is supported on set B,

$$||[U_s O_A U_s^{\dagger}, O_B]|| \le h(\operatorname{dist}(A, B))|A|||O_A||||O_B||,$$
 (7.14)

where |A| denotes the cardinality of A, for  $0 \le s \le 1$ , for some function h(l) which decays faster than any negative power of l. Similarly,

$$||[U_s^{\dagger}O_AU_s, O_B]|| \le h(\operatorname{dist}(A, B))|A|||O_A||||O_B||,$$
 (7.15)

Further, the operator  $\mathcal{D}_s$  is a sum of operators  $\mathcal{D}_s(Z)$ , with

$$\|\mathcal{D}_s(Z)\| \le \text{const.} \times \|V_Z\|,\tag{7.16}$$

and with the property that each such operator  $\mathcal{D}_s(Z)$  obeys, for any operator  $O_B$ ,

$$\|[\mathcal{D}_s(Z), O_B]\| \le h'(\operatorname{dist}(A, B))|Z|\|V_Z\|\|O_B\|,$$
 (7.17)

for some function h'(l) which decays faster than any negative power of l.

*Proof.* We have

$$\mathcal{D}_s = \sum_Z \mathcal{D}_s(Z),\tag{7.18}$$

where

$$\mathcal{D}_s(Z) = i \int dt F(t) \exp(iH_s t) V_Z \exp(-iH_s t). \tag{7.19}$$

Eq. (7.16) follows immediately from a triangle inequality because  $\int dt |F(t)|$  converges. Let  $O_B$  be an operator supported on set B. Then

$$\|[\mathcal{D}_s(Z), O_B]\| \le \int dt |F(t)| \|[\exp(iH_s t)V_Z \exp(-iH_s t), O_B]\|.$$
 (7.20)

For  $t \leq \operatorname{dist}(Z, B)/2v_{LR}$ , the above expression is exponentially small in  $\operatorname{dist}(Z, B)$  by the Lieb-Robinson bound, while for larger t, the expression is bounded by |F(t)|, and hence this commutator decays faster thany any negative power of  $\operatorname{dist}(Z, B)$ .

This decomposition and locality bound (7.20) implies that there is a Lieb-Robinson bound for evolution using  $\mathcal{D}_s$  as a Hamiltonian. Using the Lieb-Robinson bound in [20] (the use of this bound does require some geometric properties of the lattice, which hold for any finite dimensional lattice), we have that

$$||[U_s O_A U_s^{\dagger}, O_B]|| \le \exp(cs) h'(\operatorname{dist}(A, B)) |A| ||O_A|| ||O_B||,$$
 (7.21)

for some constant c which depends on  $F, J, \mu$  and on the geometric properties of the lattice, and for some function h'(l) which decays faster than any negative power of l. Since we assume  $s \leq 1$ , Eq. (7.13) follows from Eq. (7.21).

Eq. (7.13) implies that  $O_A(s) = U_s O_A U_s^{\dagger}$  can be approximated by an operator  $O_l$  localized on the set of sites within distance l of set A. To see this, define  $O_l = \int dU U O_A(s) U^{\dagger}$ , where the integral ranges over unitary rotations, with Haar measure, supported on sites with distance greater than l from set A. Then, following [2], the desired result follows.

The string operators can be used to manipulate the ground states. We will assume that the operators  $O_1$  which create defects are, in fact, unitaries. For example, if a certain product of electric operators,  $O_1O_2$ ... creates a nontrivial loop around the torus and has nontrivial action on the ground state of  $H_0$ , then the product of dressed operators  $O_1(s)O_2(s)$ ... has the same action on the ground state of  $H_s$ . However, one might wonder how to create these dressed operators; it would be preferable not to have to solve quasi-adiabatic evolution equations to calculate what the operators should be, and then to produce them by careful control of a timedependent Hamiltonian. Fortunately, we can use the gaps in the excited state spectrum to argue that it is possible to drag defects in the perturbed Hamiltonian  $H_s$  in an identical way to what is done in  $H_0$ . First, we apply some local operator in an attempt to create a defect pair on neighboring sites. If this operator is not exactly equal to the desired dressed operator, then this attempt may fail, in the sense that the energy may not fall into the desired range  $I_2$ , see Theorem 1. We can detect this failure by measuring the energy, and cooling back to the ground state: that is, if there are extra defects created, we will drag them together, using the procedure described in the next paragraph, to annihilate them. Since the dressed operator is local, then some local operator (for example, the undressed operator) is expected to have non-vanishing matrix elements between the ground state and the desired state. Eventually we will succeed in creating the defect pair.

We now try to move the defect pair by weak changes in the Hamiltonian. We add a perturbation U to the Hamiltonian with norm bounded by 1/4. Such a weak perturbation will keep the gap open to the states with 3 defects (of course, there are no such three defect states in the

toric code, but in more general models there are). So, by adding this perturbation and changing it adiabatically, we remain in the subspace with a defect pair.

In a toric code Hamiltonian, every state in this subspace is a linear combination of states which are given by acting on a ground state  $\Psi_0(s)$  by a linear combination of products of strings of dressed operators. Consider a given state created by a single string of dressed operators, with endpoints i and j which are far separated. Call this state  $\Psi(S)$ , where S is a string with endpoints at i, j. Note that using the property that any contractable loop of dressed operators has expectation value unity in the ground state of  $H_s$ , we can show that we can deform these strings in any way that leaves the endpoints fixed while leaving the state unchanged. Let us now see how the expectation value of local operators can change in the state  $\Psi(S)$ . If O is an operator which is far separated from i, j, we claim that the expectation value of O in state  $\Psi(S)$  is close to its value in the ground state. To see these, note that that we can define a state  $\Psi(S') = \Psi(S)$  where S' is a string which stays far from O. So, without loss of generality, we can assume that string S is far from O. Then, the using the locality properties of the dressed operators, they almost commute with O, and so

$$\langle \Psi(S)|O|\Psi(S)\rangle = \langle \Psi_0(s)|O_n^{\dagger}...O_1^{\dagger}OO_1...O_n|\Psi_0(s)\rangle$$

$$\approx \langle \Psi_0(s)|O_n^{\dagger}...O_1^{\dagger}O_1...O_nO|\Psi_0(s)\rangle$$

$$= \langle \Psi_0(s)|O|\Psi_0(s)\rangle.$$
(7.22)

So, the perturbation U can effect the state, but only if U acts close to one of the endpoints of the string. Further, if U is local, then one may show that U will have small matrix elements except between  $\Psi(S)$  and  $\Psi(S')$  for strings S' which differ from S only by small motion of one of the endpoints. Thus, if U is chosen to have the property that it reduces the energy when an end of the string is close to a certain point, the operator U can indeed by used to drag the defects.

So, we have established that it is possible to create a defect pair and drag it. The gap to the rest of the spectrum prevents additional defects from being created. If they are created, we can detect them by measuring the energy, and we can drag them to destroy them and correct errors.

We can drag one of the defects around the system. Since at every step we remain in a state which is a linear combination of states  $\Psi(S)$ , and eventually we succeed in dragging a defect all the way around the system, if we are able to return to a ground state  $\Psi'_0(s)$ , then  $\Psi'_0(s)$  is produced by a local operation O on a state  $\Psi(S)$  where S is a string with two nearby endpoints connected by a string that goes around the sample. Any such state  $\Psi(S)$  is a product of dressed operators acting on  $\Psi_0(s)$ . So, it is equal to  $\Psi(S) = U_s O_1 ... O_n \Psi_0(0)$ . Suppose  $\Psi'_0(s)$  is a ground state such that  $\langle \Psi'_0(s), O\Psi(S) \rangle$  is non-negligible. This expectation value is equal to

$$\langle \Psi_0'(0)| \left( U_s^{\dagger} O U_s \right) O_1 \dots O_n |\Psi_0(0)\rangle. \tag{7.23}$$

However, due to the locality properties of quasi-adiabatic continuation, the operator  $U_s^{\dagger}OU_s$  is approximately local. Hence, the ground state  $\Psi_0'(0)$  of the unperturbed system is connected by a local operator to the state  $O_1...O_n|\Psi_0(0)\rangle$ . So, the state  $\Psi_0'(0)$  must be close to the state

formed by acting on  $\Psi_0(0)$  with a noncontractible string that winds around the sample. Thus, by dragging the defects in the perturbed system, we succeed in effecting the desired transformation on the ground state sector of the perturbed system.

One thing that can go wrong in this procedure is that we might inadvertently create the wrong type of defect. We might accidentally create a magnetic defect pair rather than an electric defect pair (one may show that the perturbation U will have small amplitude to change electric into magnetic defects if the defects are far separated and U is local). However, we expect to be able to locally tell the difference between these types of defects and thus determine what kind of defect has been created, and, if the wrong type has been created, to destroy the defect pair by bringing them together.

One would like to improve the arguments here to a full formal proof that using classical control it is possible to correct local errors and perform controlled operations in this perturbed toric code. We will leave a complete proof of this to the future, but we have established the existence of operators with the necessary properties and of the gaps in the spectrum.

# Acknowledgments

We thank Barbara Terhal and David DiVincenzo for useful discussions. Part of this work was done while SB and SM were visiting the Erwin Schrödinger International Institute for Mathematical Physics at Vienna. SB was partially supported by the DARPA QUEST program under contract number HR0011-09-C-0047. SM thanks the organizers of the program on "Quantum Information Science" at the KITP at UC Santa Barbara, where part of this work was completed. SM was supported by NSF Grant DMS-07-57581 and DOE Contract DE-AC52-06NA25396.

## References

- [1] A. Kitaev, "Fault-tolerant quantum computation by anyons", Ann. Phys. 303, 2 (2003).
- [2] S. Bravyi, M. B. Hastings, and F. Verstraete, "Lieb-Robinson Bounds and the Generation of Correlations and Topological Quantum Order", Phys. Rev. Lett. 97, 050401 (2006).
- [3] M. B. Hastings and Xiao-Gang Wen, "Quasi-adiabatic Continuation of Quantum States: The Stability of Topological Ground State Degeneracy and Emergent Gauge Invariance", Phys. Rev. **B72**, 045141 (2005).
- [4] M. A. Levin and X.-G. Wen, "String-net condensation: A physical mechanism for topological phases", Phys.Rev. B 71, 045110 (2005).
- [5] L.M. Duan, E. Demler, M.D. Lukin, "Controlling Spin Exchange Interactions of Ultracold Atoms in Optical Lattices", Phys. Rev. Lett. 91, 090402 (2003)
- [6] R. Koenig, "Simplifying quantum double Hamitonians using perturbative gadgets", e-print arXiv:0901.1333.
- [7] S. Bravyi, D. DiVincenzo, D. Loss, and B. Terhal, "Simulation of Many-Body Hamiltonians using Perturbation Theory with Bounded-Strength Interactions", Phys.Rev.Lett. 101, 070503 (2008).

- [8] A. Kitaev, "Anyons in an exactly solved model and beyond", Annals of Physics **321**, 2–111 (2006).
- [9] F. D. M. Haldane, "Continuum dynamics of the 1-d Heisenberg antifer- romagnet: identification with the O(3) nonlinear sigma models", Phys. Lett. A93, pp. 464-468 (1983);
  F. D. M. Haldane, "Nonlinear field theory of large-spin Heisenberg anti-ferromagnets",
  Phys. Rev. Lett. 50, 1153-1156 (1983).
- [10] D. A. Yarotsky, "Ground states in relatively bounded quantum perturbations of classical lattice systems", Commun. Math. Phys. **261**, 799-819 (2006).
- [11] I. Affleck, T. Kennedy, E. H. Lieb, and H. Tasaki, "Valence bond ground states in isotropic quantum antiferromagnets", Commun. Math. Phys. 115, 477-528 (1987).
- [12] S. Trebst, P. Werner, M. Troyer, K. Shtengel, and C. Nayak, "Breakdown of a topological phase: Quantum phase transition in a loop gas model with tension", Phys. Rev. Lett. 98, 070602 (2007).
- [13] I. Klich, "On the stability of topological phases on a lattice", e-print arXiv:0912.0945.
- [14] S. Głazek and K. Wilson, "Renormalization of Hamiltonians, Phys. Rev. **D48**, p. 5863 (1993); F. Wegner, "Flow equations for Hamiltonians", Ann. Phys. **3**, p. 77 (1994).
- [15] S. Bravyi, D. Poulin, and B. Terhal, "Tradeoffs for reliable quantum information storage in 2D systems", e-print arXiv:0909.5200.
- [16] H. Bombin and M. Martin-Delgado, "Topological Quantum Distillation", Phys. Rev. Lett. 97, 180501 (2006).
- [17] T. Kato, "Perturbation theory for linear operators", Springer-Verlag New York (1966).
- [18] R. Bhatia, "Matrix Analysis", Springer-Verlag New York (1997)
- [19] B. Nachtergaele and R. Sims, "Locality estimates for quantum spin systems", e-print arXiv:0712.3318
- [20] M. B. Hastings and T. Koma, Commun. Math. Phys. **265**, 781 (2006).
- [21] T. J. Osborne, "Simulating adiabatic evolution of gapped spin systems", Phys. Rev. A, 75, 032321.
- [22] M. B. Hastings and S. Michalakis, "Quantization of Hall Conductance for Interacting Electrons Without Averaging Assumptions", 0911.4706.
- [23] M.B. Hastings, in preparation.