

# NORMAL FLUCTUATION IN QUANTUM ERGODICITY FOR WIGNER MATRICES

BY GIORGIO CIPOLLONI<sup>1,a</sup>, LÁSZLÓ ERDŐS<sup>2,b</sup> AND DOMINIK SCHRÖDER<sup>3,c</sup>

<sup>1</sup>*Department of Physics, Princeton University, <sup>a</sup>gc4233@princeton.edu*

<sup>2</sup>*IST Austria, <sup>b</sup>lerdos@ist.ac.at*

<sup>3</sup>*ETH Institute for Theoretical Studies, ETH Zurich, <sup>c</sup>dschoreder@ethz.ch*

We consider the quadratic form of a general high-rank deterministic matrix on the eigenvectors of an  $N \times N$  Wigner matrix and prove that it has Gaussian fluctuation for each bulk eigenvector in the large  $N$  limit. The proof is a combination of the energy method for the Dyson Brownian motion inspired by Marcinek and Yau (2021) and our recent multiresolvent local laws (*Comm. Math. Phys.* **388** (2021) 1005–1048).

**1. Introduction.** Quantum unique ergodicity (QUE) in a disordered or chaotic quantum system asserts that the eigenvectors of the Hamilton operator tend to become uniformly distributed in the phase space; see [2, 16, 27, 29, 30, 33] for the seminal results and [15] for more recent references. We study a particularly strong form of this phenomenon for Wigner random matrices, the simplest prototype of a fully chaotic Hamiltonian. These are  $N \times N$  random Hermitian matrices  $W = W^*$  with centred, independent, identically distributed (*i.i.d.*) entries up to the symmetry constraint  $w_{ab} = \overline{w_{ba}}$ . Let  $\{\mathbf{u}_i\}_{i=1}^N$  be an orthonormal eigenbasis of  $W$  corresponding to the eigenvalues  $\boldsymbol{\lambda} = (\lambda_i)_{i=1}^N$  listed in increasing order. Recently, we showed [15] that for any deterministic matrix  $A$  with  $\|A\| \leq 1$ , the eigenvector overlaps  $\langle \mathbf{u}_i, A\mathbf{u}_i \rangle$  converge to  $\langle A \rangle := \frac{1}{N} \text{Tr} A$ , the normalized trace of  $A$ , in the large  $N$  limit. More generally, we proved that

$$(1) \quad \max_{i,j} |\langle \mathbf{u}_i, A\mathbf{u}_j \rangle - \langle A \rangle \delta_{ij}| \lesssim \frac{N^\epsilon}{\sqrt{N}}$$

holds with very high probability. We note that the bound (1) is optimal for high-rank deterministic matrices  $A$  and is coined as the *Eigenstate Thermalization Hypothesis* by Deutsch [18] and Srednicki [31]; see also [17], equation (20).

The main result of the current paper, Theorem 2.2, asserts that  $\langle \mathbf{u}_i, A\mathbf{u}_i \rangle$  has a Gaussian fluctuation on scale  $N^{-1/2}$ , more precisely,

$$(2) \quad \sqrt{N}[\langle \mathbf{u}_i, A\mathbf{u}_i \rangle - \langle A \rangle]$$

converges to a normal distribution for any Hermitian observables  $A = A^*$  of high rank and for any eigenvectors  $\mathbf{u}_i$  whose eigenvalue belongs to the bulk of the spectrum.

For Gaussian ensembles and  $A$  being a projection onto macroscopically many coordinates, (2) can be proven by using the special invariance property of the eigenvectors (see [28], Theorem 2.4). Our result concerns general Wigner matrices, and it has two main features: it concerns individual eigenvectors, *and* it is valid for general high rank observables. We now explain related previous results which all addressed only one of these features. First, Gaussianity of (2), after a small averaging in the index  $i$ , has recently been established in [13],

---

Received March 2021; revised September 2021.

*MSC2020 subject classifications.* Primary 60B20; secondary 15B52.

*Key words and phrases.* Local law, Dyson Brownian motion, stochastic eigenstate equation, eigenvector moment flow.

Theorem 2.3, using resolvent methods. Second, fluctuations involving individual eigenvectors in the bulk spectrum for general Wigner matrices can only be accessed by the Dyson Brownian motion approach which has only been developed for finite rank observables [3, 11, 26]. We now explain the background of these concepts.

1.1. *Dyson Brownian motion for eigenvectors.* For the simplest rank one case,  $A = |\mathbf{q}\rangle\langle\mathbf{q}|$  with some deterministic unit vector  $\mathbf{q}$ , Bourgade and Yau [11] showed that the squared normalised overlaps  $N|\langle\mathbf{u}_i, \mathbf{q}\rangle|^2$  converge in distribution to the square of a standard Gaussian variable as  $N \rightarrow \infty$  (see also [22, 32] for the same result without DBM but under four moment matching condition in the bulk). Similar results have been obtained for deformed Wigner matrices [3], for sparse matrices [10], and for Lévy matrices [1]. Note that both the scaling and the limit distribution for the rank one case are different from (2). The basic intuition is that the coordinates of  $\mathbf{u}_i$  are roughly independent; thus the sum in  $\langle\mathbf{u}_i, \mathbf{q}\rangle = \sum_a \mathbf{u}_i(a)\mathbf{q}(a)$  obeys a central limit theorem (CLT) on scale  $N^{-1/2}$ . In fact, [11] also considers the joint distribution of finitely many eigenvectors tested against one fixed vector  $\mathbf{q}$  and the joint distribution of a single eigenvector with finitely many test vectors  $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_K$  for any fixed  $K$ , independent of  $N$ . Very recently, Marcinek and Yau [26] have established that the overlaps of finitely many eigenvectors *and* finitely many orthogonal test vectors are also asymptotically independent (squared) normal. Their method is very general and also applies to a large class of other random matrix ensembles, such as sparse or Lévy matrices.

The fundamental method behind all results involving individual eigenvectors for general Wigner matrices is the Dyson Brownian motion (DBM) for eigenvectors, also called the *stochastic eigenstate equation*, generated by a simple matrix Brownian motion for  $W$  introduced by Bourgade and Yau in [11]. We briefly summarize the key steps in [11] in order to highlight the new ideas we needed to prove the Gaussianity of (2).

For each fixed  $n$ , the evolution of the joint  $n$ th order moments of the overlaps  $N|\langle\mathbf{u}_i, \mathbf{q}\rangle|^2$  for different  $i$ 's and fixed  $\mathbf{q}$  is described by a system of parabolic evolution equations, called the *eigenvector moment flow*. Interpreting each such overlap as a particle sitting at location  $i$  in the discrete one dimensional index space  $[N] = \{1, 2, \dots, N\}$ , the moment flow naturally corresponds to a Markovian jump process of  $n$  particles. It turns out that the rate of a jump from site  $i$  to site  $j$  is proportional with  $N^{-1}(\lambda_i - \lambda_j)^{-2}$ . Different  $\mathbf{q}$ 's can be incorporated by appropriately assigning *colours* to the particles. By the fast local equilibration property of the DBM, the moments of  $N|\langle\mathbf{u}_i, \mathbf{q}\rangle|^2$  quickly become essentially independent of the index  $i$ , at least, for indices corresponding to nearby eigenvalues  $\lambda_i$ ; hence, they can be

computed by locally averaging over  $i$ . For example, in the simplest  $n = 1$  case, we have

$$(3) \quad f_i := \mathbf{E}[N|\langle\mathbf{u}_i, \mathbf{q}\rangle|^2|\boldsymbol{\lambda}] \approx f_{i'} = \mathbf{E}[N|\langle\mathbf{u}_{i'}, \mathbf{q}\rangle|^2|\boldsymbol{\lambda}], \quad |i - i'| \ll N,$$

already after a very short time  $t \gg |i - i'|/N$ . Here, we consider the conditional expectation of the eigenvectors, given that the eigenvalues are fixed. Since the global equilibrium of the DBM is the constant function  $f_i = 1$ , equilibration directly implies *smoothing* or *regularisation* in the dependence on the indices  $i$ .

On the other hand, by spectral theorem

$$(4) \quad \langle\mathbf{q}, \Im G(\lambda_i + i\eta)\mathbf{q}\rangle = \frac{1}{N} \sum_{i'=1}^N \frac{\eta}{(\lambda_i - \lambda_{i'})^2 + \eta^2} N|\langle\mathbf{u}_{i'}, \mathbf{q}\rangle|^2,$$

where  $G = G(z) = (W - z)^{-1}$  is the resolvent at a spectral parameter  $z \in \mathbf{C} \setminus \mathbf{R}$ . Using that the eigenvalues  $\lambda_{i'}$  are *rigid*, that is, they are very close the corresponding quantiles  $\gamma_{i'}$  of the

Wigner semicircle density (see (21) later), the  $i'$ -summation in (4) is a regularised averaging over indices  $|i' - i| \lesssim N\eta$ . Performing the  $i'$  summation in (4) by using (3), we obtain

$$\mathbf{E}[N|\langle \mathbf{u}_i, \mathbf{q} \rangle|^2 | \boldsymbol{\lambda}] \approx \frac{1}{\Im m_{sc}(\gamma_i)} \mathbf{E}[\langle \mathbf{q}, \Im G(\gamma_i + i\eta)\mathbf{q} \rangle | \boldsymbol{\lambda}]$$

for times  $t \gg \eta$  where  $m_{sc}$  is the Stieltjes transform of the Wigner semicircle law. Choosing  $\eta$  slightly above the local eigenvalue spacing,  $\eta = N^{-1+\epsilon}$  in the bulk of the spectrum, we have  $\langle \mathbf{q}, \Im G(\gamma_i + i\eta)\mathbf{q} \rangle \approx \Im m_{sc}(\gamma_i)$  not only in expectation but even in high probability by the *isotropic local law* for Wigner matrices [23]. Combining these inputs, we obtain  $\mathbf{E}N|\langle \mathbf{u}_i, \mathbf{q} \rangle|^2 \approx 1$  along the DBM after a short time  $t \gg N^{-1+\epsilon}$ . A similar argument holds for higher moments. Finally, the small Gaussian component added by the DBM can be removed by standard perturbation methods, by the so called *Green function comparison* theorems.

1.2. *Dyson Brownian motion for general overlaps.* Given the method to handle  $N|\langle \mathbf{u}_i, \mathbf{q} \rangle|^2$  described above, the Gaussianity of overlaps  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_i \rangle$  with a general high rank matrix  $A$  can be approached in two natural ways. We now explain both of them to justify our choice. The first approach is to write  $A = \sum_{k=1}^N a_k |\mathbf{q}_k\rangle \langle \mathbf{q}_k|$  in spectral decomposition with  $|a_k| \lesssim 1$  and an orthonormal set  $\{\mathbf{q}_k\}_{k=1}^N$  to have

$$(5) \quad \langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_i \rangle = \sum_{k=1}^N a_k |\langle \mathbf{u}_i, \mathbf{q}_k \rangle|^2.$$

If all overlaps  $|\langle \mathbf{u}_i, \mathbf{q}_k \rangle|^2$ ,  $k = 1, 2, \dots, N$  were independent, then the central limit theorem applied to the summation in (5) would prove the normality of  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_i \rangle$ . This requires that the number of nonzero summands in (5), the rank of  $A$ , also goes to infinity as  $N$  increases. Hence, via the spectral decomposition of  $A$ , the Gaussianity of  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_i \rangle$  appears rather an effect of the approximate independence of the overlaps  $|\langle \mathbf{u}_i, \mathbf{q}_k \rangle|^2$  for different  $k$ 's than their actual limit distribution. The analysis of the eigenvector moment flow [11, 26] yields this independence for finitely many  $k$ 's, but it is not well suited for tracking overlaps  $|\langle \mathbf{u}_i, \mathbf{q}_k \rangle|^2$  with a very large number of  $\mathbf{q}_k$  vectors simultaneously. Hence, we discarded this approach.

The second natural approach is to generalise the eigenvector moment flow to moments of  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_i \rangle$ ; this has been first achieved in [12]. Such flow naturally involves off-diagonal overlaps  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_j \rangle$  as well. Therefore, we need to describe conditional moments of the form  $\mathbf{E}[\prod_{r=1}^n \langle \mathbf{u}_{i_r}, \mathbf{A}\mathbf{u}_{j_r} \rangle | \boldsymbol{\lambda}]$  with different collections of *index pairs*  $(i_1, j_1), (i_2, j_2), \dots, (i_n, j_n)$  with the constraint that every index appears even number of times. Thus, the relevant moments can naturally be represented in an  $n$ -dimensional subset  $\Lambda^n$  of  $[N]^{2n}$  (see Section 4.1). Moreover, in [12], equation (2.15), a certain symmetrised linear combination of  $n$ -order moments, the *perfect matching observable* was found that satisfies a closed equation along the Dyson Brownian motion; see (20) and (23). Moments of diagonal overlaps  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_i \rangle$  can then be recovered from the perfect matching observable by setting all indices equal. Off-diagonal overlaps  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_j \rangle$ , in general, cannot be recovered (except in the  $n = 2$  case using an additional antisymmetric (“fermionic”) version of the perfect matching observable [4]).

The main obstacle along this second approach is the lack of the analogue of (4) for general overlaps  $\langle \mathbf{u}_i, \mathbf{A}\mathbf{u}_j \rangle$ . Consider the  $n = 2$  case. A (regularized) local averaging in *one* index yields

$$(6) \quad \frac{1}{N} \sum_{i'=1}^N \frac{\eta}{(\lambda_i - \lambda_{i'})^2 + \eta^2} N |\langle \mathbf{u}_{i'}, \mathbf{A}\mathbf{u}_j \rangle|^2 = \langle \mathbf{u}_j, A \Im G(\gamma_i + i\eta) \mathbf{A}\mathbf{u}_j \rangle,$$

which still involves an eigenvector  $\mathbf{u}_j$ , hence is not accessible solely by resolvent methods. Note that, for  $A = |\mathbf{q}\rangle \langle \mathbf{q}|$ , the overlap  $\langle \mathbf{u}_{i'}, \mathbf{A}\mathbf{u}_j \rangle$  factorizes, and the averaging in  $i'$  can be

done independently of  $j$ . For general  $A$  we can handle an averaging in *both* indices, that is, we will use that

$$(7) \quad \frac{1}{N^2} \sum_{i',j'=1}^N \frac{\eta}{(\lambda_i - \lambda_{i'})^2 + \eta^2} \frac{\eta}{(\lambda_j - \lambda_{j'})^2 + \eta^2} N |\langle \mathbf{u}_{i'}, A \mathbf{u}_{j'} \rangle|^2 = \langle A \mathfrak{S}G(\gamma_i + i\eta) A \mathfrak{S}G(\gamma_j + i\eta) \rangle.$$

The normalised trace in the right-hand side is accessible by resolvent methods using the recent multi- $G$  local law proven in [15], Proposition 3.4. However, the generator of the eigenvector moment flow (24) involves the *sum* of averaging operators, as in (6), in all coordinate directions and not their *product*, as needed in (7). Higher moments ( $n > 2$ ) require averaging in more than two indices simultaneously that is not apparently available in the generator. To remedy this situation, we now review how the equilibration (smoothing) property of the parabolic equation for the perfect matching observable can be manifested.

1.3. *Local smoothing of the eigenvector moment flow: An overview.* The technically simplest way to exploit the smoothing effect is via the maximum principle introduced in [11]. However, this requires that the generator is negative and itself has the necessary local averaging property to obtain a quantity computable by a local law; this is the case for eigenfunction overlaps, as in (4) but not for general overlaps  $\langle \mathbf{u}_i, A \mathbf{u}_j \rangle$  in (6). We remark that the maximum principle was also used in [12] for more general overlaps, but only for getting an a priori bound and not for establishing their distribution. For this cruder purpose a rougher bound on (6) was sufficient that could be iteratively improved, but always by an  $N^\epsilon$  factor off the optimal value.

A technically much more demanding way to exploit the equilibration of the eigenvector moment flow would be via homogenisation theory. In random matrix theory, homogenisation was originally introduced for the Dyson eigenvalue flow in [9, 24] by noticing that the generator is a discrete approximation of the one dimensional fractional Laplacian operator  $|p| = \sqrt{-\Delta}$  with translation invariant kernel  $(x - y)^{-2}$  whose heat kernel is explicitly known. Unfortunately, the eigenvector flow is more complicated, and a good approximation with a well-behaving continuous heat kernel is missing although homogenisation might also be accessible via a sequence of maximum principles as in [8].

Finally, the last and most flexible method for equilibration are the ultracontractivity estimates on the heat kernel that can be obtained by the standard Nash method from Poincaré or Sobolev inequalities for the Dirichlet form determined by the generator. In random matrix theory these ideas have been introduced in [20] for the eigenvalue gap statistics and have later been used as a priori bounds for the homogenisation theory. However, in the bulk regime they are barely not sufficiently strong to get the necessary precision for individual eigenvalues; they had to be complemented either by De Giorgi–Nash–Moser Hölder regularity estimates [20] or homogenisation [9, 24].

The recent work by Marcinek and Yau [26] remedies this shortcoming of the ultracontractivity bound by combining it with an energy method. The main motivation of [26] was to consider the joint distribution of the overlaps  $|\langle \mathbf{u}_i, \mathbf{q}_k \rangle|^2$  for several eigenvectors and several test vectors simultaneously. The generator of the resulting *coloured eigenvector moment flow* lacks the positivity preserving property rendering the simple argument via maximum principle impossible. It turns out that this lack of positivity is due to a new *exchange term* in the generator that is present only because several  $\mathbf{q}_k$ 's (distinguished by colours) are considered simultaneously. However, the generator with the problematic exchange term is still positive in  $L^2$ -sense, and its Dirichlet form satisfies the usual Poincaré inequality from which ultracontractivity bounds can still be derived. The additional smallness now comes from an effective decay of the  $L^2$ -norm of the solution where local averaging like (4) can be exploited.

1.4. *Main ideas of the current paper.* The proof of Gaussianity of (2) consists of three steps:

Step 1. We use the energy method inspired by [26] together with the recent two- $G$  local law from [15], Proposition 3.4, and more general multi- $G$  local laws, proven in Section 5, to exploit an effective averaging mechanism to reduce the  $L^2$ -norm of the solution. In particular, to understand (7) we need a two- $G$  local law instead of the single- $G$  isotropic law used in (4).

Step 2. We use an  $L^2 \rightarrow L^\infty$  ultracontractivity bound of the colourblind eigenvector moment flow from [26], Proposition 6.29.

Step 3. The first two steps prove the Gaussianity of the overlap (2) for Wigner matrices with a tiny Gaussian component. With a standard Green function comparison argument combined with the a priori bound (1) proven in [15] we remove this Gaussian component.

Step 2 and Step 3 are standard adaptations of existing previous results, so we focus only on explaining Step 1. We use the energy method in a very different way and for a very different purpose than [26] but for the same reason, its robustness. In the standard energy argument, if  $f_t$  satisfies the parabolic evolution equation  $\partial_t f_t = \mathcal{L}_t f_t$  with a (time-dependent) generator  $\mathcal{L}_t$ , then

$$\frac{1}{2} \partial_t \|f_t\|_2^2 = \langle f_t, \mathcal{L}_t f_t \rangle =: -D_t(f_t) \leq 0,$$

where  $D_t$  is the Dirichlet form (energy) associated to the generator  $\mathcal{L}_t$ . The goal is to give a good lower bound,

$$(8) \quad D_t(f) \geq c \|f\|_2^2 - \text{error},$$

and use a Gronwall argument to conclude an effective  $L^2$ -decay along the dynamics. However, at this moment the Dirichlet form may first be replaced by a smaller one,  $\tilde{D}_t(f) \lesssim D_t(f)$ , for which an effective lower bound (8) is easier to obtain. In our case, the gain comes from estimating the error term in (8) by exploiting the local averaging in all directions, as in (7), so that we could use the multi- $G$  local law.

How to find  $\tilde{D}$ ? Very heuristically, the generator of the eigenvector moment flow is a discrete analogue of  $|p_1| + |p_2| + \dots + |p_n|$ , that is, the sum of  $|p|$ -operators along all the  $n$  coordinate directions in the  $n$ -dimensional space  $\Lambda^n$ . However, the necessary averaging in (7) is rather the product of these one dimensional operators. Normally, sums of first order differential operators cannot be compared with their product since they scale differently with the length. But our operators have a short-range regularization on the scale  $\eta$ , that is, they rather correspond to  $\eta^{-1}[1 - e^{-\eta|p|}]$  than just  $|p|$  (see [25], Theorem 7.12). Therefore, we will prove the discrete analogue of the operator inequality

$$(9) \quad \frac{1}{\eta} \prod_{r=1}^n (1 - e^{-\eta|p_r|}) \leq C(n) \sum_{r=1}^n \frac{1}{\eta} [1 - e^{-\eta|p_r|}]$$

on  $\mathbf{R}^n$ , and their quadratic forms will be the two Dirichlet forms  $\tilde{D}$  and  $D$ . Since the generator of  $\tilde{D}$  now averages in all directions, these averages yield traces of products  $\mathfrak{S}GA\mathfrak{S}GA \dots \mathfrak{S}GA$  for which we have a good local law; hence, the corresponding error in (8) is smaller than its naive a priori bound using only (1). This crucial gain provides the additional smallness to overcome the general fact that ultracontractivity bounds alone are barely not sufficient to gain sufficiently precise information on individual eigenvalues and eigenvectors in the bulk.

The actual proof requires several technical steps, such as: (i) localising the dynamics by considering a short-range approximation and treating the long-range part as a perturbation;

(ii) finite speed of propagation for the short range dynamics; (iii) cutoff the initial data in a smooth way so that cutoff and time evolution almost commute. Since these steps have appeared in the literature earlier, we will not reprove them here; we just refer to [26] where they have been adapted to the eigenvector moment flow. We will give full details only for Step 1.

Parallel with but independently of the current work, Benigni and Lopatto [6] have proved the CLT for  $\langle \mathbf{u}_i, A \mathbf{u}_i \rangle$  for the observable  $A$ , projecting onto a deterministic set of orthonormal vectors  $A = \sum_{\alpha \in I} |q_\alpha\rangle \langle q_\alpha|$ , with  $N^\epsilon \leq |I| \leq N^{1-\epsilon}$ , for some small fixed  $\epsilon > 0$ . Their low-rank assumption is complementary to our condition  $\langle \dot{A}^2 \rangle \geq c$  for this class of projection operators; moreover, their result also covered the edge regime. The low-rank assumption allowed them to operate with the eigenvector moment flow from [4, 12]. However, their control can handle overlaps with at most  $N^{1-\epsilon}$  vectors  $q_\alpha$  simultaneously. It seems that this approach has a natural limitation preventing it from using it for high rank observables, for example, for  $|I| \sim N$ . In contrast, we consider overlaps  $\langle \mathbf{u}_i, A \mathbf{u}_j \rangle$  directly without relying on the spectral decomposition of  $A$ .

*Notation and conventions.* We introduce some notations we use throughout the paper. For integers  $k \in \mathbf{N}$ , we use the notation  $[k] := \{1, \dots, k\}$ . For positive quantities  $f, g$ , we write  $f \lesssim g$  and  $f \sim g$  if  $f \leq Cg$  or  $cg \leq f \leq Cg$ , respectively, for some constants  $c, C > 0$  which depend only on the constants appearing in (10). We denote vectors by bold-faced lower case Roman letters  $\mathbf{x}, \mathbf{y} \in \mathbf{C}^N$ , for some  $N \in \mathbf{N}$ . Vector and matrix norms,  $\|\mathbf{x}\|$  and  $\|A\|$ , indicate the usual Euclidean norm and the corresponding induced matrix norm. For any  $N \times N$  matrix  $A$ , we use the notation  $\langle A \rangle := N^{-1} \text{Tr } A$  to denote the normalized trace of  $A$ . Moreover, for vectors  $\mathbf{x}, \mathbf{y} \in \mathbf{C}^N$  we define

$$\langle \mathbf{x}, \mathbf{y} \rangle := \sum \bar{x}_i y_i.$$

We will use the concept of “with very high probability,” meaning that, for any fixed  $D > 0$ , the probability of the  $N$ -dependent event is bigger than  $1 - N^{-D}$  if  $N \geq N_0(D)$ . Moreover, we use the convention that  $\xi > 0$  denotes an arbitrary small positive constant which is independent of  $N$ .

**2. Main results.** Let  $W$  be an  $N \times N$  real symmetric or complex Hermitian Wigner matrix. We formulate the following assumptions on  $W$ .

ASSUMPTION 2.1. We assume that the matrix elements  $w_{ab}$  are independent up to the Hermitian symmetry  $w_{ab} = \overline{w_{ba}}$  and identically distributed in the sense that  $w_{ab} \stackrel{d}{=} N^{-1/2} \chi_{od}$ , for  $a < b$ ,  $w_{aa} \stackrel{d}{=} N^{-1/2} \chi_d$ , with  $\chi_{od}$  being a real or complex random variable and  $\chi_d$  being a real random variable such that  $\mathbf{E} \chi_{od} = \mathbf{E} \chi_d = 0$  and  $\mathbf{E} |\chi_{od}|^2 = 1$ . In the complex case we also assume that  $\mathbf{E} \chi_{od}^2 = 0$ . In addition, we assume the existence of the high moments of  $\chi_{od}, \chi_d$ , that is, that there exist constants  $C_p > 0$ , for any  $p \in \mathbf{N}$ , such that

$$(10) \quad \mathbf{E} |\chi_d|^p + \mathbf{E} |\chi_{od}|^p \leq C_p.$$

Let  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$  be its eigenvalues in increasing order, and denote by  $\mathbf{u}_1, \dots, \mathbf{u}_N$  the corresponding orthonormal eigenvectors. For any  $N \times N$  matrix  $A$ , we denote by  $\dot{A} := A - \langle A \rangle$  the traceless part of  $A$ . We now state our main result.

THEOREM 2.2 (Central limit theorem in the QUE). *Let  $W$  be a real symmetric ( $\beta = 1$ ) or complex Hermitian ( $\beta = 2$ ) Wigner matrix satisfying Assumptions (2.1). Fix small  $\delta, \delta' > 0$ , and let  $A = A^*$  be a deterministic  $N \times N$  matrix with  $\|A\| \lesssim 1$  and  $\langle \dot{A}^2 \rangle \geq \delta'$ . In the real*

symmetric case we also assume that  $A \in \mathbf{R}^{N \times N}$  is real. Then, for any  $i \in [\delta N, (1 - \delta)N]$ , it holds

$$(11) \quad \sqrt{\frac{\beta N}{2(\mathring{A}^2)}} [\langle \mathbf{u}_i, A\mathbf{u}_i \rangle - \langle A \rangle] \Rightarrow \mathcal{N} \quad \text{as } N \rightarrow \infty$$

in the sense of moments, with  $\mathcal{N}$  being a standard real Gaussian random variable. The speed of convergence is explicit; see (27).

**3. Perfect matching observables.** For definiteness, we present the proof for the real symmetric case; the analysis for the complex Hermitian case is completely analogous and so omitted. We only mention that the main difference between the two symmetry classes is that the perfect matching observables  $f_{\lambda,t}$  in (20) are defined slightly differently (see [12], equation (A.3)), but the current proof can be easily adapted to this case.

Consider the matrix flow

$$(12) \quad dW_t = \frac{d\tilde{B}_t}{\sqrt{N}}, \quad W_0 = W,$$

with  $\tilde{B}_t$  being a standard real symmetric Brownian motion (see, e.g., [11], Definition 2.1). We denote the resolvent of  $W_t$  by  $G = G_t(z) := (W_t - z)^{-1}$ , for  $z \in \mathbf{C} \setminus \mathbf{R}$ . It is well known (see, e.g., [7, 21, 23]) that as  $N \rightarrow \infty$  the resolvent  $(W - z)^{-1}$  becomes approximately deterministic; its deterministic approximation is given by the unique solution of the scalar quadratic equation

$$(13) \quad -\frac{1}{m(z)} = z + m(z), \quad \Im m(z) \Im z > 0.$$

In particular,  $m(z) = m_{\text{sc}}(z)$ ,  $m_{\text{sc}}(z)$  being the Stieltjes transform of the semicircular law  $\rho_{\text{sc}}(x) := (2\pi)^{-1} \sqrt{(4 - x^2)_+}$ . The deterministic approximation of  $G_t(z)$  is given by  $m_t(z)$ , with  $m_t$  the solution of

$$(14) \quad \partial_t m_t(z) = -m_t \partial_z m_t(z), \quad m_0 = m.$$

From now on, by  $\rho_t = \rho_t(z)$  we denote  $\rho_t(z) := \pi^{-1} \Im m_t(z)$ , for any  $t \geq 0$ . In fact, starting from the standard semicircle  $\rho_0 = \rho_{\text{sc}}$ , the density  $\rho_t(x + i0)$  is just a rescaling of  $\rho_0$  by a factor  $2 + t$ .

By [11], Definition 2.2, it follows that the eigenvectors  $\mathbf{u}_1(t), \dots, \mathbf{u}_N(t)$  of  $W_t$ , corresponding to the eigenvalues  $\lambda_1(t) \leq \lambda_2(t) \leq \dots \leq \lambda_N(t)$ , are a solution of the following system of SDE (dropping the time dependence),

$$(15) \quad d\lambda_i = \frac{dB_{ii}}{\sqrt{N}} + \frac{1}{N} \sum_{j \neq i} \frac{1}{\lambda_i - \lambda_j} dt$$

$$(16) \quad d\mathbf{u}_k = \frac{1}{\sqrt{N}} \sum_{j \neq i} \frac{dB_{ij}}{\lambda_i - \lambda_j} \mathbf{u}_j - \frac{1}{2N} \sum_{j \neq i} \frac{\mathbf{u}_i}{(\lambda_i - \lambda_j)^2} dt,$$

with  $\{B_{ij}\}_{i,j \in [N]}$  being a standard real symmetric Brownian motions; see [11], Theorem 2.3, for the existence and uniqueness of the strong solution of (15)–(16).

By (16) it follows that the flow for the diagonal overlaps  $\langle \mathbf{u}_i, A\mathbf{u}_i \rangle$  naturally depends also on the off-diagonal overlap  $\langle \mathbf{u}_i, A\mathbf{u}_j \rangle$ ; hence, our analysis will concern not only diagonal overlaps but also off-diagonal ones. Since  $\langle \mathbf{u}_i, A\mathbf{u}_i \rangle - \langle A \rangle = \langle \mathbf{u}_i, \mathring{A}\mathbf{u}_i \rangle$  and  $\langle \mathbf{u}_i, A\mathbf{u}_j \rangle = \langle \mathbf{u}_i, \mathring{A}\mathbf{u}_j \rangle$  for  $i \neq j$ , without loss of generality we may assume for the rest of the paper, that  $A$  is traceless,  $\langle A \rangle = 0$ , that is,  $A = \mathring{A}$ . For traceless  $A$  we introduce the short-hand notation

$$(17) \quad p_{ij} = p_{ij}(t) := \langle \mathbf{u}_i(t), A\mathbf{u}_j(t) \rangle, \quad i, j \in [N].$$

We are now ready to write the flow for monomials of  $p_{ii}, p_{ij}$  (see [12], Theorem 2.6, for the derivation of the flow). For any fixed  $n$ , we will only need to consider monomials of the form  $\prod_{k=1}^n p_{i_k j_k}$  where each index appears even number of times; it turns out that the linear combinations of such monomials with a fixed degree  $n$  are invariant under the flow.

To encode general monomials, we use a particle picture (introduced in [11] and developed in [12, 26]) where each particle on the set of integers  $[N]$  corresponds to two occurrences of an index  $i$  in the monomial product. We use the same notation as in [12], and we define  $\eta : [N] \rightarrow \mathbf{N}$ , where  $\eta_j := \eta(j)$  is interpreted as the number of particles at the site  $j$ , and  $n(\eta) := \sum_j \eta_j = n$  denotes the total number of particles that are conserved under the flow. The space of  $n$ -particle configurations is denoted by  $\Omega^n$ . Moreover, for any index pair  $i \neq j \in [N]$ , we define  $\eta^{ij}$  to be the configuration obtained moving a particle from the site  $i$  to the site  $j$ ; if there is no particle in  $i$ , then we define  $\eta^{ij} = \eta$ . For any configuration  $\eta$ , consider the set of vertices

$$(18) \quad \mathcal{V}_\eta := \{(i, a) : 1 \leq i \leq n, 1 \leq a \leq 2\eta_i\},$$

and let  $\mathcal{G}_\eta$  be the set of perfect matchings on  $\mathcal{V}_\eta$ . Note that every particle configuration  $\eta$  gives rise to two vertices in  $\mathcal{V}_\eta$ ; thus, the elements of  $\mathcal{V}_\eta$  represent the indices in the product  $\prod_{k=1}^n p_{i_k j_k}$ .

There is no closed equation for individual products  $\prod_{k=1}^n p_{i_k j_k}$ , but there is one for a certain symmetrized linear combination; see [12], equation (2.15). Therefore, for any perfect matching  $G \in \mathcal{G}_\eta$ , we define

$$(19) \quad P(G) := \prod_{e \in \mathcal{E}(G)} p(e), \quad p(e) := p_{i_1 i_2},$$

where  $e = \{(i_1, a_1), (i_2, a_2)\} \in \mathcal{V}_\eta$ , and  $\mathcal{E}(G)$  denotes the edges of  $G$ . For example, for  $n = 2$  and for the configuration  $\eta$  defined by  $\eta(i) = \eta(j) = 1$  with some  $i \neq j$  and zero otherwise, we have three perfect matchings corresponding to  $p_{ii} p_{jj}$  and twice  $p_{ij}^2$ . For  $n = 3$  and  $\eta$  defined by  $\eta(i) = \eta(j) = \eta(k) = 1$ , we have 15 perfect matchings;  $p_{ii} p_{jj} p_{kk}$ , two copies  $p_{ij}^2 p_{kk}, p_{ik}^2 p_{jj}, p_{jk}^2 p_{ii}$  each, and eight copies of  $p_{ij} p_{jk} p_{ki}$ .

We are now ready to define the *perfect matching observable* for any given configuration  $\eta$ ,

$$(20) \quad f_{\lambda,t}(\eta) := \frac{N^{n/2}}{[2\langle A^2 \rangle]^{n/2}} \frac{1}{(n-1)!!} \frac{1}{\mathcal{M}(\eta)} \mathbf{E} \left[ \sum_{G \in \mathcal{G}_\eta} P(G) \middle| \lambda \right], \quad \mathcal{M}(\eta) := \prod_{i=1}^N (2\eta_i - 1)!!,$$

with  $n$  being the number of particles in the configuration  $\eta$ . Here, we took the conditioning on the entire flow of eigenvalues,  $\lambda = \{\lambda(t)\}_{t \in [0, T]}$  for some fixed  $T > 0$ . The observable  $f_{\lambda,t}$  satisfies a parabolic partial differential equation; see (23) below.

REMARK 3.1. For any  $k \in \mathbf{N}$ , the double factorial  $k!!$  is defined by  $k!! = k(k-2)!!$ ,  $1!! = 0!! = (-1)!! = 1$ . We remark that, in [12, 26], the authors use a different convention for the double factorial, that is, in these papers  $k!! = (k-1)(k-2)!!$ .

Note that  $f$  in (20) is defined slightly differently compared to the definition in [12], equation (2.15), where the authors do not have the additional  $(N/(2\langle A^2 \rangle))^{n/2} [(n-1)!!]^{-1}$  factor. Our normalisation factor is dictated by the principle that, for traceless  $A$ , we expect  $\sqrt{N} p_{ii} = \sqrt{N} [ \langle \mathbf{u}_i, A \mathbf{u}_i \rangle ]$  to be approximately a centred normal random variable with variance  $2\langle A^2 \rangle$ . In particular, the  $n$ th moment of  $(N/2\langle A^2 \rangle)^{1/2} p_{ii}$  for even  $n$  is close to  $(n-1)!!$ . Therefore, if  $\eta$  is a configuration with  $n$  particles all sitting at the same site  $i$ , that is,  $\eta(i) = n$  and zero otherwise, then  $\mathcal{M}(\eta) = (2n-1)!!$  is the number of perfect matchings, and, therefore, we expect  $f_{\lambda,t}(\eta) \approx 1$ . Note that using the a priori bound  $|p_{ij}| \leq N^{-1/2+\xi}$ , for any

$\xi > 0$ , proven in [15], Theorem 2.2, we have  $|f_{\lambda,t}| \lesssim N^\xi$  with very high probability, while the analogous quantity  $f_{\lambda,t}$ , defined in [12], equation (2.15), has an a priori bound of order  $N^{-n/2+\xi}$ .

We always assume that the entire eigenvalue trajectory  $\{\lambda(t)\}_{t \in [0,T]}$  satisfies the usual rigidity estimate (see, e.g., [19] or [21], Theorem 7.6). More precisely, for any fixed  $\xi > 0$ , we define

$$(21) \quad \Omega = \Omega_\xi := \left\{ \sup_{0 \leq t \leq T} \max_{i \in [N]} N^{2/3} \widehat{t}^{1/3} |\lambda_i(t) - \gamma_i(t)| \leq N^\xi \right\},$$

where  $\widehat{t} := i \wedge (N + 1 - i)$ , then we have

$$\mathbf{P}(\Omega_\xi) \geq 1 - C(\xi, D)N^{-D}$$

for any (small)  $\xi > 0$  and (large)  $D > 0$ . Here,  $\gamma_i(t)$  are the classical eigenvalue locations (*quantiles*) defined by

$$(22) \quad \int_{-\infty}^{\gamma_i(t)} \rho_t(x) dx = \frac{i}{N}, \quad i \in [N],$$

where  $\rho_t(x) = \frac{2}{(2+t)^2\pi} \sqrt{((2+t)^2 - x^2)_+}$  is the semicircle law corresponding to  $W_t$ . Note that  $|\gamma_i(t) - \gamma_i(s)| \lesssim |t - s|$  in the bulk, for any  $t, s \geq 0$ , as a consequence of the smoothness of  $t \rightarrow \rho_t$  in the bulk.

By [12], Theorem 2.6, we have that

$$(23) \quad \partial_t f_{\lambda,t} = \mathcal{B}(t) f_{\lambda,t},$$

$$(24) \quad \mathcal{B}(t) f_{\lambda,t} = \sum_{i \neq j} c_{ij}(t) 2\eta_i (1 + 2\eta_j) (f_{\lambda,t}(\boldsymbol{\eta}^{kl}) - f_{\lambda,t}(\boldsymbol{\eta})),$$

where

$$(25) \quad c_{ij}(t) := \frac{1}{N(\lambda_i(t) - \lambda_j(t))^2}.$$

Note that  $c_{ij}$  depends on  $\{\lambda(t)\}_{t \in [0,T]}$ , for some  $T > 0$ , but we omit this fact from the notation. We note that this flow was originally derived for special observables given in [12], equation (2.6), but the same derivation immediately holds for arbitrary  $A$  (see [12], Remark 2.8).

The main technical ingredient that will be used in the proof of Theorem 2.2 is the following proposition, whose proof is postponed to Section 4.

**PROPOSITION 3.2.** *For any  $n \in \mathbf{N}$ , there exists  $c(n) > 0$  such that, for any  $\epsilon > 0$  and for any  $T \geq N^{-1+\epsilon}$ , it holds*

$$(26) \quad \sup_{\boldsymbol{\eta}} |f_T(\boldsymbol{\eta}) - \mathbf{1}(n \text{ even})| \lesssim N^{-c(n)}$$

with very high probability, where the supremum is taken over configurations  $\boldsymbol{\eta}$  such that  $\sum_i \eta_i = n$  and  $\eta_i = 0$  for  $i \notin [\delta N, (1 - \delta)N]$ , with  $\delta > 0$  from Theorem 2.2. The implicit constant in (26) depends on  $n, \epsilon, \delta$ .

**PROOF OF THEOREM 2.2.** We fix  $i \in [\delta N, (1 - \delta)N]$ , and we choose  $\boldsymbol{\eta}$  to be the configuration  $\boldsymbol{\eta}$  with  $\eta_i = n$  and all other  $\eta_j = 0$ . Then, all the terms  $P(G)$  are equal to  $p_{ii}^n$  in the definition of  $f$ ; see (20). Then, using (26) for this particular  $\boldsymbol{\eta}$ , we conclude that

$$(27) \quad \mathbf{E} \left[ \sqrt{\frac{N}{2\langle A^2 \rangle}} \langle \mathbf{u}_i(T), A \mathbf{u}_i(T) \rangle \right]^n = \mathbf{1}(n \text{ even})(n - 1)!! + \mathcal{O}(N^{-c(n)})$$

for any  $i \in [\delta N, (1 - \delta)N]$  and  $T \gg N^{-1}$ , where we used that  $\|f_T\|_\infty \leq N^{n/2}$  deterministically on the complement of the high probability set on which (26) holds. With (27) we have proved that Theorem 2.2 holds for Wigner matrices with a small Gaussian component. For the general case, Theorem 2.2 follows from (27) and a standard application of the Green function comparison theorem (GFT), relating the eigenvectors/eigenvalues of  $W_T$  to those of  $W$ ; see the Appendix where we recall the argument for completeness.  $\square$

**4. DBM analysis.** In this section we focus on the analysis of the eigenvector moment flow (23)–(24). Since in our proof we use some results proven in [26], we start giving an equivalent representation of (20) which is the same used in [26] without distinguishing the several colours.

4.1. *Equivalent representation of the flow.* Fix  $n \in \mathbf{N}$ ; then, in the remainder of this section we will consider configurations  $\eta \in \Omega^n$ , that is, such that  $\sum_j \eta_j = n$ . Following [26] (but without the extra complication involving colours), we now give an equivalent representation of the flow (23)–(24) which will be defined on the  $2n$ -dimensional lattice  $[N]^{2n}$  instead of configurations of  $n$  particles. Let  $\mathbf{x} \in [N]^{2n}$ , and define

$$(28) \quad n_i(\mathbf{x}) := |\{a \in [2n] : x_a = i\}|$$

for all  $i \in \mathbf{N}$ . We define the configuration space

$$\Lambda^n := \{\mathbf{x} \in [N]^{2n} : n_i(\mathbf{x}) \text{ is even for every } i \in [N]\}.$$

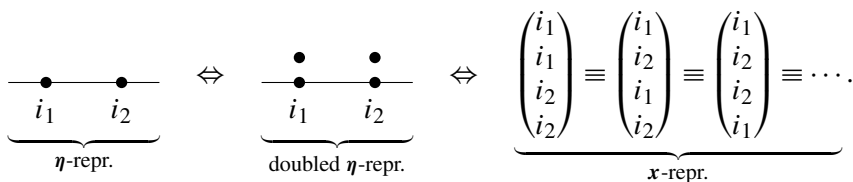
Note that  $\Lambda^n$  is an  $n$ -dimensional subset of the  $2n$  dimensional lattice  $[N]^{2n}$  in the sense that  $\Lambda^n$  is a finite union of  $n$ -dimensional sublattices of  $[N]^{2n}$ . From now on we will only consider configurations  $\mathbf{x} \in \Lambda^n$ . In particular, in this representation to each particle is associated a label  $a \in [2n]$ , that is, there is a particle at a site  $i \in [N]$  iff there exists  $a \in [2n]$  such that  $x_a = i$ . Additionally, by the definition of  $\Lambda^n$  it follows that the number of particles at a site  $i \in [N]$  is always even.

REMARK 4.1. Note that in [26] the authors consider  $\mathbf{x}$  to be an  $n$ -dimensional vector that lives in the  $n/2$ -dimensional subset  $\Lambda^n$ . For notational simplicity, in the current paper we assume that  $\mathbf{x}$  is a  $2n$ -dimensional vector and that  $\Lambda^n$  is  $n$ -dimensional.

The natural correspondence between the two representations is given by

$$(29) \quad \eta \leftrightarrow \mathbf{x}, \quad \eta_i = \frac{n_i(\mathbf{x})}{2}.$$

Note that  $\mathbf{x}$  uniquely determines  $\eta$ , but  $\eta$  determines only the coordinates of  $\mathbf{x}$  as a multiset and not its ordering. As an example, the configuration with single (or doubled) particles in  $i_1 \neq i_2$  corresponds to six  $\mathbf{x} \in \Lambda^2$ , as in



Let  $\phi: \Lambda^n \rightarrow \Omega^n$ ,  $\phi(\mathbf{x}) = \eta$  denote the map that projects the  $\mathbf{x}$ -configuration space to the  $\eta$ -configuration space using (29). This map naturally pulls back functions  $f$  of  $\eta$  to functions of  $\mathbf{x}$ ,

$$(\phi^* f)(\mathbf{x}) = f(\phi(\mathbf{x})).$$

We will always consider functions  $g$  on  $[N]^{2n}$  that are push-forward of some function  $f$  on  $\Omega^n$ ,  $g = f \circ \phi$ ; that is, they correspond to functions on the configurations

$$f(\boldsymbol{\eta}) = f(\phi(\mathbf{x})) = g(\mathbf{x}).$$

In particular,  $g$  is supported on  $\Lambda^n$ , and it is equivariant under permutation of the arguments, that is, it depends on  $\mathbf{x}$  only as a multiset. We, therefore, consider the observable

$$(30) \quad g_{\lambda,t}(\mathbf{x}) := f_{\lambda,t}(\phi(\mathbf{x})),$$

where  $f_{\lambda,t}$  was defined in (20). In the following we will often use the notation  $g_t(\mathbf{x}) = g_{\lambda,t}(\mathbf{x})$ , dropping the dependence of  $g_t(\mathbf{x})$  on the eigenvalues.

The flow (23)–(24) can be written in the  $\mathbf{x}$ -representation as follows:

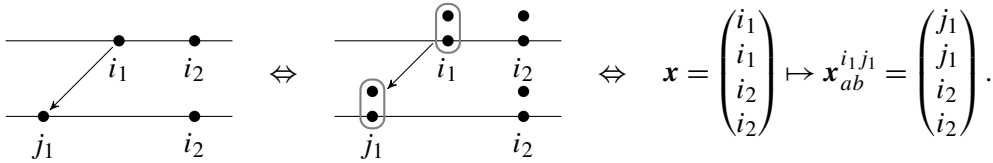
$$(31) \quad \partial_t g_t(\mathbf{x}) = \mathcal{L}(t)g_t(\mathbf{x})$$

$$(32) \quad \mathcal{L}(t) := \sum_{j \neq i} \mathcal{L}_{ij}(t), \quad \mathcal{L}_{ij}(t)g(\mathbf{x}) := c_{ij}(t) \frac{n_j(\mathbf{x}) + 1}{n_i(\mathbf{x}) - 1} \sum_{a \neq b \in [2n]} (g(\mathbf{x}_{ab}^{ij}) - g(\mathbf{x})),$$

where

$$(33) \quad \mathbf{x}_{ab}^{ij} := \mathbf{x} + \delta_{x_a i} \delta_{x_b j} (j - i)(\mathbf{e}_a + \mathbf{e}_b)$$

with  $\mathbf{e}_a(c) = \delta_{ac}$ ,  $a, c \in [2n]$ . Clearly, this flow preserves the equivariance of  $g$ ; that is, it is a map on functions defined on  $\Lambda^n$ . The jump operator  $\mathbf{x}_{ab}^{ij}$ , defined in (33), changes  $x_a, x_b$  from  $i$  to  $j$  if  $x_a = x_b = i$  and otherwise leaves  $\mathbf{x}$  unchanged. In the particle picture  $\boldsymbol{\eta}$  this corresponds in moving one particle from the site  $i$  (if there is any) to the site  $j$ ; see the following example for  $n = 2$  (with  $i = i_1, j = j_1$  and  $a = 1, b = 2$ ):



Define the measure

$$(34) \quad \pi(\mathbf{x}) := \prod_{i=1}^N ((n_i(\mathbf{x}) - 1)!)^2$$

on  $\Lambda^n$  and the corresponding  $L^2(\Lambda^n) = L^2(\Lambda^n, \pi)$  space equipped with the scalar product

$$(35) \quad \langle f, g \rangle_{\Lambda^n} = \langle f, g \rangle_{\Lambda^n, \pi} := \sum_{\mathbf{x} \in \Lambda^n} \pi(\mathbf{x}) \bar{f}(\mathbf{x}) g(\mathbf{x}).$$

We will often drop the dependence on the measure  $\pi$  in the scalar product. We also define the following norm on  $L^p(\Lambda^n)$ :

$$(36) \quad \|f\|_p := \left( \sum_{\mathbf{x} \in \Lambda^n} \pi(\mathbf{x}) |f(\mathbf{x})|^p \right)^{1/p}.$$

The measure  $\pi(\mathbf{x})$  clearly satisfies

$$(37) \quad 1 \leq \pi(\mathbf{x}) \leq (2n - 1)!,$$

uniformly in  $\mathbf{x} \in \Lambda^n$ . A direct calculation in [26], Appendix A.2, shows that the operator  $\mathcal{L} = \mathcal{L}(t)$  is symmetric with respect to the measure  $\pi$ , and it is a negative operator on the space  $L^2(\Lambda^n)$  with Dirichlet form

$$D(g) = \langle g, (-\mathcal{L})g \rangle_{\Lambda^n} = \frac{1}{2} \sum_{\mathbf{x} \in \Lambda^n} \pi(\mathbf{x}) \sum_{i \neq j} c_{ij}(t) \frac{n_j(\mathbf{x}) + 1}{n_i(\mathbf{x}) - 1} \sum_{a \neq b \in [2n]} |g(\mathbf{x}_{ab}^{ij}) - g(\mathbf{x})|^2.$$

We will often omit the time dependence of the generator  $\mathcal{L}(t)$ . We denote by  $\mathcal{U}(s, t)$  the semigroup associated to  $\mathcal{L}$  from (32); that is, for any  $0 \leq s \leq t$ , it holds

$$\partial_t \mathcal{U}(s, t) = \mathcal{L}(t)\mathcal{U}(s, t), \quad \mathcal{U}(s, s) = I.$$

4.2. *Short-range approximation.* Before proceeding we introduce a localised version of (31)–(32). Choose an ( $N$ -dependent) parameter  $1 \ll K \leq \sqrt{N}$ , and define the *averaging operator* as a simple multiplication operator by a “smooth” cut-off function,

$$(38) \quad \text{Av}(K, \mathbf{y})h(\mathbf{x}) := \text{Av}(\mathbf{x}; K, \mathbf{y})h(\mathbf{x}), \quad \text{Av}(\mathbf{x}; K, \mathbf{y}) := \frac{1}{K} \sum_{j=K}^{2K-1} \mathbf{1}(\|\mathbf{x} - \mathbf{y}\|_1 < j),$$

with  $\|\mathbf{x} - \mathbf{y}\|_1 := \sum_{a=1}^{2n} |x_a - y_a|$ . While it was denoted and called averaging operator in [12, 26], it is rather a *localization*, that is, a multiplication by a “smooth” cutoff function  $\mathbf{x} \rightarrow \text{Av}(\mathbf{x}; K, \mathbf{y})$  which is centered at  $\mathbf{y}$  and has a soft range of size  $K$ . The parameters  $K, \mathbf{y}$  are considered fixed and often omitted from the notation.

Now, we define a short-range version of the dynamics (31). Fix an integer  $\ell$  with  $1 \ll \ell \ll K$ , and define the short-range coefficients

$$(39) \quad c_{ij}^{\mathcal{S}}(t) := \begin{cases} c_{ij}(t) & \text{if } i, j \in \mathcal{J} \text{ and } |i - j| \leq \ell, \\ 0 & \text{otherwise,} \end{cases}$$

where  $c_{ij}(t)$  is defined in (25). Here,

$$(40) \quad \mathcal{J} = \mathcal{J}_\delta := \{i \in [N] : \gamma_i(0) \in \mathcal{I}_\delta\}, \quad \mathcal{I}_\delta := (-2 + \delta, 2 - \delta)$$

with  $\delta > 0$  from Theorem 2.2 so that  $\mathcal{I}_\delta$  lies entirely in the bulk spectrum.

We define  $h_t(\mathbf{x})$  as the time evolution of a localized initial data  $g_0$  by the short-range dynamics,

$$(41) \quad \begin{aligned} h_0(\mathbf{x}; \ell, K, \mathbf{y}) &= h_0(\mathbf{x}; K, \mathbf{y}) := \text{Av}(\mathbf{x}; K, \mathbf{y})(g_0(\mathbf{x}) - \mathbf{1}(n \text{ even})), \\ \partial_t h_t(\mathbf{x}; \ell, K, \mathbf{y}) &= \mathcal{S}(t)h_t(\mathbf{x}; \ell, K, \mathbf{y}), \end{aligned}$$

where

$$(42) \quad \mathcal{S}(t) := \sum_{j \neq i} \mathcal{S}_{ij}(t), \quad \mathcal{S}_{ij}(t)h(\mathbf{x}) := c_{ij}^{\mathcal{S}}(t) \frac{n_j(\mathbf{x}) + 1}{n_i(\mathbf{x}) - 1} \sum_{a \neq b \in [2n]} (h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x})).$$

Here, we used the notation  $h(\mathbf{x}) = h(\mathbf{x}; \ell, K, \mathbf{y})$  to indicate all relevant parameters:  $\ell$  indicates the short range of the dynamics,  $\mathbf{y}$  is the centre, and  $K$  is the range of the cut-off in the initial condition, and we always choose  $\ell \ll K$ . In (41) we already subtracted  $\mathbf{1}(n \text{ even})$  since in our application the initial condition  $g_0(\mathbf{x})$  after some local averaging will be close to  $\mathbf{1}(n \text{ even})$ ; hence, after longer time we expect that  $h_t$  tends to zero since the dynamics has a smoothing effect, and it is an  $L^1$  contraction.

4.3.  *$L^2$ -bound.* Define the distance on  $\Lambda^n$  as

$$(43) \quad d(\mathbf{x}, \mathbf{y}) := \sup_{a \in [2n]} |\mathcal{J} \cap [\min(x_a, y_a), \max(x_a, y_a))|$$

with  $\mathcal{J}$  defined in (40). Note that  $d$  is not a metric since it is degenerate, but it is still symmetric and satisfies the triangle inequality [26], equation (5.6). The key ingredient to prove the  $L^2$ -bound in (49) below is to show that the short-range dynamics (41)–(42) is close to the original dynamics (31)–(32). This will be achieved using the following finite speed of propagation estimate, proven in [10], Theorem 2.1, Lemma 2.4, [26], Proposition 5.2, (see also

[12], equation (3.15)), for  $\mathcal{U}_S(s, t) = \mathcal{U}_S(s, t; \ell)$  which is the transition semigroup associated to the short-range generator  $\mathcal{S}(t)$ . For any  $\mathbf{x} \in \Lambda^n$ , define the “delta-function” on  $\Lambda^n$  as

$$\delta_{\mathbf{x}}(\mathbf{u}) := \begin{cases} \pi(\mathbf{x})^{-1} & \text{if } \mathbf{u} = \mathbf{x}, \\ 0 & \text{otherwise,} \end{cases}$$

and denote the matrix entries of  $\mathcal{U}_S(s, t)$  by  $\mathcal{U}_S(s, t)_{\mathbf{x}\mathbf{y}} := \langle \delta_{\mathbf{x}}, \mathcal{U}_S(s, t)\delta_{\mathbf{y}} \rangle$ .

PROPOSITION 4.2. *Fix any small  $\epsilon > 0$  and  $\ell \geq N^\epsilon$ . Then, for any  $\mathbf{x}, \mathbf{y} \in \Lambda^n$  with  $d(\mathbf{x}, \mathbf{y}) > N^\epsilon \ell$ , it holds*

$$(44) \quad \sup_{0 \leq s_1 \leq s_2 \leq s_1 + \ell N^{-1}} |\mathcal{U}_S(s_1, s_2; \ell)_{\mathbf{x}\mathbf{y}}| \leq e^{-N^\epsilon/2}$$

on the very high probability event  $\Omega$ .

This finite speed of propagation, together with the fact that the initial condition  $h_0$  is localized in a  $K$ -neighbourhood of a fixed center  $\mathbf{y}$ , implies that  $h_t$  is supported in a  $K + N^\epsilon \ell \leq 2K$  neighbourhood of  $\mathbf{y}$  up an exponentially small tail part.

Using Proposition 4.2, by [26], Corollary 5.3, we immediately conclude the following lemma.

LEMMA 4.3. *For any times  $s_1, s_2$  such that  $0 \leq s_1 \leq s_2 \leq s_1 + \ell N^{-1}$  and for any  $\mathbf{y} \in \Lambda^n$  supported on  $\mathcal{J}$  (i.e.,  $\mathbf{y}_a \in \mathcal{J}$  for any  $a \in [2n]$ ) for the commutator of the evolution  $\mathcal{U}_S$  and the averaging operator, we have*

$$(45) \quad \|\mathcal{U}_S(s_1, s_2; \ell), \text{Av}(\mathbf{y}, K)\|_{\infty, \infty} \leq C(n) \frac{N^\epsilon \ell}{K}$$

for some constant  $C(n) > 0$  and for any small  $\epsilon > 0$ , on the very high probability event  $\Omega$ .

Another straightforward application of the finite speed of propagation estimate in Proposition 4.2 is the following bound  $\mathcal{U}(s_1, s_2) - \mathcal{U}_S(s_1, s_2; \ell)$ . This result was proven in [26], Proposition 5.7, for a specific  $f$ , but the same proof applies for a general function  $f$ .

LEMMA 4.4. *Let  $0 \leq s_1 \leq s_2 \leq s_1 + \ell N^{-1}$ , and  $f$  is a function on  $\Lambda^n$ ; then, for any  $\mathbf{x} \in \Lambda^n$  supported on  $\mathcal{J}$ , it holds*

$$(46) \quad |(\mathcal{U}(s_1, s_2) - \mathcal{U}_S(s_1, s_2; \ell))f(\mathbf{x})| \lesssim N^{1+n\xi} \frac{s_2 - s_1}{\ell} \|f\|_\infty$$

for any small  $\xi > 0$ .

PROOF. Using Proposition 4.2, the proof of (46) is completely analogous to the proof of [26], Proposition 5.7, since the only input used in [26], Proposition 5.7, is that

$$\sum_{j:|j-i|>\ell} \frac{1}{N(\lambda_i - \lambda_j)^2} \leq \frac{N^{1+\xi}}{\ell}$$

on  $\Omega$  which follows by rigidity.  $\square$

Before stating the main result of this section, we define the set  $\widehat{\Omega}$  on which the local laws for certain products of resolvents and traceless matrices  $A$  hold; that is, for a small

$\omega > 2\xi > 0$ , we define

$$\begin{aligned}
 \widehat{\Omega} = \widehat{\Omega}_{\omega, \xi} := & \bigcap_{\substack{z_t: \exists z_i \in \mathcal{I}_\delta \\ |\exists z_i| \in [N^{-1+\omega}, 10]}} \left[ \bigcap_{k=3}^n \left\{ \sup_{0 \leq t \leq T} \left| \langle G_t(z_1)A \dots G_t(z_k)A \rangle \right| \leq \frac{N^{\xi+(k-3)/2}}{\sqrt{\eta_*}} \right\} \right. \\
 (47) \quad & \cap \left\{ \sup_{0 \leq t \leq T} \left| \langle G_t(z_1)AG_t(z_2)A \rangle - m_t(z_1)m_t(z_2)\langle A^2 \rangle \right| \leq \frac{N^\xi}{\sqrt{N\eta_*}} \right\} \\
 & \left. \cap \left\{ \sup_{0 \leq t \leq T} \left| \langle G_t(z_1)A \rangle \right| \leq \frac{N^\xi}{N\sqrt{|\exists z_1|}} \right\} \right],
 \end{aligned}$$

where  $\eta_* := \min\{|\exists z_i| \mid i \in [k]\}$ . The fact that  $\widehat{\Omega}$  is a very high probability set follows by [15], Theorem 2.6, for  $k = 1$ , by [15], equation (3.9), for  $k = 2$ , and by Proposition 5.1 for  $k \geq 3$ . In particular, since  $m_t(z_1)m_t(z_2)\langle A^2 \rangle$  is bounded for  $k = 2$ , we have

$$\sup_{0 \leq t \leq T} \sup_{z_1, z_2} \langle \exists G_t(z_1)A \exists G_t(z_2)A \rangle \lesssim 1,$$

on the very high probability event  $\widehat{\Omega}_{\omega, \xi}$  which, by spectral theorem, implies

$$(48) \quad \sup_{0 \leq t \leq T} \max_{i, j \in \mathcal{I}} \left| \langle \mathbf{u}_i(t), \mathbf{A} \mathbf{u}_j(t) \rangle \right| \leq N^{-1/2+\omega} \quad \text{on } \widehat{\Omega}_{\omega, \xi}.$$

PROPOSITION 4.5. *For any scale satisfying  $N^{-1} \ll \eta \ll T_1 \ll \ell N^{-1} \ll KN^{-1}$  and any small  $\epsilon, \xi > 0$ , it holds*

$$(49) \quad \left\| h_{T_1}(\cdot; \ell, K, \mathbf{y}) \right\|_2 \lesssim K^{n/2} \mathcal{E},$$

with

$$(50) \quad \mathcal{E} := N^{n\xi} \left( \frac{N^\epsilon \ell}{K} + \frac{NT_1}{\ell} + \frac{N\eta}{\ell} + \frac{N^\epsilon}{\sqrt{N\eta}} + \frac{1}{\sqrt{K}} \right),$$

uniformly for particle configuration  $\mathbf{y} \in \Lambda^n$  supported on  $\mathcal{J}$  and eigenvalue trajectory  $\lambda$  on the high-probability event  $\Omega_\xi \cap \widehat{\Omega}_{\omega, \xi}$ .

PROOF. Before presenting the formal proof, we explain the main idea. In the sense of Dirichlet forms, we will replace the generator  $\mathcal{S}(t)$  (41)–(42), which is the *sum* of one-dimensional generators, with the generator  $\mathcal{A}(t)$  that corresponds to the *product* of such operators (see (52) below for its definition). Considering that  $c_{ij}$  decays proportionally with  $|i - j|^{-2}$  (using rigidity in (25)), it is the kernel of the discrete approximation of the one dimensional operator  $|p| = \sqrt{-\Delta}$  on  $\mathbf{R}$  but lifted to the  $n$ -dimensional space  $\Lambda^n$ . Therefore, one may think of  $\mathcal{L}(t)$ , and its short-range approximation  $\mathcal{S}(t)$ , as a discrete analogue of  $|p_1| + |p_2| + \dots + |p_n|$ , that is, the sum of  $|p|$ -operators along all the  $n$  coordinate directions. As explained in the **Introduction**, using the short distance regularisation of the underlying lattice, we really have  $\eta^{-1}[1 - e^{-\eta|p|}]$  instead of  $|p|$ , and the operator inequality (9) holds. The left-hand side of (9) corresponds to the positive operator  $(-\mathcal{A})$ , and the right-hand side corresponds to  $(-\mathcal{S})$ . The key Lemma 4.6 below asserts that  $0 \leq (-\mathcal{A}) \leq C(n)(-\mathcal{S})$  in the sense of quadratic forms. The main purpose of this replacement is that  $\mathcal{A}$  averages independently in every direction; therefore,  $\mathcal{A}$ , acting on the function  $g = f \circ \phi$ , has the effect that it averages in *all* the  $i_1, i_2, \dots$  indices in the definition of  $P(G)$ , (19). These averages yield traces of products  $\exists GA \exists GA \dots \exists GA$  for which we have a good local law on the set  $\widehat{\Omega}$ .

We now explain the origin of the errors in (50). The error in the multi- $G$  local laws give the crucial fourth error  $1/\sqrt{N\eta}$  in (50). The other errors come from various approximations:

the dynamics commute with the localization up to an error of order  $\ell/K$  by Lemma 4.3, the short-range cutoff dynamics approximates the original one up to time  $T_1$  with an error of order  $NT_1/\ell$ , while the removed long-range part contributes with an error of order  $N\eta/\ell$  to the Dirichlet form. The last  $1/\sqrt{K}$  error term is technical; we do the analysis for typical index configurations where no two indices coincide and the coinciding indices have a volume factor of order  $1/\sqrt{K}$  smaller than the total volume.

Now, we start with the actual proof. All the estimates in this proof hold uniformly for  $\mathbf{y} \in \Lambda^n$  supported on  $\mathcal{J}$ ; hence, from now on we fix a particle configuration  $\mathbf{y}$ . To make the presentation clearer, we drop the parameters  $\mathbf{y}, K, \ell$  and use the short-hand notations  $h_t(\mathbf{x}) = h_t(\mathbf{x}; \ell, K, \mathbf{y})$ ,  $\text{Av} = \text{Av}(K, \mathbf{y})$ ,  $\text{Av}(\mathbf{x}) = \text{Av}(\mathbf{x}; K, \mathbf{y})$ , etc. Moreover, for any  $\mathbf{i}, \mathbf{j} \in [N]^n$  by  $\sum_{\mathbf{i}}^*$  or  $\sum_{\mathbf{i}\mathbf{j}}^*$ , we denote the summations over indices that are all distinct; that is, the  $i_1, \dots, i_n$ , in the first sum, and  $i_1, \dots, i_n, j_1, \dots, j_n$ , in the second sum are all different. The same convention holds for summations over  $\mathbf{a}, \mathbf{b} \in [2n]^n$ .

Let

$$(51) \quad a_{ij} = a_{ij}(t) := \frac{\eta}{N((\lambda_i(t) - \lambda_j(t))^2 + \eta^2)},$$

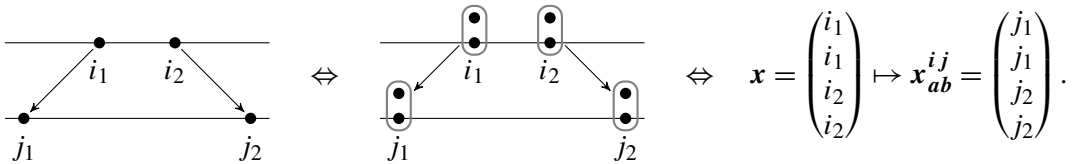
and define their short-range version  $a_{ij,S}$ , as in (39). Define the operator  $\mathcal{A} = \mathcal{A}(t)$  by

$$(52) \quad \mathcal{A}(t) := \sum_{\mathbf{i}, \mathbf{j} \in [N]^n}^* \mathcal{A}_{\mathbf{i}\mathbf{j}}(t), \quad \mathcal{A}_{\mathbf{i}\mathbf{j}}(t)h(\mathbf{x}) := \frac{1}{\eta} \left( \prod_{r=1}^n a_{i_r, j_r}^S(t) \right) \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* (h(\mathbf{x}_{\mathbf{a}\mathbf{b}}^{i\mathbf{j}}) - h(\mathbf{x})),$$

where

$$(53) \quad \mathbf{x}_{\mathbf{a}\mathbf{b}}^{i\mathbf{j}} := \mathbf{x} + \left( \prod_{r=1}^n \delta_{x_{a_r}, i_r} \delta_{x_{b_r}, i_r} \right) \sum_{r=1}^n (j_r - i_r) (\mathbf{e}_{a_r} + \mathbf{e}_{b_r}).$$

We now explain the difference between the jump operator (53) and the one defined in (33). The operator (33) changes two entries of  $\mathbf{x}$  per time, instead  $\mathbf{x}_{\mathbf{a}\mathbf{b}}^{i\mathbf{j}}$  changes all the coordinates of  $\mathbf{x}$  at the same time; that is, let  $\mathbf{i} := (i_1, \dots, i_n), \mathbf{j} := (j_1, \dots, j_n) \in [N]^n$ , with  $\{i_1, \dots, i_n\} \cap \{j_1, \dots, j_n\} = \emptyset$ , then  $\mathbf{x}_{\mathbf{a}\mathbf{b}}^{i\mathbf{j}} \neq \mathbf{x}$  iff for all  $r \in [n]$  it holds that  $x_{a_r} = x_{b_r} = i_r$ ; see, for example,



Note that  $\mu(\mathbf{x}) \equiv 1$  on  $\Lambda^n$  is a reversible measure for the generator  $\mathcal{A}(t)$  (as a consequence of  $(\mathbf{x}_{\mathbf{a}\mathbf{b}}^{i\mathbf{j}})_{\mathbf{a}\mathbf{b}}^{j\mathbf{i}} = \mathbf{x}$  for any fixed  $\mathbf{a}, \mathbf{b}$  and for any  $\mathbf{x}$  such that  $\mathbf{x}_{\mathbf{a}\mathbf{b}}^{i\mathbf{j}} \neq \mathbf{x}$ ), and that  $\pi(\mathbf{x}) \sim C(n)\mu(\mathbf{x})$  for all  $\mathbf{x} \in \Lambda^n$  (see (37)). We define the scalar product with respect to the measure  $\mu(\mathbf{x})$  analogously to (35), and we denote it by  $\langle \cdot, \cdot \rangle_{\Lambda^n, \mu}$ .

We now analyse the time evolution of  $\|h_t\|_2^2$ ,

$$(54) \quad \partial_t \|h_t\|_2^2 = 2\langle h_t, \mathcal{S}(t)h_t \rangle_{\Lambda^n}.$$

The main ingredient to give an upper bound on (54) is the following lemma, whose proof is postponed at the end of this section.

LEMMA 4.6. *Let  $\mathcal{S}(t), \mathcal{A}(t)$  be the generators defined in (42) and (52), respectively. Then there, exists a constant  $C(n) > 0$ , which depends only on  $n$ , such that*

$$(55) \quad \langle h, \mathcal{S}(t)h \rangle_{\Lambda^n, \pi} \leq C(n) \langle h, \mathcal{A}(t)h \rangle_{\Lambda^n, \mu} \leq 0$$

for any  $h \in L^2(\Lambda^n)$ , on the very high probability set  $\Omega$ .

From now on, by  $C(n)$  we denote a constant that depends only on  $n$  and that may change from line to line.

Next, combining (54)–(55) and using that  $\mathbf{x}_{ab}^{ij} = \mathbf{x}$  unless  $\mathbf{x}_{a_r} = \mathbf{x}_{b_r} = i_r$  for all  $r \in [n]$ , we conclude that

$$\begin{aligned} \partial_t \|h_t\|_2^2 &\leq C(n) \langle h_t, \mathcal{A}(t)h_t \rangle_{\Lambda^n, \mu} \\ (56) \quad &= \frac{C(n)}{2\eta} \sum_{\mathbf{x} \in \Lambda^n} \sum_{\mathbf{i}, \mathbf{j} \in [N]^n}^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \right) \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \bar{h}_t(\mathbf{x}) (h_t(\mathbf{x}_{ab}^{ij}) - h_t(\mathbf{x})) \Psi(\mathbf{x}), \end{aligned}$$

where for any fixed  $\mathbf{i}, \mathbf{a}, \mathbf{b}$  we defined

$$\Psi(\mathbf{x}) = \Psi_{\mathbf{i}, \mathbf{a}, \mathbf{b}}(\mathbf{x}) := \left( \prod_{r=1}^n \delta_{x_{a_r} i_r} \delta_{x_{b_r} i_r} \right).$$

Define

$$(57) \quad \Gamma := \{ \mathbf{x} \in \Lambda^n : d(\mathbf{x}, \mathbf{y}) \leq 3K \} \subset \Lambda^n,$$

and note that, by the finite speed of propagation estimate in Proposition 4.2 and the support of  $h_0(\mathbf{x})$ , the function  $h_t(\mathbf{x})$  is supported on  $\Gamma$  up to an exponentially small error term (see [26], equations (5.76)–(5.77), for a more detailed calculation). For simplicity, for the rest of the proof we treat  $h_t(\mathbf{x})$  as if it were supported on  $\Gamma$ , neglecting the exponentially small error term of size  $\pi(\Lambda^n \setminus \Gamma)e^{-N^\epsilon} \leq N^{2n}e^{-N^\epsilon}$ . Since the dynamics is a linear contraction in  $L^\infty$ , this small error term remains small throughout the whole evolution.

Now, we consider the term with  $|h_t(\mathbf{x})|^2$  in (56) (here, we use the notation  $\Psi(\mathbf{x}) = \Psi_{\mathbf{i}, \mathbf{a}, \mathbf{b}}(\mathbf{x})$ ),

$$\begin{aligned} & - \sum_{\mathbf{x} \in \Gamma} |h_t(\mathbf{x})|^2 \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \sum_{\mathbf{i}, \mathbf{j}}^* \Psi(\mathbf{x}) \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \right) \\ (58) \quad &= - \sum_{\mathbf{x} \in \Gamma} |h_t(\mathbf{x})|^2 \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \sum_{\mathbf{i}}^* \Psi(\mathbf{x}) \prod_{r=1}^n \left( \sum_{j_r} a_{i_r j_r}^S(t) + \mathcal{O}\left(\frac{N^\xi}{N\eta}\right) \right) \\ &= - \sum_{\mathbf{x} \in \Gamma} |h_t(\mathbf{x})|^2 \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \sum_{\mathbf{i}}^* \Psi(\mathbf{x}) \prod_{r=1}^n \left( \sum_{j_r} a_{i_r j_r}^S(t) + \mathcal{O}\left(\frac{N^\xi}{N\eta} + \frac{N^{1+\xi}\eta}{\ell}\right) \right) \\ &\leq -C(n) \sum_{\mathbf{x} \in \Gamma} |h_t(\mathbf{x})|^2 \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \sum_{\mathbf{i}}^* \Psi(\mathbf{x}) \end{aligned}$$

on the very high probability event  $\Omega_\xi$ , where the error term in the second line comes from adding back the finitely many excluded summands  $j_r \in \{i_1, \dots, i_n\}$  and  $j_r \in \{j_1, \dots, j_{r-1}, j_{r+1}, \dots, j_n\}$ . The new error in the third line comes from removing the short-range restriction from  $a_{i_r j_r}^S$ , that is, adding back the regimes  $|j_r - i_r| > \ell$  using

$$(59) \quad \sum_{j_r: |j_r - i_r| > \ell} a_{i_r j_r}(t) \leq \frac{N\eta}{\ell}.$$

Finally, to go from the third to the fourth line in (58) we used the local law

$$(60) \quad \sum_{j_r} a_{i_r j_r}(t) = \langle \Im G_t(\lambda_{i_r} + i\eta) \rangle = \Im m_t(\lambda_{i_r} + i\eta) + \mathcal{O}(N^\xi (N\eta)^{-1})$$

with very high probability on the event  $\Omega_\xi$ , and that  $\Im m_t(\lambda_{i_r} + i\eta) \sim 1$  in the bulk whenever  $\eta \geq N^{-1+\xi}$ .

We now bound the last line in (58) in terms of  $-\|h_t\|^2$  plus a small error term by removing the restriction from the  $i$ -summation,

$$\begin{aligned}
 (61) \quad & - \sum_{x \in \Gamma} |h_t(x)|^2 \sum_{a,b \in [2n]^n}^* \sum_i^* \Psi(x) = - \sum_{x \in \Gamma} |h_t(x)|^2 \sum_{a,b \in [2n]^n}^* \sum_i^* \Psi(x) \\
 & + \sum_{x \in \Gamma} |h_t(x)|^2 \sum_{a,b \in [2n]^n}^* \left( \sum_i^* - \sum_i^* \right) \Psi(x) \\
 & \leq -C(n)\|h_t\|_2^2 + C(n)N^\xi K^{n-1}.
 \end{aligned}$$

To estimate the first term in the right-hand side, we used that  $\sum_{ab}^* \sum_i \Psi_{i,a,b}(x) \geq 1$ ,  $1 \geq C(n)\pi(x)$  for all  $x \in \Gamma$ , and that we can add back the regime  $\Lambda^n \setminus \Gamma$  at the price of a negligible  $N^n e^{-N^\epsilon}$  error term by finite speed of propagation. For the second term in the right-hand side of (61), we estimated  $\|h_t\|_\infty \leq N^\xi$  as a consequence of  $\|h_0\|_\infty \leq N^\xi$  and the fact that the evolution is an  $L^\infty$ -contraction. Finally we used the fact that

$$\sum_{a,b \in [2n]^n}^* \left( \sum_i^* - \sum_i^* \right) \Psi_{i,a,b}(x) \neq 0,$$

only if there exist  $a, b, c, d \in [2n]$ , all distinct, such that  $x_a = x_b = x_c = x_d$ . The volume of this one codimensional subset of  $\Gamma$  is  $C(n)K^{n-1}$ , that is, by factor  $K^{-1}$  smaller than the volume of  $\Gamma$  which is of order  $K^n$ .

Finally, combining (58) and (61), we conclude the estimate for the term containing  $|h_t(x)|^2$  in (56),

$$(62) \quad - \sum_{x \in \Gamma} |h_t(x)|^2 \sum_{a,b \in [2n]^n}^* \sum_{i,j}^* \Psi(x) \left( \prod_{i=1}^r a_{i_r j_r}^S(t) \right) \leq -C_1(n)\|h_t\|_2^2 + C(n)N^\xi K^{n-1}.$$

Then, using (56) together with (62), we conclude that

$$\begin{aligned}
 (63) \quad & \partial_t \|h_t\|_2^2 \leq -\frac{C_1(n)}{\eta} \|h_t\|_2^2 + \frac{C(n)N^\xi K^{n-1}}{\eta} \\
 & + \frac{C_2(n)}{\eta} \sum_{x \in \Gamma} |h_t(x)| \sum_{a,b \in [2n]^n}^* \sum_i^* \Psi(x) \left| \sum_j^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \right) h_t(x_{ab}^{ij}) \right|
 \end{aligned}$$

for some constants  $C_1(n), C_2(n) > 0$ , on the event  $\Omega_\xi$  with  $\xi > 0$  arbitrarily small. In order to conclude the bound of  $\partial_t \|h_t\|_2^2$ , we are now left with the estimate of the last line in (63).

In the remainder of the proof, we will show that

$$\begin{aligned}
 (64) \quad & \frac{C_2(n)}{\eta} \sum_{x \in \Gamma} |h_t(x)| \sum_{a,b \in [2n]^n}^* \sum_i^* \Psi(x) \left| \sum_j^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \right) h_t(x_{ab}^{ij}) \right| \\
 & \leq \frac{C_1(n)}{2\eta} \|h_t\|_2^2 + \frac{C_3(n)}{\eta} \mathcal{E}^2 K^n
 \end{aligned}$$

with  $\mathcal{E}$  defined in (50) and  $C_1(n)$  being the constant from the first line of (63). Note that, using (64), we readily conclude the proof of (49) by

$$(65) \quad \partial_t \|h_t\|_2^2 \leq -\frac{C_1(n)}{2\eta} \|h_t\|_2^2 + \frac{C_3(n)}{\eta} \mathcal{E}^2 K^n,$$

which implies  $\|h_{T_1}\|_2^2 \leq C(n)\mathcal{E}^2 K^n$ , by a simple Gronwall inequality, using that  $T_1 \gg \eta$ .

We now conclude the proof of (49) proving the bound in (64). We start with the analysis of

$$(66) \quad \sum_{j \in [N]^n}^* \left( \prod_{r=1}^n a_{i_r, j_r}^{\mathcal{S}}(t) \right) h_t(\mathbf{x}_{ab}^{ij})$$

for any fixed  $\mathbf{x} \in \Gamma$ ,  $\mathbf{i} \in [N]^n$ ,  $\mathbf{a}, \mathbf{b} \in [2n]^n$  with all distinct coordinates such that  $\Psi(\mathbf{x}) \neq 0$ . It will be very important that the configuration  $\phi(\mathbf{x}_{ab}^{ij})$  contains exactly one particle at every index  $j_r$ ; that is, we have

$$(67) \quad \prod_{l=1}^N (n_l(\mathbf{x}_{ab}^{ij}) - 1)!! = 1.$$

Similarly to [26], equations (5.89)–(5.91), equations (5.95)–(5.97), using that the function  $f(\mathbf{x}) \equiv \mathbf{1}(n \text{ even})$  is in the kernel of  $\mathcal{S}(t)$ , for any fixed  $\mathbf{x} \in \Gamma$ , and for any fixed  $\mathbf{i}, \mathbf{a}, \mathbf{b}$ , we conclude that

$$(68) \quad \begin{aligned} & h_t(\mathbf{x}_{ab}^{ij}) \\ &= \mathcal{U}_{\mathcal{S}}(0, t)((\text{Av } g_0)(\mathbf{x}_{ab}^{ij}) - (\text{Av } \mathbf{1}(n \text{ even}))(\mathbf{x}_{ab}^{ij})) \\ &= \text{Av}(\mathbf{x}_{ab}^{ij})(\mathcal{U}_{\mathcal{S}}(0, t)g_0(\mathbf{x}_{ab}^{ij}) - \mathbf{1}(n \text{ even})) + \mathcal{O}\left(\frac{N^{\epsilon+n\xi}\ell}{K}\right) \\ &= \left(\text{Av}(\mathbf{x}) + \mathcal{O}\left(\frac{\ell}{K}\right)\right)\left(\mathcal{U}(0, t)g_0(\mathbf{x}_{ab}^{ij}) - \mathbf{1}(n \text{ even}) + \mathcal{O}\left(\frac{N^{1+n\xi}t}{\ell}\right)\right) \\ &\quad + \mathcal{O}\left(\frac{N^{\epsilon+n\xi}\ell}{K}\right) \\ &= \text{Av}(\mathbf{x})(g_t(\mathbf{x}_{ab}^{ij}) - \mathbf{1}(n \text{ even})) + \mathcal{O}\left(\frac{N^{\epsilon+n\xi}\ell}{K} + \frac{N^{1+n\xi}t}{\ell}\right), \end{aligned}$$

where the error terms are uniform in  $\mathbf{x} \in \Gamma$ . Note that to go from the first to the second line in (68) we used Lemma 4.3, to go from the second to the third line we used Lemma 4.4 together with the a priori bound  $\|g_t\|_{\infty} \leq N^{n\xi}$  for any  $0 \leq t \leq T$  on the very high probability event  $\widehat{\Omega}_{\omega, \xi}$ , and that

$$|\text{Av}(\mathbf{x}) - \text{Av}(\mathbf{x}_{ab}^{ij})| \leq \frac{1}{K} \|\mathbf{x} - \mathbf{x}_{ab}^{ij}\|_1 \leq \frac{2n\ell}{K},$$

where  $\|\mathbf{x}\|_1 = \sum_{c=1}^{2n} |x_c|$ . To go from the third to the fourth line in (68), we used that  $|\text{Av}(\mathbf{x})| \leq 1$  and again that  $\|g_t\|_{\infty} \leq N^{n\xi}$ . Then, from (68), we conclude that

$$(69) \quad \begin{aligned} \sum_j^* \left( \prod_{r=1}^n a_{i_r, j_r}^{\mathcal{S}}(t) \right) h_t(\mathbf{x}_{ab}^{ij}) &= \text{Av}(\mathbf{x}) \sum_j^* \left( \prod_{r=1}^n a_{i_r, j_r}^{\mathcal{S}}(t) \right) (g_t(\mathbf{x}_{a,b}^{ij}) - \mathbf{1}(n \text{ even})) \\ &\quad + \mathcal{O}\left(\frac{N^{\epsilon+n\xi}\ell}{K} + \frac{N^{1+n\xi}T_1}{\ell}\right). \end{aligned}$$

From now on, we will omit the  $\text{Av}(\mathbf{x})$  prefactor in (69), since  $|\text{Av}(\mathbf{x})| \leq 1$ .

Using the definition of  $g_t$  from (30) and (20), for any  $\mathbf{x} \in \Gamma$  such that  $\Psi(\mathbf{x}) \neq 0$ , and for any fixed  $\mathbf{i}, \mathbf{a}, \mathbf{b}$ , dropping the  $t$ -dependence of the eigenvalues  $\lambda_i = \lambda_i(t)$ , we have

$$\begin{aligned}
 & \sum_j^* \left( \prod_{r=1}^n a_{i_r j_r}^{\mathcal{S}}(t) \right) (g_t(\mathbf{x}_{ab}^{ij}) - \mathbf{1}(n \text{ even})) \\
 &= \sum_j^* \left( \prod_{r=1}^n a_{i_r j_r}(t) \right) \left( \frac{N^{n/2}}{\langle A^2 \rangle^{n/2} 2^{n/2} (n-1)!!} \sum_{G \in \mathcal{G}_{\eta^j}} P(G) - \mathbf{1}(n \text{ even}) \right) \\
 (70) \quad &+ \mathcal{O}\left(\frac{N^{1+\xi} \eta}{\ell}\right) \\
 &= \sum_j \left( \prod_{r=1}^n a_{i_r j_r}(t) \right) \left( \frac{N^{n/2}}{\langle A^2 \rangle^{n/2} 2^{n/2} (n-1)!!} \sum_{G \in \mathcal{G}_{\eta^j}} P(G) - \mathbf{1}(n \text{ even}) \right) \\
 &+ \mathcal{O}\left(\frac{N^\xi}{N\eta} + \frac{N^{1+\xi} \eta}{\ell}\right).
 \end{aligned}$$

Note that in (70) we used the notation  $\eta^j := \phi(\mathbf{x}_{ab}^{ij})$  to denote the particle configuration which has exactly one particle at each site  $\{j_1, \dots, j_n\}$ . Note that in the last line of (70) we do not exclude the possibility that two indices  $j$  may assume the same value, since the sum is unrestricted. In the second and third lines of (70), we simply omitted the conditional expectation  $\mathbf{E}[\cdot \cdot | \lambda]$  to shorten the formulas. Since all subsequent estimates hold with high probability, the conditional expectation does not play a role. When going from the first to the second line of (70), we removed the short-range restriction, as in (59), by adding back the summations over the regimes  $|j_r - i_r| > \ell$ , and we also used (67) since the coordinates of  $\mathbf{j}$  are all distinct, and so that  $\mathcal{M}(\eta^j) = 1$  in the definition of  $g_t$  in (30) and (20). Additionally, the error term in the third line of (70) comes from adding back the missing  $j_r$ -summations; in this bound we used the a priori bound  $|P(G)| \leq N^{\xi-n/2}$  on the very high probability event  $\widehat{\Omega}_{\omega, \xi}$  and (60).

We now use the definition of  $P(G)$  in (19) on the right-hand side of (70). Since every particle is doubled, we may rewrite the sum over perfect matchings as

$$(71) \quad \sum_{G \in \mathcal{G}_{\eta^j}} P(G) = \sum_{G \in \text{Gr}_2[n]} \prod_{(v_1 \dots v_k) \in \text{Cyc}(G)} (2k-2)!! p_{j_{v_1} j_{v_2}} \cdots p_{j_{v_k} j_{v_1}},$$

where  $\text{Gr}_2[n]$  denotes the set of 2-regular multigraphs (possibly with loop-edges) on  $[n]$  and  $\text{Cyc}(G)$  denoting the collection of cycles in any such graph  $G \in \text{Gr}_2[n]$ . The combinatorial factor  $(2k-2)!!$  is due to the fact that, for each cycle in  $G$ , there are  $(2k-2)!!$  equivalent perfect matchings giving the very same cyclic monomial. For example, for  $n=2$  there are two 2-regular multi-graphs, (11), (22) and (12), (12), and thus  $\sum_{G \in \mathcal{G}_{\eta^j}} P(G) = 2p_{j_1 j_2}^2 + p_{j_1 j_1} p_{j_2 j_2}$ . Similarly, for  $n=3$ , there are the graphs

$$\begin{aligned}
 & \{(11), (22), (33)\}, \quad \{(12), (12), (33)\}, \quad \{(13), (13), (22)\}, \\
 & \{(23), (23), (11)\}, \quad \{(12), (23), (13)\}
 \end{aligned}$$

yielding

$$\begin{aligned}
 \sum_{G \in \mathcal{G}_{\eta^j}} P(G) &= p_{j_1 j_1} p_{j_2 j_2} p_{j_3 j_3} + 2p_{j_1 j_1} p_{j_2 j_3}^2 + 2p_{j_2 j_2} p_{j_1 j_3}^2 \\
 &+ 2p_{j_3 j_3} p_{j_1 j_2}^2 + 8p_{j_1 j_2} p_{j_2 j_3} p_{j_1 j_3}.
 \end{aligned}$$

For each graph  $G \in \text{Gr}_2[n]$ , we may use the spectral theorem to perform the  $\mathbf{j}$  summation as

$$(72) \quad \sum_{j_{v_1}, \dots, j_{v_k}} \left( \prod_{r \in [k]} a_{i_{v_r} j_{v_r}}(t) \right) p_{j_{v_1} j_{v_2}} \cdots p_{j_{v_k} j_{v_1}} = N^{1-k} F_k(v_1, \dots, v_k)$$

with

$$F_k(v_1, \dots, v_k) := \langle \mathfrak{S}G_t(\lambda_{i_{v_1}} + i\eta)A \cdots \mathfrak{S}G_t(\lambda_{i_{v_k}} + i\eta)A \rangle.$$

Since each vertex appears in exactly one cycle, we can use (72) to perform the summation for the indices corresponding to any cycle separately and obtain

$$(73) \quad \sum_j \left( \prod_{r=1}^n a_{i_r j_r}(t) \right) \sum_{G \in \mathcal{G}_{\eta, j}} P(G) = \sum_{E \in \text{Gr}_2[n]} \prod_{(v_1 \cdots v_k) \in \text{Cyc}(E)} (2k-2)!! N^{1-k} F_k(v_1, \dots, v_k).$$

We note that, from (47) for each  $k \geq 1$ , we have the estimate

$$(74) \quad F_k(v_1, \dots, v_k) = \mathbf{1}(k=2) \langle A^2 \rangle \mathfrak{S}m(z_{i_{v_1}}) \mathfrak{S}m(z_{i_{v_k}}) + \mathcal{O}\left(N^\xi \frac{N^{k/2-1}}{\sqrt{N\eta}}\right)$$

on the high-probability set  $\widehat{\Omega}$ . By using (74) within (73) and using the fact that there are  $\mathbf{1}(n \text{ even})(n-1)!!$  graphs in  $\text{Gr}_2[n]$  all of which cycles have length two, it follows that

$$(75) \quad (73) = \mathbf{1}(n \text{ even})(n-1)!! 2^{n/2} N^{-n/2} \langle A^2 \rangle^{n/2} \prod_{r \in [n]} \mathfrak{S}m(z_{i_r}) + \mathcal{O}\left(N^\xi \frac{N^{-n/2}}{\sqrt{N\eta}}\right),$$

and from (70) we conclude

$$(76) \quad \Psi(\mathbf{x}) \sum_j^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \right) (g_t(\mathbf{x}_{ab}^{ij}) - \mathbf{1}(n \text{ even})) = \Psi(\mathbf{x}) \mathcal{O}\left(\frac{N^\xi}{\sqrt{N\eta}} + \frac{N^\xi}{N\eta} + \frac{N^{1+\xi}\eta}{\ell}\right).$$

Combining (69) and (76), we get that

$$(77) \quad \Psi(\mathbf{x}) \left| \sum_j^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \right) h_t(\mathbf{x}_{ab}^{ij}) \right| \leq C(n) \Psi(\mathbf{x}) \mathcal{E},$$

and finally, by (77), we conclude that

$$(78) \quad \text{l.h.s. (64)} \leq \frac{C_1(n)}{2\eta} \|h_t\|_2^2 + \frac{C_3(n)}{\eta} \mathcal{E}^2 K^n,$$

where we used that, for any fixed  $\mathbf{x} \in \Lambda^n$ , we have

$$\sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \sum_i^* \Psi_{i, \mathbf{a}, \mathbf{b}}(\mathbf{x}) \leq C(n),$$

and that

$$\left| \sum_{\mathbf{x} \in \Gamma} h_t(\mathbf{x}) \mathcal{E} \right| \leq \frac{C_1(n)}{2} \sum_{\mathbf{x} \in \Gamma} \pi(\mathbf{x}) |h_t(\mathbf{x})|^2 + C_3(n) \mathcal{E}^2 K^n,$$

by the Schwarz inequality, the bound  $1 \leq \pi(\mathbf{x})$  from (37), and  $\sum_{\mathbf{x} \in \Gamma} \pi(\mathbf{x}) \leq C(n) K^n$ . Note that, by balancing between the two terms in the Schwarz inequality, we could achieve the

same constant  $C_1(n)$  with an additional  $1/2$  factor in front of the  $\|h_t\|_2^2$  term as in the leading term in (63) with a minus sign. This concludes the proof of the bound in (64).  $\square$

**PROOF OF LEMMA 4.6.** All along the proof  $C(n) > 0$  is a constant that depends only on  $n$  and that may change from line to line.

We consider

$$\begin{aligned}
 \langle h, \mathcal{S}(t)h \rangle_{\Lambda^n, \pi} &= -\frac{1}{2} \sum_{\mathbf{x} \in \Lambda^n} \pi(\mathbf{x}) \sum_{j \neq i} c_{ij}^{\mathcal{S}}(t) \frac{n_j(\mathbf{x}) + 1}{n_i(\mathbf{x}) - 1} \sum_{a \neq b \in [2n]} |h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x})|^2 \\
 (79) \qquad \qquad \qquad &\leq -\frac{C(n)}{\eta} \sum_{\mathbf{x} \in \Lambda^n} \sum_{j \neq i} a_{ij}^{\mathcal{S}}(t) \sum_{a \neq b \in [2n]} |h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x})|^2
 \end{aligned}$$

and

$$(80) \qquad \langle h, \mathcal{A}(t)h \rangle_{\Lambda^n, \mu} = -\frac{1}{2\eta} \sum_{\mathbf{x} \in \Lambda^n} \sum_{i, j}^* \left( \prod_{r=1}^n a_{i_r j_r}^{\mathcal{S}}(t) \right) \sum_{a, b \in [2n]^n}^* |h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x})|^2.$$

Note that in (79) we used that  $a_{ij}^{\mathcal{S}}(t) \leq \eta c_{ij}^{\mathcal{S}}(t)$  to compare the kernels, that  $\pi(\mathbf{x}) \geq 1$  uniformly in  $\mathbf{x} \in \Lambda^n$ , and finally that  $n_j(\mathbf{x}) + 1 \geq 1$ ,  $1 \leq n_i(\mathbf{x}) - 1 \leq n$  for  $\mathbf{x}$  and  $i$  such that  $h(\mathbf{x}_{ab}^{ij}) \neq h(\mathbf{x})$ .

We start with the bound

$$\begin{aligned}
 (81) \qquad \sum_{\mathbf{x} \in \Lambda^n} \sum_{i, j \in [N]^n}^* \left( \prod_{r=1}^n a_{i_r j_r}^{\mathcal{S}}(t) \mathbf{1}(n_{i_r}(\mathbf{x}) > 0) \right) \sum_{a, b \in [2n]^n}^* |h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x})|^2 \\
 \leq C(n) \sum_{\mathbf{x} \in \Lambda^n} \sum_{i, j}^* \left( \prod_{r=1}^n a_{i_r j_r}^{\mathcal{S}}(t) \mathbf{1}(n_{i_r}(\mathbf{x}) > 0) \right) \sum_{l=1}^n \sum_{a, b \in [2n]^n}^* |h((\mathbf{y}_{l-1})_{ab}^{i_l j_l}) - h(\mathbf{y}_{l-1})|^2,
 \end{aligned}$$

where we recursively defined  $\mathbf{y}_0 = \mathbf{x}$ ,  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n = \mathbf{x}_{ab}^{ij}$  by performing the jumps  $i_1 \rightarrow j_1$ ,  $i_2 \rightarrow j_2$ , etc., one by one (assuming that the choice of  $(a_l, b_l)$  allows it), otherwise  $\mathbf{y}_l = \mathbf{y}_{l-1}$ ,

$$(82) \qquad \mathbf{y}_0 = \mathbf{y}_0(\mathbf{x}) := \mathbf{x}, \qquad \mathbf{y}_l = \mathbf{y}_l(\mathbf{x}) := (\mathbf{y}_{l-1})_{ab}^{i_l j_l}.$$

In the first line of (81), we could add the indicator  $\mathbf{1}(n_{i_r}(\mathbf{x}) > 0)$  since in case  $n_{i_r}(\mathbf{x}) = 0$  for some  $r$  it holds that  $\mathbf{x}_{ab}^{ij} = \mathbf{x}$ . Note that, to go from the first to the second line of (81), we wrote a telescopic sum

$$h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x}) = \sum_{l=1}^n [h((\mathbf{y}_{l-1})_{ab}^{i_l j_l}) - h(\mathbf{y}_{l-1})]$$

and used Schwarz inequality.

Next, we consider

$$\begin{aligned}
 & \sum_{l=1}^n \sum_{\mathbf{x} \in \Lambda^n} \sum_{i,j}^* \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \mathbf{1}(n_{i_r}(\mathbf{x}) > 0) \right) |h(\mathbf{y}_{l-1})_{a_l b_l}^{i_l j_l} - h(\mathbf{y}_{l-1})|^2 \\
 &= \sum_{l=1}^n \sum_{\mathbf{w} \in \Lambda^n} \sum_{i,j}^* \sum_{\mathbf{a}, \mathbf{b} \in [2n]^n}^* \left( \prod_{r=1}^n a_{i_r j_r}^S(t) \mathbf{1}(n_{i_r}(\mathbf{z}_{l-1}) > 0) \right) |h(\mathbf{w}_{a_l b_l}^{i_l j_l}) - h(\mathbf{w})|^2 \\
 (83) \quad &\leq C(n) \sum_{\mathbf{w} \in \Lambda^n} \sum_{l=1}^n \sum_{i_l \neq j_l} a_{i_l j_l}^S(t) \sum_{a_l \neq b_l \in [2n]} |h(\mathbf{w}_{a_l b_l}^{i_l j_l}) - h(\mathbf{w})|^2 \\
 &\quad \times \left( \prod_{r \neq l} \sum_{i_r, j_r} a_{i_r j_r} [\mathbf{1}(n_{i_r}(\mathbf{w}) > 0) + \mathbf{1}(n_{j_r}(\mathbf{w}) > 0)] \right) \\
 &\leq C(n) \sum_{\mathbf{w} \in \Lambda^n} \sum_{l=1}^n \sum_{i_l \neq j_l} a_{i_l j_l}^S(t) \sum_{a_l \neq b_l \in [2n]} |h(\mathbf{w}_{a_l b_l}^{i_l j_l}) - h(\mathbf{w})|^2 \\
 &\leq C(n) \sum_{\mathbf{w} \in \Lambda^n} \sum_{i \neq j} a_{ij}^S(t) \sum_{a \neq b \in [2n]} |h(\mathbf{w}_{ab}^{ij}) - h(\mathbf{w})|^2.
 \end{aligned}$$

Note that, to go from the first to the second line, we did the change of variables  $\mathbf{w} = \mathbf{y}_{l-1}(\mathbf{x})$ , we used that  $(\mathbf{x}_{a_l b_l}^{i_l j_l})_{a_l b_l}^{j_l i_l} = \mathbf{x}$  for any  $\mathbf{x} \in \Lambda^n$  such that  $\prod_r \mathbf{1}(n_{i_r}(\mathbf{x}) > 0)$ , and we defined  $\mathbf{z}_{l-1} = ((\mathbf{w}_{a_{l-1} b_{l-1}}^{j_{l-1} i_{l-1}}) \dots)_{a_1 b_1}^{j_1 i_1}$ . Moreover, to go from the second to the third line in (83) we used that

$$\begin{aligned}
 (84) \quad & \prod_{r \in [n] \setminus \{l\}} \mathbf{1}(n_{i_r}(\mathbf{z}_{l-1}) > 0) \leq C(n) \left( \prod_{r=1}^{l-1} [\mathbf{1}(n_{j_r}(\mathbf{w}) > 0) + \mathbf{1}(n_{i_r}(\mathbf{w}) > 0)] \right) \\
 & \quad \times \left( \prod_{r=l+1}^n \mathbf{1}(n_{i_r}(\mathbf{w}) > 0) \right)
 \end{aligned}$$

for  $i_1, \dots, i_n, j_1, \dots, j_n$  all distinct, which follows by  $n_{i_r}(\mathbf{z}_{l-1}) = n_{i_r}(\mathbf{w})$  if  $r \geq l + 1$  and

$$\mathbf{1}(n_{i_r}(\mathbf{z}_{l-1}) > 0) \leq \mathbf{1}(n_{i_r}(\mathbf{w}) > 0) + \mathbf{1}(n_{j_r}(\mathbf{w}) > 0)$$

for  $r \leq l - 1$ . In the penultimate inequality in (83), we also used that

$$(85) \quad \prod_{r \neq l} \sum_{i_r, j_r} a_{i_r j_r} [\mathbf{1}(n_{i_r}(\mathbf{w}) > 0) + \mathbf{1}(n_{j_r}(\mathbf{w}) > 0)] \leq C(n)$$

on the very high probability event  $\widehat{\Omega}$ . Combining (79)–(80), (81), and (83), we finally conclude (55).  $\square$

4.4. *Proof of Proposition 3.2.* Fix  $1 \ll NT_1 \ll \ell_1 \ll K$  and  $1 \ll NT_2 \ll \ell_2 \ll K$ , with  $T_1 \leq T_2/2$ . Define the lattice generator  $\mathcal{W}(t)$  by

$$(86) \quad \mathcal{W}(t) := \sum_{i \neq j \in [N]} \mathcal{W}_{ij}(t), \mathcal{W}_{ij}(t) := c_{ij}^{\mathcal{W}}(t) \frac{n_j(\mathbf{x}) + 1}{n_i(\mathbf{x}) - 1} \sum_{a \neq b \in [2n]} (h(\mathbf{x}_{ab}^{ij}) - h(\mathbf{x}))$$

with

$$(87) \quad c_{ij}^{\mathcal{W}}(t) := \begin{cases} c_{ij}(t) & \text{if } i, j \in \mathcal{J} \text{ and } 1 \leq |i - j| \leq \ell_2, \\ \frac{N}{|i - j|^2} & \text{otherwise.} \end{cases}$$

Denote by  $\mathcal{U}_{\mathcal{W}}(s, t)$  the semigroup associated to the generator  $\mathcal{W}(t)$ . Note that  $\mathcal{W}(t)$  is the original generator of the Dyson eigenvector flow  $\mathcal{L}$  from (32) on short scales and in the interval  $\mathcal{J}$  well inside the bulk, while on large scales it has an equidistant jump rate. In [26] this replacement made up for the missing rigidity (regularity) control of the eigenvalues outside of a local interval  $\mathcal{J}$ ; in our case its role is just to handle the somewhat different scaling of the eigenvalues near the edges. We follow the setup of [26] for convenience.

On the event  $\Omega_\xi$ , the coefficients  $c_{ij}^{\mathcal{W}}(t)$  satisfy [26], Assumption 6.8, with a rate  $v = N^{1-\xi}$ , for any arbitrary small  $\xi > 0$ ; hence, all the results in [26], Section 6, apply to the generator  $\mathcal{W}(t)$ . Most importantly, the Dirichlet form of  $\mathcal{W}(t)$  satisfies a Poincaré inequality and, consequently, we have an  $L^2 \rightarrow L^\infty$  ultracontractive decay bound for the corresponding semigroup. Their scaling properties confirm the intuition that  $\mathcal{W}(t)$  is a discrete analogue of the  $|p| = \sqrt{-\Delta}$  operator in  $\mathbf{R}^{2n}$ . In the continuous setting, standard Sobolev inequality combined with the Nash method implies that

$$(88) \quad \|e^{-t|p|} f\|_{L^\infty(\mathbf{R}^{2n})} \leq \frac{C(n)}{t^{n/2}} \|f\|_{L^2(\mathbf{R}^{2n})}$$

holds for any  $L^2$  function on  $\mathbf{R}^{2n}$ . The same decay holds for the semigroup generated by  $\mathcal{W}(t)$  by [26], Proposition 6.29, (recall that [26] uses  $n$  to denote the dimension of the space of  $\mathbf{x}$ 's, we use  $2n$ ). We remark that the proofs in [26], Section 6, are designed for the more involved *coloured* dynamics; here, we need only its simpler *colourblind* version which immediately follows from the coloured version by ignoring the colors. In particular, in our case the *exchange operator*  $\mathcal{E}_{ij}$  is identically zero. While a direct proof of the colourblind version is possible and it would require less combinatorial complexity, for brevity, we directly use the results of [26], Section 6.

For each  $\mathbf{y}$  supported on  $\mathcal{J}$ , let  $q_t(\mathbf{x}) = q_t(\mathbf{x}; \mathbf{y})$  be the solution of

$$(89) \quad \begin{cases} q_0(\mathbf{x}) = Av(\mathbf{x}; K, \mathbf{y})(g_0(\mathbf{x}) - \mathbf{1}(n \text{ even})), \\ \partial_t q_t(\mathbf{x}) = \mathcal{S}(t)q_t(\mathbf{x}) & \text{for } 0 < t \leq T_1, \\ \partial_t q_t(\mathbf{x}) = \mathcal{W}(t)q_t(\mathbf{x}) & \text{for } T_1 < t \leq T_2 \end{cases}$$

with  $\mathcal{S}(t)$  being the short-range generator on a scale  $\ell = \ell_1$  from (42). Note that  $q_t = h_t$ , for any  $0 \leq t \leq T_1$ , with  $h_t$  being the solution of (41).

By Proposition 4.5, choosing  $\eta = N^{-\epsilon} T_1$ , we have

$$(90) \quad \sup_{\mathbf{y}: y_a \in \mathcal{J}} \|q_{T_1}(\cdot; \mathbf{y})\|_2 \lesssim N^\epsilon K^{n/2} \left( \frac{\ell_1}{K} + \frac{NT_1}{\ell_1} + \frac{1}{\sqrt{NT_1}} + \frac{1}{\sqrt{K}} \right)$$

for any arbitrary small  $\epsilon > 0$ , where the supremum is over all the  $\mathbf{y}$  supported on  $\mathcal{J}$ . We recall that, by the finite speed of propagation estimate in Proposition 4.2, together with [26], equation (7.12), the function  $q_t$  is supported on the subset of  $\Gamma \subset \Lambda^n$  such that  $d(\mathbf{x}, \mathbf{y}) \leq 3K$  for any  $\mathbf{x} \in \Gamma$  (modulo a negligible exponentially small error term). Then, using the ultracontractivity bound for the dynamics of  $\mathcal{W}(t)$  from [26], Proposition 6.29, with  $v = N^{1-\xi}$ , we get that

$$(91) \quad \begin{aligned} \sup_{\mathbf{y}} \|q_{T_2}(\cdot; \mathbf{y})\|_\infty &= \sup_{\mathbf{y}} \|\mathcal{U}_{\mathcal{W}}(T_1, T_2)q_{T_1}(\cdot; \mathbf{y})\|_\infty \\ &\leq \sup_{\mathbf{y}} \|(1 - \Pi)\mathcal{U}_{\mathcal{W}}(T_1, T_2)q_{T_1}(\cdot; \mathbf{y})\|_\infty \\ &\quad + \sup_{\mathbf{y}} \|\Pi\mathcal{U}_{\mathcal{W}}(T_1, T_2)q_{T_1}(\cdot; \mathbf{y})\|_\infty \\ &\lesssim \frac{N^{n\xi}}{[N(T_2 - T_1)]^{n/2}} \sup_{\mathbf{y}} \|q_{T_1}(\cdot; \mathbf{y})\|_2 + \frac{K^n}{N^{n(1-\xi)}}, \end{aligned}$$

where  $\Pi$  is the orthogonal projection into the kernel of  $\mathcal{L}$ ,  $\ker(\mathcal{L}) = \bigcap_{i \neq j} \ker(\mathcal{L}_{ij})$ , defined in [26], Lemma 4.17. Note that in (91) we used that by [26], Corollary 4.20, it holds

$$\|\Pi q_{T_2}\|_\infty \leq C(n)N^{-n} \|q_{T_2}\|_1 \leq K^n N^{-n+n\xi},$$

since  $\|q_{T_2}\|_\infty \leq N^{n\xi}$  on the very high probability set  $\widehat{\Omega}$ . We remark that in [26], Proposition 6.29,  $\mathcal{U}_{\mathcal{W}}$  is replaced by  $\mathcal{U}$ , but this does not play any role since the only assumption on  $\mathcal{L}_{ij}$ , used in [26], Section 6, is that  $c_{ij}(t) \geq N^{1-\xi} |i - j|^{-2}$  (see [26], Definition 6.8). Combining (90)–(91), we conclude

$$(92) \quad \sup_y \|q_{T_2}(\cdot; y)\|_\infty \lesssim N^{2\epsilon} \left(\frac{K}{NT_2}\right)^{n/2} \left(\frac{\ell_1}{K} + \frac{NT_1}{\ell_1} + \frac{1}{\sqrt{NT_1}} + \frac{1}{\sqrt{K}}\right),$$

where we used that  $T_1 \leq T_2/2$ .

Now, we compare the solution  $q_t$  from (89) with the original dynamics  $g_t$  from (31). This is done, after several steps, using [26], Proposition 7.2, with  $F_t(\mathbf{y}; \mathbf{y})$  replaced by  $\mathbf{1}(n \text{ even})$ , asserting that

$$(93) \quad \sup_y |q_{T_2}(\mathbf{y}; \mathbf{y}) - (g_{T_2}(\mathbf{y}) - \mathbf{1}(n \text{ even}))| \lesssim N^\epsilon \left(\frac{\ell_1}{K} + \frac{NT_1}{\ell_1} + \frac{\ell_2}{K} + \frac{NT_2}{\ell_2}\right).$$

In particular, the only thing used about  $F_t(\mathbf{y}; \mathbf{y})$  in the proof of [26], Proposition 7.2, is that  $F_t$  is in the kernel of all  $\mathcal{L}_{ij}$ , and this is clearly the case for  $\mathbf{1}(n \text{ even})$  as well. The origins of the error terms in (93) are as follows. The smooth cutoff given by the Av localising operator in the initial condition (89) commutes with the time evolution generated by  $\mathcal{S}$  up an error of order  $\ell_1/K$ ; see Lemma 4.3. The difference between the original dynamics and the short-range dynamics in the time interval  $t \in [0, T_1]$  yields the error  $NT_1/\ell_1$ ; see Lemma 4.4. Similar errors hold for the approximation of the original dynamics by the time evolution generated by  $\mathcal{W}$  on the time interval  $t \in [T_1, T_2]$ , giving rise to the errors  $\ell_2/K$  and  $N(T_2 - T_1)/\ell_2 \leq NT_2/\ell_2$ .

Combining (92)–(93), we conclude that

$$(94) \quad \begin{aligned} & \sup_y |g_{T_2}(\mathbf{y}) - \mathbf{1}(n \text{ even})| \\ & \lesssim N^{2\epsilon} \left(\frac{\ell_1}{K} + \frac{NT_1}{\ell_1} + \frac{\ell_2}{K} + \frac{NT_2}{\ell_2} + \left(\frac{K}{NT_2}\right)^{n/2} \left(\frac{\ell_1}{K} + \frac{NT_1}{\ell_1} + \frac{1}{\sqrt{NT_1}} + \frac{1}{\sqrt{K}}\right)\right) \\ & \lesssim N^{-c/(20n)}, \end{aligned}$$

on the with very high probability event  $\Omega_\xi \cap \widehat{\Omega}_{\xi, \epsilon}$  with choosing a very small  $\xi$ . In the last step we optimised the error terms in the second line of (94) with the choice of

$$K = N^c, \quad T_2 = N^{-1-c/(10n)} K, \quad \ell_2 = \sqrt{NKT_2}, \quad \ell_1 = \sqrt{NKT_1}, \quad T_1 = \frac{\sqrt{K}}{N}$$

with some small fixed  $0 < c \leq 1/2$ . Finally, using that

$$\sup_{y: y_a \in \mathcal{J}} |g_{T_2}(\mathbf{y}) - \mathbf{1}(n \text{ even})| = \sup_\eta |f_{T_2}(\boldsymbol{\eta}) - \mathbf{1}(n \text{ even})|$$

by (30), where the supremum in the right-hand side is taken over configurations  $\boldsymbol{\eta}$  such that  $\eta_i = 0$  for  $i \in [\delta N, (1 - \delta)N]^c$  and  $\sum_i \eta_i = n$ . The bound in (94) concludes the proof of Proposition 3.2.

**5. Local law bounds.** In this section we prove the local laws needed to estimate the probability of the event  $\widehat{\Omega}$  in (47). We recall [21] that the resolvent  $G = (W - z)^{-1}$  of the Wigner matrix  $W$  is approximately equal,

$$(95) \quad G_{ab} = \delta_{ab}m + \mathcal{O}\left(\frac{N^\xi}{\sqrt{N\Im z}}\right), \quad \langle G \rangle = m + \mathcal{O}\left(\frac{N^\xi}{N\Im z}\right)$$

to the Stieltjes transform  $m = m_{\text{sc}}(z)$  of the semicircular distribution  $\rho_{\text{sc}} = \sqrt{4 - x^2}/2\pi$  which solves the equation

$$(96) \quad -\frac{1}{m} = m + z.$$

**PROPOSITION 5.1.** *Let  $k \geq 3$  and  $z_1, \dots, z_k \in \mathbf{C} \setminus \mathbf{R}$  with  $N \min_i (\rho_i \eta_i) \geq N^\epsilon$  for some  $\epsilon > 0$  with  $\eta_i := |\Im z_i|$  and  $\rho_i := \rho(z_i)$ ,  $\rho(z) := |\Im m(z)|/\pi$ . Then, for arbitrary traceless matrices  $A_1, \dots, A_k$  with  $\|A_i\| \lesssim 1$ , we have*

$$(97) \quad |\langle G_1 A_1 \dots G_k A_k \rangle| \lesssim N^{\xi + (k-3)/2} \sqrt{\frac{\rho^*}{\eta_*}}$$

with very high probability for any  $\xi > 0$ , where  $\rho^* := \max_i \rho_i$  and  $\eta_* := \min \eta_i$ .

**PROOF.** Using  $WG - zG = I$  and (96), we write

$$(98) \quad G = m - m \underline{WG} + m \langle G - m \rangle G,$$

where

$$\underline{WG} = WG + \langle G \rangle G$$

denotes a renormalization of  $WG$ . More generally, for functions  $f(W)$  we define

$$\underline{Wf(W)} := Wf(W) - \widetilde{\mathbf{E}} \widetilde{W} (\partial_{\widetilde{W}} f)(W)$$

with  $\partial_{\widetilde{W}}$  denoting the directional derivative in direction  $\widetilde{W}$  and  $\widetilde{W}$  being an independent GUE-matrix with expectation  $\widetilde{\mathbf{E}}$ . We now use (98) and (95) for  $G_1 = G(z_1)$  and  $m_1 = m(z_1)$  to obtain

$$(99) \quad \left(1 - \mathcal{O}\left(\frac{N^\xi}{N\eta_*}\right)\right) \left\langle \prod_{i=1}^k (G_i A_i) \right\rangle = m_1 \left\langle A_1 \prod_{i=2}^k (G_i A_i) \right\rangle - m_1 \left\langle \underline{WG_1 A_1} \prod_{i=2}^k (G_i A_i) \right\rangle.$$

Together with

$$\begin{aligned} \left\langle \underline{W \prod_{i=1}^k (G_i A_i)} \right\rangle &= \left\langle W \prod_{i=1}^k (G_i A_i) \right\rangle + \sum_{j=1}^k \widetilde{\mathbf{E}} \left\langle \widetilde{W} \left[ \prod_{i=1}^{j-1} (G_i A_i) \right] G_j \widetilde{W} \prod_{i=j}^k (G_i A_i) \right\rangle \\ &= \left\langle \underline{WG_1 A_1} \prod_{i=2}^k (G_i A_i) \right\rangle + \sum_{j=2}^k \left\langle \left[ \prod_{i=1}^{j-1} (G_i A_i) \right] G_j \right\rangle \left\langle \prod_{i=j}^k (G_i A_i) \right\rangle, \end{aligned}$$

we thus have

$$(100) \quad \begin{aligned} \left(1 - \mathcal{O}\left(\frac{N^\xi}{N\eta_*}\right)\right) \left\langle \prod_{i=1}^k (G_i A_i) \right\rangle &= m_1 \left\langle A_1 \prod_{i=2}^k (G_i A_i) \right\rangle - m_1 \left\langle \underline{W \prod_{i=1}^k (G_i A_i)} \right\rangle \\ &\quad + \sum_{j=2}^k \left\langle \left[ \prod_{i=1}^{j-1} (G_i A_i) \right] G_j \right\rangle \left\langle \prod_{i=j}^k (G_i A_i) \right\rangle. \end{aligned}$$

We now apply the inequality [15], equation (5.35),

$$|\langle XY \rangle| \leq [(X^* X (Y Y^*)^{1/2}) \langle (Y^* Y)^{1/2} \rangle]^{1/2}$$

for arbitrary matrices  $X, Y$  to  $X = \prod_{i=1}^{j-1} (G_i A_i), Y = G_j$  to obtain

$$\left| \left\langle \prod_{i=1}^{j-1} (G_i A_i) G_j \right\rangle \right| \leq \eta_1^{-1/2} \langle (A_{j-1}^* G_{j-1}^* \cdots A_1^* \mathfrak{S} G_1 A_1 \cdots G_{j-1} A_{j-1} | G_j |) \rangle^{1/2}$$

from  $G^* G = (\mathfrak{S} G) / \eta$ . By spectral decomposition we may further estimate with very high probability, for any  $\xi > 0$ ,

$$\begin{aligned} & \langle (A_{j-1}^* G_{j-1}^* \cdots A_1^* \mathfrak{S} G_1 A_1 \cdots G_{j-1} A_{j-1} | G_j |) \rangle \\ &= \left| \frac{1}{N} \sum_a \frac{\langle \mathbf{u}_a, A_{j-1}^* \mathbf{u}_{a_{j-1}} \rangle \cdots \langle \mathbf{u}_{a_2}, A_1^* \mathbf{u}_{a_1} \rangle \langle \mathbf{u}_{a_1}, A_1 \mathbf{u}_{a_2} \rangle \cdots \langle \mathbf{u}_{a_{j-1}}, A_{j-1} \mathbf{u}_{a_j} \rangle}{[(\lambda_{a_{j-1}} - z_{j-1})(\lambda_{a_{j-1}} - \bar{z}_{j-1}) \cdots (\lambda_{a_2} - z_2)(\lambda_{a_2} - \bar{z}_2)] |\lambda_{a_j} - z_j|} \mathfrak{S} \frac{1}{\lambda_{a_1} - z_1} \right| \\ &\lesssim N^{\xi+j-2} \rho(z_1) \end{aligned}$$

from the overlap bound  $|\langle \mathbf{u}_a, A \mathbf{u}_b \rangle| \lesssim N^{\xi-1/2}$ , and where  $\sum_a$  is the summation over the  $2j - 2$  indices  $a_1, a_{\pm 2}, \dots, a_{\pm(j-1)}, a_j$  and conclude

$$(101) \quad \left| \left\langle \prod_{i=1}^{j-1} (G_i A_i) G_j \right\rangle \right| \lesssim N^{\xi+j/2-1} \frac{\sqrt{\rho_1}}{\sqrt{\eta_1}}$$

Similarly, we also have

$$\left| \left\langle \prod_{i=j+1}^k (G_i A_i) \right\rangle \right| \lesssim N^{\xi+(k-j)/2-1},$$

and the claim follows from (100) and the bound

$$(102) \quad \left| \underbrace{W \prod_{i=1}^k (G_i A_i)} \right| \lesssim N^{\xi} N^{(k-3)/2} \frac{\sqrt{\rho^*}}{\sqrt{\eta^*}}$$

on the underlined term in [15], Theorem 4.1, Remark 4.3.  $\square$

### APPENDIX: GREEN FUNCTION COMPARISON

Here, we briefly recall the standard Green function comparison method for eigenvector statistics. The only novelty is that, in addition to the standard entrywise local law,  $|G_{ab}(z)| \lesssim N^{\zeta+\xi}$  for  $\mathfrak{S}z \sim N^{-1-\zeta}$ , we also need an analogous a priori bound for  $(GAG)_{ab}$  that exploits the fact that  $A$  is traceless; see (113) later. Consider the Ornstein–Uhlenbeck flow

$$(103) \quad d\widehat{W}_t = -\frac{1}{2} \widehat{W}_t dt + \frac{d\widehat{B}_t}{\sqrt{N}}, \quad \widehat{W}_0 = W$$

with  $\widehat{B}_t$  a real symmetric Brownian motion. The OU-flow (103) has the effect of adding a small Gaussian component to  $W$  so that, for any fixed  $T$ , we can decompose

$$(104) \quad \widehat{W}_T \stackrel{d}{=} \sqrt{1-cT} \widetilde{W} + \sqrt{cT} U$$

with  $c = c(T) > 0$  a constant very close to one as long as  $T \ll 1$ , and  $U, \widetilde{W}$  being independent GOE/Wigner matrices. Now, let  $W_t$  be the solution of the flow (12) with initial condition  $W_0 = \sqrt{1-cT} \widetilde{W}$  so that

$$(105) \quad W_{cT} \stackrel{d}{=} \widehat{W}_T.$$

LEMMA A.1. *Let  $\widehat{W}_t$  be the solution of (103), and let  $\widehat{\mathbf{u}}_i(t)$  be its eigenvectors. Then, for any smooth test function  $\theta$  of, at most, polynomial growth and any fixed  $\epsilon \in (0, 1/2)$ , there exists an  $\omega = \omega(\theta, \epsilon) > 0$  such that, for any  $i \in [\delta N, (1 - \delta)N]$  (with  $\delta > 0$  from Theorem 2.2) and  $t = N^{-1+\epsilon}$ , it holds that*

$$(106) \quad \mathbf{E} \theta \left( \sqrt{\frac{N}{2\langle A^2 \rangle}} \langle \widehat{\mathbf{u}}_i(t), A\widehat{\mathbf{u}}_i(t) \rangle \right) = \mathbf{E} \theta \left( \sqrt{\frac{N}{2\langle A^2 \rangle}} \langle \widehat{\mathbf{u}}_i(0), A\widehat{\mathbf{u}}_i(0) \rangle \right) + \mathcal{O}(N^{-\omega}).$$

With  $T = N^{-1+\epsilon}$  and  $\theta(x) = x^n$  for some integer  $n \in \mathbf{N}$ , it now follows that

$$(107) \quad \begin{aligned} \mathbf{E} \left[ \sqrt{\frac{N}{2\langle A^2 \rangle}} \langle \mathbf{u}_i, A\mathbf{u}_i \rangle \right]^n &= \mathbf{E} \left[ \sqrt{\frac{N}{2\langle A^2 \rangle}} \langle \widehat{\mathbf{u}}_i(T), A\widehat{\mathbf{u}}_i(T) \rangle \right]^n + \mathcal{O}(N^{-c}) \\ &= \mathbf{E} \left[ \sqrt{\frac{N}{2\langle A^2 \rangle}} \langle \mathbf{u}_i(cT), A\mathbf{u}_i(cT) \rangle \right]^n + \mathcal{O}(N^{-c}) \\ &= \mathbf{1}(n \text{ even})(n - 1)!! + \mathcal{O}(N^{-c}) \end{aligned}$$

for some small  $c = c(n, \epsilon) > 0$ , with  $\mathbf{u}_i, \widehat{\mathbf{u}}_i(t), \mathbf{u}_i(t)$  being the eigenvectors of  $W, \widehat{W}_t, W_t$ , respectively, concluding the proof of Theorem 2.2. Note that in (107) we used Lemma A.1 in the first, (105) in the second, and (27) in the third step, using that in distribution the eigenvectors of  $W_{cT}$  are equal to those of  $\widetilde{W}_{cT/(1-cT)}$  with  $\widetilde{W}_t$  being the solution to the DBM flow with initial condition  $\widetilde{W}_0 = \widetilde{W}$ .

PROOF OF LEMMA A.1. The proof of Lemma A.1 follows from comparing expectations of products of resolvents  $G(z)$  at scales slightly below the eigenvalue spacing, that is, for  $\Im z \sim N^{-1-\zeta}$ . Green function comparison for eigenvectors has been presented in [23] in details and has been used in [11, 12, 26]. Since this is a standard argument, we only give an outline. Let  $W_t$  be the solution of (103), with  $W_0 = W$ , where  $W$  is a Wigner matrix satisfying Assumption 2.1. Here, we dropped the hat compared to the notation used in (103) to make the presentation clearer; that is, we use  $W_t$  instead of  $\widehat{W}_t$ ,  $\mathbf{u}_i(t)$  instead of  $\widehat{\mathbf{u}}_i(t)$ , etc. From now on, by  $G_t = G_t(z)$  we denote the resolvent of  $W_t$ . Note that along the flow (103) the first two moments of  $W$  are preserved.

Due to level repulsion, as in [22], Lemma 5.2, to understand  $\sqrt{N} \langle \mathbf{u}_i, A\mathbf{u}_i \rangle$  it is sufficient to understand functions of  $\sqrt{N} \langle \Im G(z) A \rangle$  with  $\Im z$  slightly below  $N^{-1}$ , that is, the local eigenvalue spacing. In order to prove (106), it is enough to show that

$$(108) \quad \sup_{E \in (-2+\delta, 2-\delta)} |\mathbf{E} \theta(\sqrt{N} \langle \Im G_t(z) A \rangle) - \mathbf{E} \theta(\sqrt{N} \langle \Im G_0(z) A \rangle)| \lesssim N^{-\omega}$$

for  $z = E + i\eta$  for some  $\zeta > 0, \omega > 0$  and all  $\eta \geq N^{-1-\zeta}$ , cf. [3], Section 4, and [11], Appendix A. Define

$$(109) \quad R_t := \theta(\sqrt{N} \langle \Im G_t(z) A \rangle);$$

then, by Itô’s formula we have

$$(110) \quad \mathbf{E} \frac{dR_t}{dt} = \mathbf{E} \left[ -\frac{1}{2} \sum_{\alpha} w_{\alpha}(t) \partial_{\alpha} R_t + \frac{1}{2} \sum_{\alpha, \beta} \kappa_t(\alpha, \beta) \partial_{\alpha} \partial_{\beta} R_t \right],$$

where  $\alpha, \beta \in [N]^2$  are double indices,  $w_{\alpha}(t)$  are the entries of  $W_t$ , and  $\partial_{\alpha} := \partial_{w_{\alpha}}$ . Here,

$$(111) \quad \kappa_t(\alpha_1, \dots, \alpha_l) := \kappa(w_{\alpha_1}(t), \dots, w_{\alpha_l}(t))$$

denotes the joint cumulant of  $w_{\alpha_1}(t), \dots, w_{\alpha_l}(t)$ , with  $l \in \mathbf{N}$ . Note that by (10) it follows that  $|\kappa_t(\alpha_1, \dots, \alpha_l)| \lesssim N^{-l/2}$  uniformly in  $t \geq 0$ . Performing a cumulant expansion in (110) (see [14], equation (25), for more details), we are left with

$$(112) \quad \mathbf{E} \frac{dR_t}{dt} = \sum_{l=3}^R \sum_{\alpha_1, \dots, \alpha_l} \kappa_t(\alpha_1, \dots, \alpha_l) \mathbf{E}[\partial_{\alpha_1} \cdots \partial_{\alpha_l} R_t] + \Omega(R),$$

where  $\Omega(R)$  is an error term, easily seen to be negligible as every additional derivative gains a further factor of  $N^{-1/2}$ . In order to estimate (112), we use  $|(G_t)_{ab}| \leq N^\zeta$  and

$$(113) \quad \begin{aligned} |(G_t(z_1)AG_t(z_2))_{ab}| &= \left| \sum_{ij} \frac{\mathbf{u}_i(a)\langle \mathbf{u}_i, A\mathbf{u}_j \rangle \mathbf{u}_j(b)}{(\lambda_i - z_1)(\lambda_j - z_2)} \right| \\ &\lesssim N^{1/2+\xi} \left( \frac{1}{N} \sum_i \frac{1}{|\lambda_i - z_1|} \right) \left( \frac{1}{N} \sum_i \frac{1}{|\lambda_i - z_2|} \right) \lesssim N^{1/2+\xi+2\zeta} \end{aligned}$$

which holds with very high probability for  $z_1, z_2 \in \{z, \bar{z}\}$ . In (113) we used the eigenvector delocalisation  $\|\mathbf{u}_i\|_\infty \lesssim N^{-1/2+\xi}$  (see [21], or [5] for the optimal bound) and the optimal a priori bound  $|\langle \mathbf{u}_i, A\mathbf{u}_j \rangle| \lesssim N^{-1/2+\xi}$  for traceless  $A$  by [15], Theorem 2.2 (note that this step crucially uses that  $\langle A \rangle = 0$ , the analogous bound for a general  $A$  would be larger by a factor  $\sqrt{N}$ ). We claim that, for any  $l \geq 0$ , it holds that

$$(114) \quad |\partial_{\alpha_1} \dots \partial_{\alpha_l} \sqrt{N} \langle \mathfrak{S} G_t A \rangle| \leq N^{(l+3)(\zeta+\xi)},$$

for any arbitrary small  $\xi > 0$ , with very high probability. Together with

$$\sum_{\alpha_1, \dots, \alpha_l} |\kappa_t(\alpha_1, \dots, \alpha_l)| \lesssim N^{2-l/2},$$

we are then able to estimate  $|\mathbf{E} dR_t/dt|$  by the chain rule to finally obtain (108).

The bound (114) for  $l = 0$  follows immediately from the a priori bound  $\sqrt{N} |\langle G_t A \rangle| \lesssim N^{\xi+\zeta}$ . For  $l = 1$ , the first derivative yields  $|\partial_{ab} \sqrt{N} \langle G_t A \rangle| = N^{-1/2} |(G_t A G_t)_{ba}|$ , and each additional derivative creates a single factor of  $G_t$ , and (114) follows from estimating each factor entrywise by  $|(G_t)_{ab}| \lesssim N^{\xi+\zeta}$  and (113).  $\square$

**Acknowledgments.** L.E. would like to thank Zhigang Bao for many illuminating discussions in an early stage of this research. The authors are also grateful to Paul Bourgade for his comments on the manuscript and the anonymous referee for several useful suggestions.

### REFERENCES

- [1] AGGARWAL, A., LOPATTO, P. and MARCINEK, J. (2021). Eigenvector statistics of Lévy matrices. *Ann. Probab.* **49** 1778–1846. MR4260468 <https://doi.org/10.1214/20-aop1493>
- [2] ANANTHARAMAN, N. and SABRI, M. (2019). Quantum ergodicity on graphs: From spectral to spatial delocalization. *Ann. of Math. (2)* **189** 753–835. MR3961083 <https://doi.org/10.4007/annals.2019.189.3.3>
- [3] BENIGNI, L. (2020). Eigenvectors distribution and quantum unique ergodicity for deformed Wigner matrices. *Ann. Inst. Henri Poincaré Probab. Stat.* **56** 2822–2867. MR4164858 <https://doi.org/10.1214/20-AIHP1060>
- [4] BENIGNI, L. (2021). Fermionic eigenvector moment flow. *Probab. Theory Related Fields* **179** 733–775. MR4242625 <https://doi.org/10.1007/s00440-020-01018-0>
- [5] BENIGNI, L. and LOPATTO, P. (2020). Optimal delocalization for generalized Wigner matrices. Preprint. Available at arXiv:2007.09585.
- [6] BENIGNI, L. and LOPATTO, P. (2021). Fluctuations in local quantum unique ergodicity for generalized Wigner matrices. Preprint. Available at arXiv:2103.12013.

- [7] BLOEMENDAL, A., ERDŐS, L., KNOWLES, A., YAU, H.-T. and YIN, J. (2014). Isotropic local laws for sample covariance and generalized Wigner matrices. *Electron. J. Probab.* **19** no. 33, 53 pp. MR3183577 <https://doi.org/10.1214/ejp.v19-3054>
- [8] BOURGADE, P. (2021). Extreme gaps between eigenvalues of Wigner matrices. *J. Eur. Math. Soc.* **23**. <https://doi.org/10.4171/JEMS/1141>
- [9] BOURGADE, P., ERDŐS, L., YAU, H.-T. and YIN, J. (2016). Fixed energy universality for generalized Wigner matrices. *Comm. Pure Appl. Math.* **69** 1815–1881. MR3541852 <https://doi.org/10.1002/cpa.21624>
- [10] BOURGADE, P., HUANG, J. and YAU, H.-T. (2017). Eigenvector statistics of sparse random matrices. *Electron. J. Probab.* **22** Paper No. 64, 38 pp. MR3690289 <https://doi.org/10.1214/17-EJP81>
- [11] BOURGADE, P. and YAU, H.-T. (2017). The eigenvector moment flow and local quantum unique ergodicity. *Comm. Math. Phys.* **350** 231–278. MR3606475 <https://doi.org/10.1007/s00220-016-2627-6>
- [12] BOURGADE, P., YAU, H.-T. and YIN, J. (2020). Random band matrices in the delocalized phase I: Quantum unique ergodicity and universality. *Comm. Pure Appl. Math.* **73** 1526–1596. MR4156609 <https://doi.org/10.1002/cpa.21895>
- [13] CIPOLLONI, G., ERDŐS, L. and SCHRÖDER, D. (2020). Functional central limit theorems for Wigner matrices. Preprint. Available at [arXiv:2012.13218](https://arxiv.org/abs/2012.13218).
- [14] CIPOLLONI, G., ERDŐS, L. and SCHRÖDER, D. (2021). Edge universality for non-Hermitian random matrices. *Probab. Theory Related Fields* **179** 1–28. MR4221653 <https://doi.org/10.1007/s00440-020-01003-7>
- [15] CIPOLLONI, G., ERDŐS, L. and SCHRÖDER, D. (2021). Eigenstate thermalization hypothesis for Wigner matrices. *Comm. Math. Phys.* **388** 1005–1048. MR4334253 <https://doi.org/10.1007/s00220-021-04239-z>
- [16] COLIN DE VERDIÈRE, Y. (1985). Ergodicité et fonctions propres du laplacien. *Comm. Math. Phys.* **102** 497–502. MR0818831
- [17] D’ALESSIO, L., KAFRI, Y., POLKOVNIKOV, A. and RIGOL, M. (2016). From quantum chaos and eigenstate thermalization to statistical mechanics and thermodynamics. *Adv. Phys.* **65** 239–362. Available at [arXiv:1509.06411](https://arxiv.org/abs/1509.06411). <https://doi.org/10.1080/00018732.2016.1198134>
- [18] DEUTSCH, J. (1991). Quantum statistical mechanics in a closed system. *Phys. Rev. A* **43** 2046–2049. <https://doi.org/10.1103/PhysRevA.43.2046>
- [19] ERDŐS, L., KNOWLES, A., YAU, H.-T. and YIN, J. (2013). The local semicircle law for a general class of random matrices. *Electron. J. Probab.* **18** no. 59, 58 pp. MR3068390 <https://doi.org/10.1214/EJP.v18-2473>
- [20] ERDŐS, L. and YAU, H.-T. (2015). Gap universality of generalized Wigner and  $\beta$ -ensembles. *J. Eur. Math. Soc. (JEMS)* **17** 1927–2036. MR3372074 <https://doi.org/10.4171/JEMS/548>
- [21] ERDŐS, L., YAU, H.-T. and YIN, J. (2012). Rigidity of eigenvalues of generalized Wigner matrices. *Adv. Math.* **229** 1435–1515. MR2871147 <https://doi.org/10.1016/j.aim.2011.12.010>
- [22] KNOWLES, A. and YIN, J. (2013). Eigenvector distribution of Wigner matrices. *Probab. Theory Related Fields* **155** 543–582. MR3034787 <https://doi.org/10.1007/s00440-011-0407-y>
- [23] KNOWLES, A. and YIN, J. (2013). The isotropic semicircle law and deformation of Wigner matrices. *Comm. Pure Appl. Math.* **66** 1663–1750. MR3103909 <https://doi.org/10.1002/cpa.21450>
- [24] LANDON, B., SOSOE, P. and YAU, H.-T. (2019). Fixed energy universality of Dyson Brownian motion. *Adv. Math.* **346** 1137–1332. MR3914908 <https://doi.org/10.1016/j.aim.2019.02.010>
- [25] LIEB, E. H. and LOSS, M. (2001). *Analysis*, 2nd ed. *Graduate Studies in Mathematics* **14**. Amer. Math. Soc., Providence, RI. MR1817225 <https://doi.org/10.1090/gsm/014>
- [26] MARCINEK, J. and YAU, H.-T. (2020). High dimensional normality of noisy eigenvectors. Preprint. Available at [arXiv:2005.08425](https://arxiv.org/abs/2005.08425).
- [27] MARKLOF, J. and RUDNICK, Z. (2000). Quantum unique ergodicity for parabolic maps. *Geom. Funct. Anal.* **10** 1554–1578. MR1810753 <https://doi.org/10.1007/PL00001661>
- [28] O’ROURKE, S., VU, V. and WANG, K. (2016). Eigenvectors of random matrices: A survey. *J. Combin. Theory Ser. A* **144** 361–442. MR3534074 <https://doi.org/10.1016/j.jcta.2016.06.008>
- [29] RUDNICK, Z. and SARNAK, P. (1994). The behaviour of eigenstates of arithmetic hyperbolic manifolds. *Comm. Math. Phys.* **161** 195–213. MR1266075
- [30] ŠNIREL’MAN, A. I. (1974). Ergodic properties of eigenfunctions. *Uspekhi Mat. Nauk* **29** 181–182. MR0402834
- [31] SREDNICKI, M. (1994). Chaos and quantum thermalization. *Phys. Rev. E* **50** 888–901. <https://doi.org/10.1103/PhysRevE.50.888>
- [32] TAO, T. and VU, V. (2012). Random matrices: Universal properties of eigenvectors. *Random Matrices Theory Appl.* **1** 1150001, 27 pp. MR2930379 <https://doi.org/10.1142/S2010326311500018>
- [33] ZELDITCH, S. (1987). Uniform distribution of eigenfunctions on compact hyperbolic surfaces. *Duke Math. J.* **55** 919–941. MR0916129 <https://doi.org/10.1215/S0012-7094-87-05546-3>