

# Random band matrices in the delocalized phase, II: Generalized resolvent estimates

P. Bourgade <i>Courant Institute</i> <i>bourgade@cims.nyu.edu</i>	F. Yang <i>U. of California, Los Angeles</i> <i>fyang75@math.ucla.edu</i>	H.-T. Yau <i>Harvard University</i> <i>htyau@math.harvard.edu</i>	J. Yin <i>U. of California, Los Angeles</i> <i>jyin@math.ucla.edu</i>
---	---	---	---

This is the second part of a three part series of papers. In this paper, we consider a general class of  $N \times N$  random band matrices  $H = (H_{ij})$  whose entries are centered random variables, independent up to a symmetry constraint. We assume that the variances  $\mathbb{E}|H_{ij}|^2$  form a band matrix with typical band width  $1 \ll W \ll N$ . We consider the generalized resolvent of  $H$  defined as  $G(Z) := (H - Z)^{-1}$ , where  $Z$  is a deterministic diagonal matrix such that  $Z_{ij} = (z\mathbb{1}_{1 \leq i \leq W} + \tilde{z}\mathbb{1}_{i > W})\delta_{ij}$ , with two distinct spectral parameters  $z \in \mathbb{C}_+ := \{z \in \mathbb{C} : \text{Im } z > 0\}$  and  $\tilde{z} \in \mathbb{C}_+ \cup \mathbb{R}$ . In this paper, we prove a sharp bound for the local law of the generalized resolvent  $G$  for  $W \gg N^{3/4}$ . This bound is a key input for the proof of delocalization and bulk universality of random band matrices in [2]. Our proof depends on a fluctuations averaging bound on certain averages of polynomials in the resolvent entries, which will be proved in [10].

1	The model and the results.....	1
2	Tools for the proof of Theorem 1.4 .....	6
3	Proof of Theorem 1.4 .....	13
4	Properties of $M$ .....	16

## 1 THE MODEL AND THE RESULTS.

**1.1 The model.** Our goal in this paper is to establish estimates on Green's functions which were used in the proof of delocalization conjecture and bulk universality for random band matrices. All results in this paper apply to both real and complex band matrices. For simplicity of notations, we consider only the real symmetric case. Random band matrices are characterized by the property that the matrix element  $H_{ij}$  becomes negligible if  $\text{dist}(i, j)$  exceeds the band width  $W$ . We shall restrict ourselves to the convention that  $i, j \in \mathbb{Z}_N = \mathbb{Z} \cap (-N/2, N/2]$ , and  $i - j$  is defined modular  $N$ . More precisely, we consider the following matrix ensembles.

**Definition 1.1** (Band matrix  $H_N$  with bandwidth  $W_N$ ). *Let  $H_N$  be an  $N \times N$  matrix with real centered entries  $(H_{ij} : i, j \in \mathbb{Z}_N)$  which are independent up to the condition  $H_{ij} = H_{ji}$ . We say that  $H_N$  is a random band matrix with (typical) bandwidth  $W = W_N$  if*

$$s_{ij} := \mathbb{E}|H_{ij}|^2 = f(i - j) \tag{1.1}$$

---

The work of P.B. is partially supported by the NSF grant DMS#1513587. The work of H.-T. Y. is partially supported by NSF Grant DMS-1606305 and a Simons Investigator award. The work of J.Y. is partially supported by the NSF grant DMS#1552192.

for some non-negative symmetric function  $f : \mathbb{Z}_N \rightarrow \mathbb{R}_+$  satisfying

$$\sum_{x \in \mathbb{Z}_N} f(x) = 1, \quad (1.2)$$

and there exist some (small) positive constant  $c_s$  and (large) positive constant  $C_s$  such that

$$c_s W^{-1} \cdot \mathbf{1}_{|x| \leq W} \leq f(x) \leq C_s W^{-1} \cdot \mathbf{1}_{|x| \leq C_s W}, \quad i, j \in \mathbb{Z}_N. \quad (1.3)$$

The method in this paper also allows to treat cases with exponentially small mass away from the band width (e.g.  $f(x) \leq C_s W^{-1} e^{-c_s |x|^2 / W^2}$ ). We work under the hypothesis (1.3) mainly for simplicity.

We assume that the random variables  $H_{ij}$  have arbitrarily high moments, in the sense that for any fixed  $p \in \mathbb{N}$ , there is a constant  $\mu_p > 0$  such that

$$\max_{i,j} (\mathbb{E}|H_{ij}|^p)^{1/p} \leq \mu_p \text{Var}(H_{ij})^{1/2} \quad (1.4)$$

uniformly in  $N$ .

In this paper, we will *not* need the following moment condition assumed in Part I of this series [2]: there is fixed  $\varepsilon_m > 0$  such that for  $|i - j| \leq W$ ,  $\min_{|i-j| \leq W} (\mathbb{E} \xi_{ij}^4 - (\mathbb{E} \xi_{ij}^3)^2 - 1) \geq N^{-\varepsilon_m}$ , where  $\xi_{ij} := H_{ij}(s_{ij})^{-1/2}$  is the normalized random variable with mean zero and variance one.

All the results in this paper will depend on the parameters  $c_s, C_s$  in (1.3) and  $\mu_p$  in (1.4). But we will not track the dependence on  $c_s, C_s$  and  $\mu_p$  in the proof.

Denote the eigenvalues of  $H_N$  by  $\lambda_1 \leq \dots \leq \lambda_N$ . It is well-known that the empirical spectral measure  $\frac{1}{N} \sum_{k=1}^N \delta_{\lambda_k}$  converges almost surely to the Wigner semicircle law with density

$$\rho_{\text{sc}}(x) = \frac{1}{2\pi} \sqrt{(4 - x^2)_+}.$$

The aim of this paper is to estimate “the generalized resolvent”  $G(z, \tilde{z})$  of  $H_N$  defined by

$$G(z, \tilde{z}) := \left( H_N - \begin{pmatrix} zI_{W \times W} & 0 \\ 0 & \tilde{z}I_{(N-W) \times (N-W)} \end{pmatrix} \right)^{-1}, \quad z, \tilde{z} \in \mathbb{C}^+ \cup \mathbb{R}, \quad (1.5)$$

where  $\mathbb{C}^+$  denotes the upper half complex plane  $\mathbb{C}^+ := \{z \in \mathbb{C} : \text{Im } z > 0\}$ . The generalized resolvent is an important quantity used in Part I of this series [2]. The key point of this generalization, compared with the usual resolvent, is the freedom to choose different  $z$  and  $\tilde{z}$ . To the best of our knowledge, the local law for this type of generalized resolvent has only been studied in the preceding paper [1], where it was assumed that  $W \geq cN$  for some constant  $c > 0$ .

To understand the role of the generalized resolvent, we block-decompose the band matrix  $H_N$  and its eigenvectors as

$$H_N = \begin{pmatrix} A & B^* \\ B & D \end{pmatrix}, \quad \psi_j := \begin{pmatrix} \mathbf{w}_j \\ \mathbf{p}_j \end{pmatrix},$$

where  $A$  is a  $W \times W$  Wigner matrix. From the eigenvector equation  $H\psi_j = \lambda_j \psi_j$ , we get

$$Q_{\lambda_j} \mathbf{w}_j = \lambda_j \mathbf{w}_j, \quad Q_e := A - B^* \frac{1}{D - e} B.$$

Thus  $\mathbf{w}_j$  is an eigenvector of  $Q_e := A - B^*(D - e)^{-1}B$  with eigenvalue  $\lambda_j$  when  $e = \lambda_j$ . A key input to the proof of universality and QUE for random band matrices is an estimate on the Green’s function of  $Q_e$ . Since some eigenvalues of  $D$  can be very close to  $e$ , the matrix  $(D - e)^{-1}$  can be very singular. It is thus very difficult (if possible) to estimate the Green’s function of  $Q_e$  directly. On the other hand, the Green’s function of  $Q_e$  is just the  $W \times W$  minor of the generalized resolvent  $G(z, e)$  of  $H_N$ , which we find to be relatively more doable.

Due to the need in Part I, we will consider generalized resolvent for a general class of band matrices. More precisely, we introduce the following Definition 1.2. Here and throughout the rest of this paper, we will use the notation that for any  $a, b \in \mathbb{Z}$ ,

$$\llbracket a, b \rrbracket := [a, b] \cap \mathbb{Z}.$$

**Definition 1.2** (Definition of  $H_\zeta^{\mathbf{g}}$ ). For any sufficiently small  $\zeta > 0$  and any  $\mathbf{g} = (g_1, g_2, \dots, g_N) \in \mathbb{R}^N$ ,  $H_\zeta$  and  $H_\zeta^{\mathbf{g}}$  will denote  $N \times N$  real symmetric matrices satisfying the following properties. The entries  $(H_\zeta)_{ij}$  are centered and independent up to the symmetry condition, satisfy (1.4), and have variances

$$\mathbb{E}|(H_\zeta)_{ij}|^2 = (s_\zeta)_{ij} := s_{ij} - \frac{\zeta(1 + \delta_{ij})}{W} \mathbf{1}_{i,j \in [1,W]},$$

where  $s_{ij}$ ,  $i, j \in \mathbb{Z}_N$ , satisfy the conditions in Definition 1.1. Then the matrix  $H_\zeta^{\mathbf{g}}$  is defined by

$$(H_\zeta^{\mathbf{g}})_{ij} := (H_\zeta)_{ij} - g_i \delta_{ij}.$$

We denote by  $S_0$  and  $\Sigma$  the matrices with entries  $(S_0)_{ij} = s_{ij}$  and  $\Sigma_{ij} = \frac{(1 + \delta_{ij})}{W} \mathbf{1}_{i,j \in [1,W]}$ , respectively. Then the matrix of variances is

$$S_\zeta := S_0 - \zeta \Sigma, \quad (S_\zeta)_{ij} = (s_\zeta)_{ij}.$$

**1.2 The results.** The generalized resolvent  $G_\zeta^{\mathbf{g}}(z, \tilde{z})$  of  $H_\zeta^{\mathbf{g}}$  is defined similarly as in (1.5) by

$$G_\zeta^{\mathbf{g}}(z, \tilde{z}) := \left( H_\zeta^{\mathbf{g}} - \begin{pmatrix} zI_{W \times W} & 0 \\ 0 & \tilde{z}I_{(N-W) \times (N-W)} \end{pmatrix} \right)^{-1}.$$

Define  $((M_\zeta^{\mathbf{g}})_i(z, \tilde{z}))_{i=1}^N$  as the solution vector to the system of self-consistent equations

$$\left( (M_\zeta^{\mathbf{g}})_i(z, \tilde{z}) \right)^{-1} = -z \mathbf{1}_{i \in [1,W]} - \tilde{z} \mathbf{1}_{i \notin [1,W]} - g_i - \sum_j (s_\zeta)_{ij} (M_\zeta^{\mathbf{g}})_j(z, \tilde{z}), \quad (1.6)$$

for  $z, \tilde{z} \in \mathbb{C}^+ \cup \mathbb{R}$  and  $i \in \mathbb{Z}_N$ , with the constraint that

$$(M_0^{\mathbf{0}})_i(\tilde{z}, \tilde{z}) = m_{\text{sc}}(\tilde{z} + i0^+),$$

where  $m_{\text{sc}}$  denotes the Stieltjes transform of the semicircle law

$$m_{\text{sc}}(z) := \frac{-z + \sqrt{z^2 - 4}}{2}, \quad z \in \mathbb{C}^+. \quad (1.7)$$

(The existence, uniqueness and continuity of the solution is given by Lemma 1.3 below.) For simplicity of notations, we denote by  $M_\zeta^{\mathbf{g}}(z, \tilde{z})$  the diagonal matrix with entries

$$(M_\zeta^{\mathbf{g}})_{ij} := (M_\zeta^{\mathbf{g}})_i \delta_{ij}.$$

We will show that  $M_\zeta^{\mathbf{g}}(z, \tilde{z})$  is the asymptotic limit of the generalized resolvent  $G_\zeta^{\mathbf{g}}(z, \tilde{z})$ . We now list some properties of  $M_\zeta^{\mathbf{g}}$  needed for the proof of local law stated in Theorem 1.4. Its proof is delayed to Section 4.

**Lemma 1.3.** Assume  $|\text{Re } \tilde{z}| \leq 2 - \kappa$  and  $|\tilde{z}| \leq \kappa^{-1}$  for some (small) constant  $\kappa > 0$ . Then there exist constants  $c, C > 0$  such that the following statements hold.

- (Existence and Lipschitz continuity) If

$$\zeta + \|\mathbf{g}\|_\infty + |z - \tilde{z}| \leq c, \quad (1.8)$$

then there exist  $(M_\zeta^{\mathbf{g}})_i(z, \tilde{z})$ ,  $i \in \mathbb{Z}_N$ , which satisfy (1.6) and

$$\max_i \left| (M_\zeta^{\mathbf{g}})_i(z, \tilde{z}) - m_{\text{sc}}(\tilde{z} + i0^+) \right| \leq C (\zeta + \|\mathbf{g}\|_\infty + |z - \tilde{z}|). \quad (1.9)$$

If, in addition, we have  $\zeta' + \|\mathbf{g}'\|_\infty + |z' - \tilde{z}'| \leq c$ , then

$$\max_i \left| (M_{\zeta'}^{\mathbf{g}'})_i(z', \tilde{z}') - (M_\zeta^{\mathbf{g}})_i(z, \tilde{z}) \right| \leq C (\|\mathbf{g} - \mathbf{g}'\|_\infty + |z' - z| + |\tilde{z}' - \tilde{z}| + |\zeta' - \zeta|). \quad (1.10)$$

- (Uniqueness) The solution vector  $((M_\zeta^{\mathbf{g}})_i(z, \tilde{z}))_{i=1}^N$  to (1.6) is unique under (1.8) and the constraint

$$\max_i \left| (M_\zeta^{\mathbf{g}})_i(z, \tilde{z}) - m_{\text{sc}}(\tilde{z} + i0^+) \right| \leq c.$$

We now state our results on the generalized resolvent of  $H_\zeta^{\mathbf{g}}$ . In this paper, we will always use  $\tau$  to denote an arbitrarily small positive constant independent of  $N$ , and  $D$  to denote an arbitrarily large positive constant independent of  $N$ . Define for any matrix  $X$  the max norm

$$\|X\|_{\max} := \max_{i,j} |X_{ij}|.$$

The notations  $\eta_*, \eta^*$  and  $r$  in next theorem were used in Assumptions 2.3 and 2.4 of Part I of this series [2]. Their meanings are not important for this paper and the reader can simply view them as some parameters. In this paper, all the statements hold for sufficiently large  $N$  and we will not repeat it everywhere.

**Theorem 1.4** (Local law). *Define a set of parameters with some constants  $\varepsilon_*, \varepsilon^* > 0$ :*

$$\eta_* := N^{-\varepsilon_*}, \quad \eta^* := N^{-\varepsilon^*}, \quad r := N^{-\varepsilon_* + 3\varepsilon^*}, \quad T := N^{-\varepsilon_* + \varepsilon^*}, \quad 0 < \varepsilon^* \leq \varepsilon_*/20. \quad (1.11)$$

*Fix any  $|e| < 2 - \kappa$  for some constant  $\kappa > 0$ . Then for any deterministic  $z, \zeta, \mathbf{g}$  satisfying*

$$|\operatorname{Re} z - e| \leq r, \quad \eta_* \leq \operatorname{Im} z \leq \eta^*, \quad 0 \leq \zeta \leq T, \quad \|\mathbf{g}\|_\infty \leq W^{-3/4}, \quad (1.12)$$

*and  $W, \varepsilon_*, \varepsilon^*$  satisfying*

$$\log_N W \geq \max \left\{ \frac{6}{7} + \varepsilon^*, \frac{3}{4} + \frac{3}{4}\varepsilon_* + \varepsilon^* \right\}, \quad (1.13)$$

*we have that for any fixed  $\tau > 0$  and  $D > 0$ ,*

$$\mathbb{P} \left( \|G_\zeta^{\mathbf{g}}(z, e) - M_\zeta^{\mathbf{g}}(z, e)\|_{\max} \geq N^\tau \left( \frac{1}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2}}{W} \right) \right) \leq N^{-D}. \quad (1.14)$$

*In fact, the last estimate holds under the weaker assumption*

$$\log_N W \geq \max \left\{ \frac{3}{4} + \varepsilon^*, \frac{1}{2} + \varepsilon_* + \varepsilon^* \right\}. \quad (1.15)$$

We will refer to the first statement, i.e., (1.14) under the assumption (1.13), as the weak form of this theorem, and the statement (1.14) under assumption (1.15) as the strong form. This paper gives a full and self-contained proof for the weak form, which helps the reader understand the basic strategy of our proof. On the other hand, the proof for the strong form is much more involved, and we include a substantial part into a separate paper [10]. Only the strong form of Theorem 1.4 was used in part I of this series [2], where we took  $\log_N W > 3/4$ ,  $\varepsilon_* < 1/4$  and  $\varepsilon^*$  to be a sufficiently small constant.

The main purpose of this part and part III [10] of this series is to prove the above Theorem 1.4. In fact, the bound (1.14) is almost optimal under our setting in the sense that it (at least) gives the correct size of  $\mathbb{E}|(G_\zeta^{\mathbf{g}})_{ij}|^2$  for  $i \neq j$  up to an  $N^\tau$  factor. This sharp bound is very important for the proof of the complete delocalization of eigenvectors and the bulk universality of random band matrices in part I [2]. As explained there, the bound must be of order  $o(W/N)$  to allow the application of the so-called *mean field reduction* method, which was introduced in [1] and is the starting point of this series. Compared with the local law for regular resolvents, the main difficulty in proving the local law for the generalized resolvents is due to the small and even vanishing imaginary part of  $\tilde{z}$ . As a result, some key inputs, such as Ward's identity (see (3.2)) for the regular resolvents estimates are missing. In fact, as discussed before, the case  $\|G(z, \tilde{z})\|_{\max} = \infty$  could occur when  $\tilde{z} = e$  is real. This difficulty has already appeared in the case  $W \geq cN$  in [1], where some "uncertainty principle" was introduced to solve this problem. Unfortunately, this method seems difficult to apply in the  $W \ll N$  case. Instead, in this paper, we shall use a totally different strategy, i.e, the  $T$ -equation method, which was introduced in [4]. Moreover, we have to improve the induction (bootstrap) argument used in [4], as explained below. We remark that the proofs of the weak form and strong form of Theorem 1.4 are completely parallel, except that we will apply a stronger  $T$ -equation estimate (Lemma 2.14) than the one (Lemma 2.8) used in the proof of the weak form. We shall give a simple proof of the weak  $T$ -equation estimate using the standard fluctuation averaging mechanism as in the previous proof of local semicircle law [5, 8]. The proof of the strong  $T$ -equation estimate is based on an improved (and substantially more involved) fluctuation averaging result, whose proof is delayed to part III of this series [10].

**1.3 Sketch of proof.** In the following discussion, for two random variables  $X$  and  $Y$ , we shall use the notation  $X \prec Y$  if for any fixed  $\tau > 0$ ,  $|X| \leq N^\tau |Y|$  with high probability for large enough  $N$ .

We define the  $T$  matrix with entries

$$T_{ij} := \sum_k S_{ik} |G_{kj}|^2, \quad G \equiv G_\zeta^{\mathbf{g}}, \quad S_{ik} \equiv (S_\zeta)_{ik}, \quad (1.16)$$

With a standard self-consistent equation estimate (see Lemma 2.1), one can show that

$$\|G - M\|_{\max}^2 \prec \|T\|_{\max}, \quad M \equiv M_\zeta^{\mathbf{g}}. \quad (1.17)$$

Our proof of Theorem 1.4 is based on an induction argument combined with a self-consistent  $T$ -equation estimate as explained below. We introduce the following notation:

$$\|G\|^2(z, \tilde{z}) := \max_j \sum_{1 \leq i \leq N} |G_{ij}(z, \tilde{z})|^2, \quad \Lambda(z, \tilde{z}) := \|G - M\|_{\max}(z, \tilde{z}). \quad (1.18)$$

Fix  $z$  and  $\operatorname{Re} \tilde{z} = e$ . We perform the induction with respect to the imaginary part of  $\tilde{z}$ . Define a sequence of  $\tilde{z}_n$  such that

$$\operatorname{Im} \tilde{z}_n = N^{-n\varepsilon} \operatorname{Im} z, \quad \operatorname{Re} \tilde{z}_n = e,$$

for small enough constant  $\varepsilon > 0$ . In the  $n = 0$  case with  $\operatorname{Im} \tilde{z}_0 = \operatorname{Im} z$ , using the methods in [5, 8], we can obtain the local law (1.14) for  $G(z, \tilde{z}_0)$ . Suppose we has proved the local law for  $G(z, \tilde{z}_{n-1})$ :

$$\Lambda(z, \tilde{z}_{n-1}) \prec \Phi_{\text{goal}}, \quad \Phi_{\text{goal}} := \frac{1}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2}}{W}. \quad (1.19)$$

Then with  $\operatorname{Im} \tilde{z}_n = N^{-\varepsilon} \operatorname{Im} z_{n-1}$  and a simple (but quite sharp up to an  $N^{2\varepsilon}$  factor)  $L^2$ -estimate, we get a bound on the  $n$ -th level:

$$\|G\|^2(z, \tilde{z}_n) \prec N \tilde{\Phi}^2, \quad \tilde{\Phi}^2 := N^{2\varepsilon} \Phi_{\text{goal}}^2, \quad (1.20)$$

which gives a rough bound  $\Phi^{(0)}$  by the self-consistent equation estimate (1.17):

$$\|T\|_{\max}(z, \tilde{z}_n) \leq \frac{C_s}{W} \|G\|^2(z, \tilde{z}_n) \prec (\Phi^{(0)})^2 \Rightarrow \Lambda(z, \tilde{z}_n) \prec \Phi^{(0)}, \quad \Phi^{(0)} := \sqrt{\frac{N}{W}} \tilde{\Phi}, \quad (1.21)$$

where  $C_s$  is the constant from (1.3). Note that  $\tilde{\Phi}$  is very close to  $\Phi_{\text{goal}}$ , while  $\Phi^{(0)}$  is not. Now with the strong  $T$ -equation estimate (see Lemma 2.14), one can get an improved bound  $(\Phi^{(1)})^2$  on  $T$  as follows:

$$\|T\|_{\max}(z, \tilde{z}_n) \prec (\Phi^{(1)})^2 \Rightarrow \Lambda(z, \tilde{z}_n) \prec \Phi^{(1)}, \quad \Phi^{(1)} := \Phi_{\text{goal}}^2 + \left( \frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2} \right) (\tilde{\Phi}^2 + N^{-1/2}) (\Phi^{(0)})^2, \quad (1.22)$$

where we used (1.17) to get a better bound  $\Lambda(z, \tilde{z}_n) \prec \Phi^{(1)}$ . With (1.15), one can verify that  $\Phi^{(1)} \leq \Phi_{\text{goal}} + N^{-\varepsilon'} \Phi^{(0)}$  for some constant  $\varepsilon' > 0$ . After at most  $l := 1/\varepsilon'$  many iterations with (1.22) and (1.17), i.e.  $\Phi^{(0)} \rightarrow \Phi^{(1)} \rightarrow \dots \rightarrow \Phi^{(l)}$ , we can obtain the local law (1.19) for  $G(z, \tilde{z}_n)$ , which is used as the input for the next induction. The key point of this induction argument is that one has a good  $L^2$ -bound (1.20) inherited from the local law on the upper level, and this  $L^2$ -bound can be used in the  $T$ -equation estimate (1.22) to give an improved bound for  $\Lambda(z, \tilde{z}_n)$  on this level. Finally, after finitely many inductions in  $n$ , we can obtain the local law (1.14) for, say,  $G(z, e + iN^{-10})$ . Then with a continuity argument, we can prove the local law (1.14) for  $G(z, e)$ . In Fig. 1, we illustrate the flow of the induction argument with a diagram.

We remark that the above induction argument is not a continuity argument, as used e.g. in the works [5, 7, 8] on local semicircle law of regular resolvents. The multiplicative steps  $\operatorname{Im} \tilde{z}_n \rightarrow N^{-\varepsilon} \operatorname{Im} \tilde{z}_n$  that we made are far too large for a continuity argument to work. The main reason for choosing this multiplicative step is that the  $T$ -equation estimate can only be applied for  $O(1)$  number of times due to the degrade of the probability set (see Remark 2.9).

The main difficulty of our proof lies in establishing the  $T$ -equation estimate (1.22). The starting point is a self-consistent equation for the  $T$  matrix, i.e. the  $T$ -equation, see (2.14) below. In this paper, we focus on proving the stability of the  $T$ -equation, i.e. bounding  $\|(1 - S|M|^2)^{-1} S\|_{\max}$  in (2.14), where we abbreviate  $S \equiv S_\zeta$ . For regular resolvent of generalized Wigner matrices (i.e.  $\tilde{z} = z$ ,  $\zeta = 0$  and  $\mathbf{g} = \mathbf{0}$ ), we have

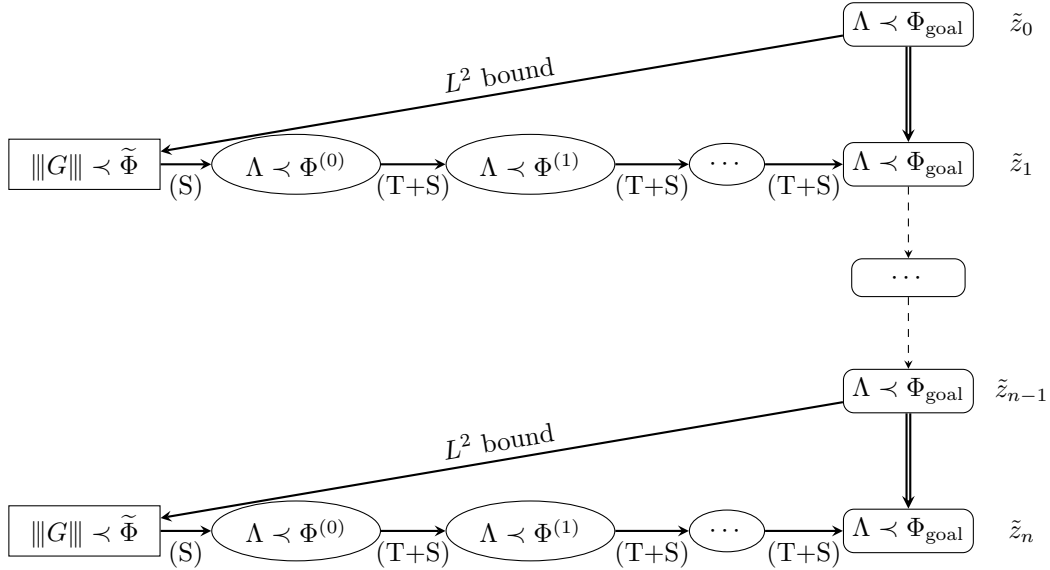


Figure 1: The diagram for the induction argument with respect to  $n$ . At each level  $n - 1$ , we obtain the local law (1.19), which gives the rough bound  $\Phi^{(0)}$  on level  $n$  through (1.20) and (1.21). Applying (1.22) and (1.17) iteratively, one can improve the initial bound  $\Phi^{(0)}$  to the sharp bound  $\Phi_{\text{goal}}$ . In the diagram, (S) stands for an application of the self-consistent equation estimate (1.17), and (T+S) stands for an application of the  $T$ -equation estimate (1.22) followed by a self-consistent equation estimate (1.17).

$|M| \leq 1 - c \text{Im } z$  for some constant  $c > 0$ . However, in our general setting and in particular when  $\text{Im } \tilde{z}$  is small, we actually have  $\|M\|_\infty > 1$  and  $\|S|M|^2\|_{l^\infty \rightarrow l^\infty} > 1$ . Therefore, the usual Taylor expansion approach cannot be used (in fact, it is not even easy to see that 1 is outside the spectrum of  $|M|^2 S$ ). In this paper, we will establish the following bound

$$\|(1 - S|M|^2)^{-1} S\|_{\max} = O\left(\frac{1}{W \text{Im } z} + \frac{N}{W^2}\right).$$

One important component for the proof is the estimate  $\sum_i (|M_i|^2 - 1) \leq -cW \text{Im } z$  for some constant  $c > 0$ . To see this bound is useful, we can intuitively view  $(|M|^2 S)^n$  as an  $n$ -step inhomogeneous random walk on  $\mathbb{Z}_N$  with annihilation, where the average annihilation rate is  $-W \text{Im } z/N$  by the above bound. This shows that we can explore some decay properties of  $(|M|^2 S)^n$  as  $n$  increase, which may give some useful bounds on the Taylor expansion of  $(1 - S|M|^2)^{-1}$ . However, our proof actually will not follow this heuristic argument, see Section 4.

Finally, to finish the proof of the strong version of the  $T$ -equation estimate (Lemma 2.14), we need a fluctuation averaging results for a quantity of the form  $N^{-1} \sum_k \mathcal{E}_k$ , where  $\mathcal{E}_k$ 's are some polynomials of the generalized resolvent entries. The proof involves a new graphical method and we include it in part III of this series [10].

## 2 TOOLS FOR THE PROOF OF THEOREM 1.4

The basic strategy to prove Theorem 1.4 is to apply the self-consistent equation estimate: Lemma 2.1, and the  $T$ -equation estimate: Lemma 2.8 or 2.14, in turns. We collect these results in this section, and use them to prove Theorem 1.4 in next section.

For simplicity, we will often drop the superscripts  $\zeta$  and  $\mathbf{g}$  from our notations. In particular,  $G$  and  $M$  are always understood as  $G_\zeta^{\mathbf{g}}$  and  $M_\zeta^{\mathbf{g}}$ , while  $H$  and  $S$  are understood as  $H_\zeta$  and  $S_\zeta$  in the rest of this paper.

In the proof, for quantities  $A_N$  and  $B_N$ , we will use the notations  $A_N = O(B_N)$  and  $A_N \asymp B_N$  to mean that  $|A_N| \leq C|B_N|$  and  $C^{-1}|B_N| \leq |A_N| \leq C|B_N|$ , respectively, for some constant  $C > 0$ .

**2.1 The self-consistent equation estimate.** The self-consistent equation estimate is the starting point of almost every proof of the local law of the (generalized) resolvents of random matrices. We now state the self-consistent equation estimate for our model.

**Lemma 2.1** (Self-consistent equation estimate). *Suppose that  $|\operatorname{Re} \tilde{z}| \leq 2 - \kappa$  for some constant  $\kappa > 0$ . Then there exists constant  $c_0 > 0$  such that if*

$$\zeta + \|\mathbf{g}\|_\infty + |z - \tilde{z}| \leq c_0,$$

*then the following statement holds. If there exist some fixed  $\delta > 0$  and some deterministic parameter  $\Phi \geq W^{-1/2}$  such that*

$$\|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max} \leq N^{-\delta}, \quad \|T\|_{\max} \leq \Phi^2, \quad (2.1)$$

*in a subset  $\Omega$  of the sample space of the random matrices, then for any fixed  $\tau > 0$  and  $D > 0$ ,*

$$\mathbb{P}(\mathbf{1}_\Omega \|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max} \geq N^\tau \Phi) \leq N^{-D}. \quad (2.2)$$

Note that by the definition of  $T$ -matrix in (1.16), we have

$$\|T\|_{\max} \leq \|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max}^2 + O(W^{-1}).$$

Hence we can always choose  $\Phi = O(N^{-\delta})$  in (2.1). The proof of Lemma 2.1 follows the standard idea of using a vector-level self consistent equation method [5, 8]. In preparation for the proof, we recall the following definition of minors.

**Definition 2.2** (Minors). *For any  $N \times N$  matrix  $A$  and  $\mathbb{T} \subset \{1, \dots, N\}$ , we define the minor of the first kind  $A^{[\mathbb{T}]}$  as the  $(N - |\mathbb{T}|) \times (N - |\mathbb{T}|)$  matrix with*

$$(A^{[\mathbb{T}]})_{ij} := A_{ij}, \quad i, j \notin \mathbb{T}.$$

*For any  $N \times N$  invertible matrix  $B$ , we define the minor of the second kind  $B^{(\mathbb{T})}$  as the  $(N - |\mathbb{T}|) \times (N - |\mathbb{T}|)$  matrix with*

$$(B^{(\mathbb{T})})_{ij} = \left( (B^{-1})^{[\mathbb{T}]} \right)_{ij}^{-1}, \quad i, j \notin \mathbb{T},$$

*whenever  $(B^{-1})^{[\mathbb{T}]}$  is invertible. Note that we keep the names of indices when defining the minors. By definition, for any sets  $\mathbb{U}, \mathbb{T} \subset \{1, \dots, N\}$ , we have*

$$(A^{[\mathbb{T}]})^{[\mathbb{U}]} = A^{[\mathbb{T} \cup \mathbb{U}]}, \quad (B^{(\mathbb{T})})^{(\mathbb{U})} = B^{(\mathbb{T} \cup \mathbb{U})}.$$

*For convenience, we shall also adopt the convention that for  $i \in \mathbb{T}$  or  $j \in \mathbb{T}$ ,*

$$(A^{[\mathbb{T}]})_{ij} = 0, \quad (B^{(\mathbb{T})})_{ij} = 0.$$

*For  $\mathbb{T} = \{a\}$  or  $\mathbb{T} = \{a, b\}$ , we shall abbreviate  $(\{a\}) \equiv (a)$  and  $(\{a, b\}) \equiv (ab)$ .*

**Remark 2.3.** *In previous works, e.g. [6, 8], we have used the notation  $(\cdot)$  for both the minor of the first kind and the minor of the second kind. Here we try to distinguish between  $(\cdot)$  and  $[\cdot]$  in order to be more rigorous.*

The following identities were proved in Lemma 4.2 of [8] and Lemma 6.10 of [6].

**Lemma 2.4** (Resolvent identities). *For an invertible matrix  $B \in \mathbb{C}^{N \times N}$  and  $k \notin \{i, j\}$ , we have*

$$B_{ij} = B_{ij}^{(k)} + \frac{B_{ik} B_{kj}}{B_{kk}}, \quad \frac{1}{B_{ii}} = \frac{1}{B_{ii}^{(k)}} - \frac{B_{ik} B_{ki}}{B_{ii}^{(k)} B_{ii} B_{kk}}, \quad (2.3)$$

*and*

$$\frac{1}{B_{ii}} = (B^{-1})_{ii} - \sum_{k,l}^{(i)} (B^{-1})_{ik} B_{kl}^{(i)} (B^{-1})_{li}. \quad (2.4)$$

*Moreover, for  $i \neq j$  we have*

$$B_{ij} = -B_{ii} \sum_k^{(i)} (B^{-1})_{ik} B_{kj}^{(i)} = -B_{jj} \sum_k^{(j)} B_{ik}^{(j)} (B^{-1})_{kj}. \quad (2.5)$$

*The above equalities are understood to hold whenever the expressions in them make sense.*

Since the  $N^\tau$  factor and the  $N^{-D}$  bound for small probability event appear very often in our proof, we introduce the following notations.

**Definition 2.5.** For any non-negative  $A$ , we denote

$$\mathcal{O}_\tau(A) := \mathcal{O}(N^{\mathcal{O}(\tau)} A).$$

We shall say an event  $\mathcal{E}_N$  holds with high probability (w.h.p.) if for any fixed  $D > 0$ ,

$$\mathbb{P}(\mathcal{E}_N) \geq 1 - N^{-D}$$

for sufficiently large  $N$ . Moreover, we say  $\mathcal{E}_N$  holds with high probability in  $\Omega$  if for any fixed  $D > 0$ ,

$$\mathbb{P}(\Omega \setminus \mathcal{E}_N) \leq N^{-D}$$

for sufficiently large  $N$ .

The following lemma gives standard large deviation bounds that will be used in the proof of Lemma 2.1.

**Lemma 2.6** (Lemma 3.5 of [9]). Let  $(X_i)$  be a family of independent random variables and  $(b_i)$ ,  $(B_{ij})$  be deterministic families of complex numbers, where  $i, j = 1, \dots, N$ . Suppose the entries  $X_i$  satisfy  $\mathbb{E}X_i = 0$ ,  $\mathbb{E}|X_i|^2 = 1$  and the bound(1.4). Then for any fixed  $\tau > 0$ , we have

$$\left| \sum_i b_i X_i \right| \leq N^\tau \left( \sum_i |b_i|^2 \right)^{1/2}, \quad \left| \sum_{i,j} \bar{X}_i B_{ij} X_j \right| \leq N^\tau \left( \sum_{i,j} |B_{ij}|^2 \right)^{1/2},$$

with high probability.

The following lemma provides estimates on the entries of  $(1 - M^2 S)^{-1}$  and  $(1 - S|M|^2)^{-1} S$ . It will be used in the proof of Lemma 2.1 and Theorem 1.4, and its proof is delayed until Section 4.

**Lemma 2.7.** Suppose that the assumptions for the strong form of Theorem 1.4, i.e., (1.11), (1.12) and (1.15), hold. If  $\tilde{z}$  satisfies

$$\operatorname{Re} \tilde{z} = e, \quad 0 \leq \operatorname{Im} \tilde{z} \leq \operatorname{Im} z,$$

then we have for  $M \equiv M_\zeta^{\mathbf{g}}(z, \tilde{z})$  and  $S \equiv S_\zeta$ ,

$$[(1 - M^2 S)^{-1}]_{ij} = \begin{cases} \delta_{ij} + \mathcal{O}(W^{-1}), & \text{if } |i - j| \leq (\log N)^2 W \\ \mathcal{O}(N^{-c \log N}), & \text{if } |i - j| > (\log N)^2 W \end{cases}, \quad (2.6)$$

and

$$\left\| (1 - S|M|^2)^{-1} S \right\|_{\max} = \mathcal{O} \left( \frac{1}{W \operatorname{Im} z} + \frac{N}{W^2} \right). \quad (2.7)$$

Now we can give the proof of Lemma 2.1.

*Proof of Lemma 2.1.* The following proof is fairly standard in random matrix theory and we will omit some details. For simplicity, we drop  $\zeta$  and  $\mathbf{g}$  in superscripts. Using (2.5), we have  $G_{ij} = -G_{ii} \sum_k^{(i)} H_{ik} G_{kj}^{(i)}$  for  $i \neq j$ . Since the elements in  $\{H_{ik}\}_{k=1}^N$  are independent of  $G^{(i)}$ , by the standard large deviations estimates in Lemma 2.6, we have that for any fixed  $\tau > 0$  and  $D > 0$ ,

$$\mathbb{P} \left( |G_{ij}|^2 \leq N^\tau |G_{ii}|^2 \sum_k S_{ik} |G_{kj}^{(i)}|^2 \right) \geq 1 - N^{-D}, \quad i \neq j. \quad (2.8)$$

Since  $G_{ii} \asymp 1$  in  $\Omega$ , (2.8) implies that

$$\mathbb{P} \left( \mathbf{1}_\Omega |G_{ij}|^2 = \mathcal{O}_\tau \left( \sum_k S_{ik} |G_{kj}^{(i)}|^2 \right) \right) \geq 1 - N^{-D}, \quad i \neq j.$$



By (2.3), the definition of  $T$  in (1.16), and the bound for  $T$  in (2.1), we have

$$\sum_k S_{ik} |G_{kj}^{(i)}|^2 \leq 2 \sum_k S_{ik} |G_{kj}|^2 + 2 \sum_k S_{ik} \frac{|G_{ki} G_{ij}|^2}{|G_{ii}|^2} = O(\Phi^2) \quad \text{in } \Omega.$$

Therefore, we obtain (2.2) for the  $i \neq j$  case.

For the diagonal case, we define

$$\mathcal{Z}_i := Q_i \left( \sum_{kl}^{(i)} H_{ik} H_{il} G_{kl}^{(i)} \right) - H_{ii}.$$

Using (2.4), (2.3), the off-diagonal case for (2.2) we just proved, and the standard large deviations estimates in Lemma 2.6, we can get that for any fixed  $\tau > 0$ ,

$$\frac{1}{G_{ii}} = -z \mathbb{1}_{i \in [1, W]} - \tilde{z} \mathbb{1}_{i \notin [1, W]} - g_i - \sum_j S_{ij} G_{jj} - \mathcal{Z}_i + O_\tau(\Phi^2), \quad \text{with } \mathcal{Z}_i = O_\tau(\Phi),$$

holds with high probability in  $\Omega$ . With the definition of  $M_i$  in (1.6), we have

$$G_{ii}^{-1} - M_i^{-1} = - \sum_j S_{ij} (G_{jj} - M_j) + O_\tau(\Phi), \quad \text{w.h.p. in } \Omega,$$

which implies

$$M_i - G_{ii} = - \sum_j M_i^2 S_{ij} (G_{jj} - M_j) + O_\tau(\Phi) + O\left(\max_i |G_{ii} - M_i|^2\right), \quad \text{w.h.p. in } \Omega.$$

We rewrite the above estimate as

$$\sum_j (1 - M^2 S)_{ij} (G_{jj} - M_j) = O_\tau(\Phi) + O\left(\max_i |G_{ii} - M_i|^2\right).$$

Then with (2.6) and the first bound in (2.1), we can get (2.2) for the diagonal entries and complete the proof of Lemma 2.1  $\square$

**2.2 The  $T$ -equation estimate.** A key component for the proof of Theorem 1.4 is the self-consistent equation for the  $T$  variables. It leads to a self-improved bound on  $\|G - M\|_{\max}$ . This kind of approach was also used in [4] to prove a weak type delocalization result for random band matrices. To help the reader understand the proof, we first prove a weak  $T$ -equation estimate, i.e. Lemma 2.8, which will give the weak form of Theorem 1.4. The stronger  $T$ -equation estimate will be stated in Lemma 2.14, and its proof is put in the companion paper [10].

**Lemma 2.8** (Weak  $T$ -equation estimate). *Under the assumptions of Theorem 1.4 (i.e., (1.11), (1.12), (1.15) and the assumption on  $e$ ), the following statements hold provided  $\varepsilon_* > 0$  is a sufficiently small constant. Let  $\tilde{z}$  satisfy*

$$\operatorname{Re} \tilde{z} = e, \quad N^{-10} \leq \operatorname{Im} \tilde{z} \leq \operatorname{Im} z, \quad (2.9)$$

and  $\Phi$  be any deterministic parameter satisfying

$$W^{-1} \leq \Phi^2 \leq N^{-\delta}$$

for some fixed  $\delta > 0$ . Fix some  $z$  and  $\tilde{z}$  (which can depend on  $N$ ). If for any constants  $\tau' > 0$  and  $D' > 0$ ,

$$\mathbb{P}\left(\|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max} \geq N^{\tau'} \Phi\right) \leq N^{-D'}, \quad (2.10)$$

then for any fixed (small)  $\tau > 0$  and (large)  $D > 0$ , we have

$$\mathbb{P}\left(\|T(z, \tilde{z})\|_{\max} \geq N^\tau (\Phi_{\#}^w)^2\right) \leq N^{-D}, \quad (\Phi_{\#}^w)^2 := \left(\frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2}\right) (\Phi^3 + N^{-1}). \quad (2.11)$$

Furthermore, if the parameter  $\Phi$  satisfies

$$\Phi \leq \min \left\{ \frac{W}{N^{1+\varepsilon_*+\varepsilon^*}}, \frac{W^2}{N^{2+\varepsilon^*}} \right\}, \quad (2.12)$$

then for any fixed  $\tau > 0$  and  $D > 0$  we have

$$\|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max} \leq \Phi N^{-\frac{1}{3}\varepsilon^*} + N^\tau \left( \frac{1}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2}}{W} \right) \quad (2.13)$$

with probability at least  $1 - N^{-D}$ .

**Remark 2.9.** *The above statements should be understood as follows. For any small constant  $\tau > 0$  and large constant  $D > 0$ , (2.11) and (2.13) hold if (2.10) holds for some constants  $\tau', D'$  that depend on  $\tau$  and  $D$ . In general, we need to take  $\tau' < \tau$  to be sufficiently small and  $D' > D$  to be sufficiently large. Compared with Lemma 2.1, we lose a much “larger” portion of the probability set. Hence Lemma 2.8 can only be iterated for  $O(1)$  number of times, while Lemma 2.1 can be applied for  $O(N^C)$  times for any fixed  $C > 0$ .*

*Proof of Lemma 2.8.* From the defining equation (1.16) of  $T$ , we add and subtract  $\sum_k S_{ik} |M_k|^2 T_{kj}$  so that

$$T_{ij} = \sum_k S_{ik} |M_k|^2 T_{kj} + \sum_k S_{ik} (|G_{kj}|^2 - |M_k|^2 T_{kj}).$$

Therefore, we have

$$T_{ij} = \sum_k \left[ (1 - S|M|^2)^{-1} S \right]_{ik} (|G_{kj}|^2 - |M_k|^2 T_{kj}). \quad (2.14)$$

Isolating the diagonal terms, we can write the  $T$ -equation as

$$T_{ij} = T_{ij}^0 + \sum_{k \neq j} \left[ (1 - S|M|^2)^{-1} S \right]_{ik} (|G_{kj}|^2 - |M_k|^2 T_{kj}), \quad T_{ij}^0 := \left[ (1 - S|M|^2)^{-1} S \right]_{ij} (|G_{jj}|^2 - |M_j|^2 T_{jj}). \quad (2.15)$$

By the definition of  $T$ , the assumption (2.10) and the estimate (1.9) on  $M_i$ , we can get the simple bounds  $G_{jj} = O(1)$  and  $T_{jj} = O_\tau(\Phi^2)$ . Applying these bounds to the definition of  $T_{ij}^0$ , we get

$$T_{ij}^0 = O \left( \left[ (1 - S|M|^2)^{-1} S \right]_{ij} \right), \quad (2.16)$$

which will be shown to be the main term of  $T_{ij}$  up to an  $N^\tau$  factor. By (2.7) and the condition (1.12) on  $\operatorname{Im} z$ , we have

$$\left[ (1 - S|M|^2)^{-1} S \right]_{ij} = O \left( \frac{1}{W \operatorname{Im} z} + \frac{N}{W^2} \right). \quad (2.17)$$

**Definition 2.10** ( $\mathbb{E}_k, P_k$  and  $Q_k$ ). *We define  $\mathbb{E}_k$  as the partial expectation with respect to the  $k$ -th row and column of  $H$ , i.e.  $\mathbb{E}_k(\cdot) := \mathbb{E}(\cdot | H^{[k]})$ . For simplicity, we will also use the notations*

$$P_k := \mathbb{E}_k, \quad Q_k := 1 - \mathbb{E}_k. \quad (2.18)$$

Using this definition and the bound (2.17), we rewrite the off-diagonal terms in (2.15) into two parts:

$$\begin{aligned} & \sum_{k \neq j} \left[ (1 - S|M|^2)^{-1} S \right]_{ik} (|G_{kj}|^2 - |M_k|^2 T_{kj}) \\ &= \left( \frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2} \right) \left( \sum_{k \neq j} c_k (\mathbb{E}_k |G_{kj}|^2 - |M_k|^2 T_{kj}) + \sum_{k \neq j} c_k Q_k |G_{kj}|^2 \right), \end{aligned} \quad (2.19)$$

where  $c_k$  is a sequence of deterministic numbers satisfying

$$c_k := \left[ (1 - S|M|^2)^{-1} S \right]_{ik} \left( \frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2} \right)^{-1} = O(N^{-1}).$$

The following two lemmas provide estimates for the two parts in (2.19), where Lemma 2.12 is a standard fluctuation averaging lemma.

**Lemma 2.11.** *Suppose that  $b_k, k \in \mathbb{Z}_N$ , are deterministic coefficients satisfying  $\max_k |b_k| = O(N^{-1})$ . Then under the assumptions of Lemma 2.8, we have that for any fixed (small)  $\tau > 0$ ,*

$$\sum_{k \neq j} b_k (\mathbb{E}_k |G_{kj}|^2 - |M_k|^2 T_{kj}) = O_\tau(\Phi^3), \quad j \in \mathbb{Z}_N, \quad (2.20)$$

with high probability.

*Proof.* By (2.5) and (2.10), we have  $-\sum_l^{(k)} H_{kl} G_{lj}^{(k)} = G_{kj}/G_{kk} = O_\tau(\Phi)$  and  $G_{kk} - M_k = O_\tau(\Phi)$  (w.h.p.). Then we can obtain that for  $k \neq j$ ,

$$\mathbb{E}_k |G_{kj}|^2 = \mathbb{E}_k |M_k|^2 \left| \sum_l^{(k)} H_{kl} G_{lj}^{(k)} \right|^2 + O_\tau(\Phi^3) = |M_k|^2 \sum_l^{(k)} s_{kl} \left| G_{lj}^{(k)} \right|^2 + O_\tau(\Phi^3) \quad (2.21)$$

with high probability. Using (2.3), we have

$$G_{lj}^{(k)} = G_{lj} + O_\tau(|G_{lk}| |G_{kj}|) = G_{lj} + O_\tau(\Phi^2), \quad l, j \neq k,$$

with high probability. Inserting it into (2.21) and using the definition (1.16), we can obtain (2.20).  $\square$

**Lemma 2.12.** *Suppose that  $b_k, k \in \mathbb{Z}_N$  are deterministic coefficients satisfying  $\max_k |b_k| = O(N^{-1})$ . Then under the assumptions of Lemma 2.8, we have for any fixed (large)  $p \in 2\mathbb{N}$  and (small)  $\tau > 0$ ,*

$$\mathbb{E} \left| \sum_{k \neq j} b_k Q_k |G_{kj}|^2 \right|^p \leq (N^\tau \Phi^3)^p, \quad j \in \mathbb{Z}_N. \quad (2.22)$$

*Proof.* Our proof follows the arguments in [5, Appendix B]. We consider the decomposition of the space of random variables using  $P_k$  and  $Q_k$  defined in (2.18). It is evident that  $P_k$  and  $Q_k$  are projections,  $P_k + Q_k = 1$ ,  $P_k Q_k = 0$ , and all of these projections commute with each other. For a set  $A \subset \mathbb{Z}_N$ , we denote  $P_A := \prod_{k \in A} P_k$  and  $Q_A := \prod_{k \in A} Q_k$ . Now fix any  $j \in \mathbb{Z}_N$ , we set  $X_k := Q_k |G_{kj}|^2$ . Then for  $p \in 2\mathbb{N}$ , we can write

$$\begin{aligned} \mathbb{E} \left| \sum_{k \neq j} b_k X_k \right|^p &= \sum_{k_1, k_2, \dots, k_p}^* c_{\mathbf{k}} \mathbb{E} \prod_{s=1}^p X_{k_s} = \sum_{\mathbf{k}} c_{\mathbf{k}} \mathbb{E} \prod_{s=1}^p \left( \prod_{r=1}^p (P_{k_r} + Q_{k_r}) X_{k_s} \right) \\ &= \sum_{\mathbf{k}} c_{\mathbf{k}} \sum_{A_1, \dots, A_p \subset [\mathbf{k}]} \mathbb{E} \prod_{s=1}^p (P_{A_s^c} Q_{A_s} X_{k_s}), \end{aligned}$$

where  $\mathbf{k} := (k_1, k_2, \dots, k_p)$ ,  $[\mathbf{k}] := \{k_1, k_2, \dots, k_p\}$ ,  $\sum^*$  means summation with indices not equal to  $j$ , and  $c_{\mathbf{k}}$  are deterministic coefficients satisfying  $c_{\mathbf{k}} = O(N^{-p})$ . Then with the same arguments as in [5] (more specifically, the ones between (B.21)-(B.24)), we see that to conclude (2.22), it suffices to prove that for  $k \in A \subset \mathbb{Z}_N \setminus \{j\}$  and any fixed  $\tau > 0$ ,

$$|Q_A X_k| = O_\tau(\Phi^{|A|+1}) \quad w.h.p. \quad (2.23)$$

We first recall the following simple bound for partial expectations, which is proved in Lemma B.1 of [5]. Given a nonnegative random variable  $X$  and a deterministic control parameter  $\Psi$  such that  $X \leq \Psi$  with high probability. Suppose  $\Psi \geq N^{-C}$  and  $X \leq N^C$  almost surely for some constant  $C > 0$ . Then for any fixed  $\tau > 0$ , we have

$$\max_i P_i X = O_\tau(\Psi) \quad w.h.p. \quad (2.24)$$

In fact, (2.24) follows from Markov's inequality, using high-moments estimates combined with the definition of high probability events in Definition 2.5 and Jensen's inequality for partial expectations. In the application to resolvent entries, the deterministic bound follows from  $\|G\| \leq (\text{Im } \tilde{z})^{-1} \leq N^{10}$  by (2.9).

Now the bound (2.23) in the case  $|A| = 1$  follows from (2.24) directly. For the case  $|A| = n \geq 2$ , we assume without loss of generality that  $j = 1, k = 2$  and  $A = \{2, \dots, n+1\}$ . It suffices to prove that

$$Q_{n+1} \cdots Q_3 |G_{21}|^2 = O_\tau(\Phi^{n+1}). \quad (2.25)$$

Using the identity (2.3), we can write

$$Q_3 |G_{21}|^2 = Q_3 \left( G_{21}^{(3)} + \frac{G_{23}G_{31}}{G_{33}} \right) \overline{\left( G_{21}^{(3)} + \frac{G_{23}G_{31}}{G_{33}} \right)} = Q_3 \left( \frac{\overline{G_{21}^{(3)}} G_{23} G_{31}}{G_{33}} + G_{21}^{(3)} \frac{\overline{G_{23} G_{31}}}{\overline{G_{33}}} + \left| \frac{G_{23} G_{31}}{G_{33}} \right|^2 \right).$$

Note that the leading term  $Q_3 \left| G_{21}^{(3)} \right|^2$  vanishes since  $G_{21}^{(3)}$  is independent of the 3rd row and column of  $H$ , and the rest of the three terms have at least three off-diagonal resolvent entries. We now act  $Q_4$  on these terms, apply (2.3) with  $k = 4$  to each resolvent entry, and multiply everything out. This gives a sum of fractions, where all the entries in the numerator are off-diagonal and all the entries in the denominator are diagonal. Moreover, the leading order terms vanish,

$$Q_4 Q_3 \left( \frac{\overline{G_{21}^{(34)}} G_{23}^{(4)} G_{31}^{(4)}}{G_{33}^{(4)}} + G_{21}^{(34)} \frac{\overline{G_{23}^{(4)} G_{31}^{(4)}}}{\overline{G_{33}^{(4)}}} \right) = 0,$$

and each of the surviving term has at least four off-diagonal resolvent entries. We then continue in this manner, and at each step the number of off-diagonal resolvent entries in the numerator increases at least by one. Finally,  $Q_{n+1} \cdots Q_3 |G_{kj}|^2$  is a sum of fractions where each of them contains at least  $n+1$  off-diagonal entries in the numerator. Together with (2.24), this gives the estimate (2.25), which further proves (2.23).  $\square$

**Remark 2.13.** *Lemma 2.12 asserts that the  $Q_k$  operation yields an improvement by a factor  $\Phi$ . In fact, for the regular resolvents of band matrices, a stronger version of averaging fluctuation results was proved in [3]. We believe that following the methods there, the bounds in Lemma 2.11 and Lemma 2.12 can be improved to*

$$O_\tau(\Phi^4 + W^{-1/2} \Phi^2). \quad (2.26)$$

*In this paper, however, we will skip the discussion on the strategy in [3], since its proof is rather involved, and more importantly, we will prove an even stronger bound, i.e., (2.30) below, in Part III of this series [10]. With (2.26), the  $\Phi_\#^w$  in (2.11) can be improved to*

$$(\Phi_\#^w)^2 = \left( \frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2} \right) (\Phi^4 + W^{-1/2} \Phi^2 + N^{-1}),$$

and the condition (2.12) becomes

$$\Phi^2 \leq \min \left\{ \frac{W}{N^{1+\varepsilon_*+\varepsilon^*}}, \frac{W^2}{N^{2+\varepsilon^*}} \right\}. \quad (2.27)$$

Using this estimate, the conditions (1.13) can be weakened to

$$\log_N W \geq \max \left\{ \frac{4}{5} + \varepsilon^*, \frac{2}{3} + \frac{2}{3} \varepsilon_* + \varepsilon^* \right\}. \quad (2.28)$$

Now we finish the proof of Lemma 2.8. Using (2.19), Lemma 2.11, Lemma 2.12 and Markov's inequality, we can get that

$$\sum_{k \neq j} [(1 - S|M|^2)^{-1} S]_{ik} (|G_{kj}|^2 - |M_k|^2 T_{kj}) = O_\tau \left( \left( \frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2} \right) \Phi^3 \right)$$

with high probability. Note that it only includes the off-diagonal terms, i.e.  $k \neq j$  terms. Now plugging it into the  $T$ -equation (2.15) and using (2.16), we obtain (2.11).

Finally, we need to prove (2.13). Clearly, if (2.12) holds, then  $\Phi \leq N^{-\delta}$  and  $(\Phi_\#^w)^2 \leq N^{-2\delta}$  for some constant  $\delta > 0$ . Thus (2.1) is satisfied, and then (2.13) follows from an application of (2.11) and Lemma 2.1. This completes the proof of Lemma 2.8.  $\square$

The following lemma gives a stronger form of Lemma 2.8. It will be proved in the companion paper [10]. Here we recall the notation in (1.18).

**Lemma 2.14** (Strong  $T$ -equation estimate). *Suppose the assumptions of Theorem 1.4 (i.e., (1.11), (1.12), (1.15) and the assumption on  $e$ ) and (2.9) hold. Let  $\Phi$  and  $\tilde{\Phi}$  be deterministic parameters satisfying*

$$W^{-1} \leq \tilde{\Phi}^2 \leq \Phi^2 \leq \tilde{\Phi} \leq N^{-\delta} \quad (2.29)$$

for some constant  $\delta > 0$ . Fix some  $z$  and  $\tilde{z}$  (which can depend on  $N$ ). If for any constants  $\tau' > 0$  and  $D' > 0$ ,

$$\mathbb{P}\left(\|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max} \geq N^{\tau'} \Phi\right) + \mathbb{P}\left(\|G\|^2(z, \tilde{z}) \geq N^{1+\tau'} \tilde{\Phi}^2\right) \leq N^{-D'},$$

then for any fixed (small)  $\tau > 0$  and (large)  $D > 0$ , we have

$$\mathbb{P}\left(\|T(z, \tilde{z})\|_{\max} \geq N^\tau \Phi_{\#}^2\right) \leq N^{-D}, \quad \Phi_{\#}^2 := \left(\frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2}\right) \left(\Phi^2 \tilde{\Phi}^2 + \Phi^2 N^{-1/2} + N^{-1}\right). \quad (2.30)$$

Furthermore, if the parameter  $\tilde{\Phi}$  satisfies

$$\tilde{\Phi}^2 \leq \min\left\{\frac{W}{N^{1+\varepsilon_*+\varepsilon^*}}, \frac{W^2}{N^{2+\varepsilon^*}}\right\}, \quad (2.31)$$

then for any fixed  $\tau > 0$  and  $D > 0$  we have

$$\|G(z, \tilde{z}) - M(z, \tilde{z})\|_{\max} \leq \Phi N^{-\frac{1}{3}\varepsilon^*} + N^\tau \left(\frac{1}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2}}{W}\right) \quad (2.32)$$

with probability at least  $1 - N^{-D}$ .

The Remark 2.9 also applies to this lemma. Note that (2.13) or (2.32) gives a self-improved bound on  $\|G - M\|_{\max}$ , which explains how we can improve the estimate on  $G$  (from  $\Phi$  to  $\Phi_{\#}$ ) via  $T$  equations. As long as we have an initial estimate such that (2.12) or (2.31) holds, we can then iterate the proof and improve the estimate on  $G$  to  $\Phi_{\text{goal}} = \left(\frac{1}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2}}{W}\right)$  in (1.14).

*Proof of Lemma 2.14.* See the proof of Theorem 2.7 in part III of this series [10].  $\square$

### 3 PROOF OF THEOREM 1.4

Fix a parameter  $0 < \varepsilon_0 < \varepsilon_*/5$ . We define

$$\tilde{z}_n := \operatorname{Re} \tilde{z} + iN^{-n\varepsilon_0} \operatorname{Im} z,$$

so that  $\operatorname{Im} \tilde{z}_{n+1} = N^{-\varepsilon_0} \operatorname{Im} \tilde{z}_n$ . The basic idea in proving Theorem 1.4 is to use mathematical induction on  $n \in \mathbb{N}$ .

The proofs of the weak form and strong form of Theorem 1.4 are completely parallel. In the following proof, we will only remark on the minor differences between them.

**Step 0:** The special case with  $\tilde{z} = z$  and  $\zeta = 0$ ,  $\mathbf{g} = \mathbf{0}$  (i.e.  $G(H, z)$  is the ordinary resolvent of a generalized Wigner matrix) was proved in [5]. The proof given there can be carried over to our case without changes under the assumptions of Theorem 1.4 when  $\tilde{z} = z$  and  $\operatorname{Im} z \geq W^{-1+\delta}$  for some fixed  $\delta > 0$ .

This gives that

$$\mathbb{P}\left(\|G(z, z) - M(z, z)\|_{\max} \geq \frac{N^\tau}{\sqrt{W \operatorname{Im} z}}\right) \leq N^{-D},$$

for any fixed  $\tau > 0$ . This bound is clearly stronger than the one in (1.14).

**Step 1:** Consider the case  $n = 0$ , i.e.,  $G(z, \tilde{z}_0)$ , where we have

$$\operatorname{Re} \tilde{z}_0 = \operatorname{Re} \tilde{z}, \quad \operatorname{Im} \tilde{z}_0 = \operatorname{Im} z.$$

We claim that for any  $w, \tilde{w} \in \mathbb{C}_+$ ,

$$\|G(w, \tilde{w})\|_{L^2 \rightarrow L^2} \leq \frac{1}{\min(\operatorname{Im} w, \operatorname{Im} \tilde{w})}. \quad (3.1)$$

To prove it, we first assume that  $\operatorname{Im} w = a + \operatorname{Im} \tilde{w}$  with  $a \geq 0$ . We write

$$G(w, \tilde{w}) = (A - iaJ - i \operatorname{Im} \tilde{w})^{-1}, \quad J_{kl} = \mathbf{1}_{k \in [1, W]} \delta_{kl},$$

where  $A$  is a symmetric matrix. Then

$$(A - iaJ - i \operatorname{Im} \tilde{w})^* (A - iaJ - i \operatorname{Im} \tilde{w}) = (A - iaJ)^* (A - iaJ) + 2a(\operatorname{Im} \tilde{w})J + (\operatorname{Im} \tilde{w})^2 \geq (\operatorname{Im} \tilde{w})^2.$$

Obviously, we have a similar estimate with  $\operatorname{Im} \tilde{w}$  replaced by  $\operatorname{Im} w$  when  $\operatorname{Im} w \leq \operatorname{Im} \tilde{w}$ . This proves the claim (3.1).

Now by the definition of  $T$  and (1.3), we know

$$|T_{ij}(z, \tilde{z}_0)| \leq \frac{C_s}{W} \sum_k |G_{kj}(z, \tilde{z}_0)|^2 = \frac{C_s \operatorname{Im} G_{jj}(z, \tilde{z}_0)}{W \operatorname{Im} z},$$

where in the second step we used the so-called Ward identity that for any symmetric matrix  $A$  and  $\eta > 0$ ,

$$\sum_k |R_{kj}(A, i\eta)|^2 = \frac{\operatorname{Im} R_{jj}(A, i\eta)}{\eta}, \quad R(A, i\eta) := (A - i\eta)^{-1}. \quad (3.2)$$

Obviously, the same argument gives that

$$\|T(z, \tilde{z}_0(t))\|_{\max} \leq \frac{C_s \max_j \operatorname{Im} G_{jj}(z, \tilde{z}_0(t))}{W \operatorname{Im} z}, \quad \tilde{z}_0(t) := (1-t)z + t\tilde{z}_0, \quad t \in [0, 1]. \quad (3.3)$$

Now we claim that for any small enough  $\tau > 0$ ,

$$\sup_{s \in [0, 1]} \mathbb{P} \left( \|G(z, \tilde{z}_0(t)) - M(z, \tilde{z}_0(t))\|_{\max} \geq \frac{N^\tau}{\sqrt{W \operatorname{Im} z}} \right) \leq N^{-D}. \quad (3.4)$$

To prove (3.4), we first note that for any  $w, w' \in \mathbb{C}$ ,

$$G(z, w) = G(z, w') + G(z, w)(w - w') \tilde{J} G(z, w'), \quad \tilde{J}_{kl} = \mathbf{1}_{k \notin [1, W]} \delta_{kl}. \quad (3.5)$$

This implies that

$$\|\partial_{\tilde{z}} G(z, \tilde{z})\|_{\max} \leq \sqrt{N} \|G(z, \tilde{z})\|_{L^2 \rightarrow L^2} \|G(z, \tilde{z})\|_{\max} \leq \frac{\sqrt{N}}{\min(\operatorname{Im} z, \operatorname{Im} \tilde{z})} \|G\|_{\max}.$$

In particular, in this step we have

$$\|\partial_s G(z, \tilde{z}_0(t))\|_{\max} \leq CN^{1/2+\varepsilon_*} |z - \tilde{z}_0| \|G(z, \tilde{z}_0(t))\|_{\max}. \quad (3.6)$$

This provides some continuity estimate on  $G(z, \tilde{z}_0(t))$ , which shows that (3.4) can be obtained from the following estimate:

$$\max_{k \in [0, N^5]} \mathbb{P} \left( \|G(z, \tilde{z}_0(kN^{-5})) - M(z, \tilde{z}_0(kN^{-5}))\|_{\max} \geq \frac{N^\tau}{\sqrt{W \operatorname{Im} z}} \right) \leq N^{-D}. \quad (3.7)$$

From Step 0, this estimate holds for  $k = 0$ . By induction, we assume that (3.7) holds for  $k = k_0$ . Then using (3.6) and (1.10), we know that the first estimate of (2.1) holds for  $G(z, \tilde{z}_0(t))$  with  $t = (k_0 + 1)N^{-5}$ .

Then by (3.3) and applying Lemma 2.1, we obtain (3.7) for  $k = k_0 + 1$ . This completes the proof of (3.7) and (3.4). Note that the estimate (3.4) applied to  $G(z, \tilde{z}_0(1))$  is the result we want for this step

**Step 2:** Suppose that for some  $n \in \mathbb{N}$  with  $\text{Im } \tilde{z}_n \geq N^{-10}$ , (1.14) holds for  $G(z, \tilde{z}_n)$  and  $M(z, \tilde{z}_n)$  for any large  $D > 0$ . We first prove the following estimate for  $G(z, \tilde{z}_{n+1}) - M(z, \tilde{z}_{n+1})$ , which is weaker than (1.14):

$$\mathbb{P} \left( \|G(z, \tilde{z}_{n+1}) - M(z, \tilde{z}_{n+1})\|_{\max} \geq N^\tau \left( \frac{N^{1/2+\varepsilon_0}}{W\sqrt{\text{Im } z}} + \frac{N^{1+\varepsilon_0}}{W^{3/2}} \right) \right) \leq N^{-D} \quad (3.8)$$

for any fixed  $\tau > 0$ .

For any  $w, w' \in \mathbb{C}^+$  satisfying

$$\text{Re } w = \text{Re } w', \quad N^{-\varepsilon_0} \text{Im } w' \leq \text{Im } w \leq \text{Im } w', \quad (3.9)$$

using (3.9) and (3.5), we have

$$\begin{aligned} \sum_i |G_{ij}(z, w)|^2 &\leq 2 \left( 1 + |w - w'|^2 \|G(z, w)\|_{L^2 \rightarrow L^2}^2 \right) \sum_i |G_{ij}(z, w')|^2 \\ &\leq 2 \left( 1 + \frac{(\text{Im } w')^2}{(\text{Im } w)^2} \right) \sum_i |G_{ij}(z, w')|^2 \leq 3N^{2\varepsilon_0} \|G(z, w')\|^2, \end{aligned}$$

where we have used (3.1) to bound  $\|G(z, w)\|_{L^2 \rightarrow L^2}^2$ . We apply this inequality with  $w' = \tilde{z}_n$  and  $w$  satisfying (3.9). Using (1.14) and the definition (1.18), we can bound  $\|G(z, \tilde{z}_n)\|^2$  as

$$\sup_{\substack{\text{Re } w = \text{Re } \tilde{z}_n, \\ \text{Im } \tilde{z}_{n+1} \leq \text{Im } w \leq \text{Im } \tilde{z}_n}} \|T(z, w)\|_{\max} \leq \sup_{\substack{\text{Re } w = \text{Re } \tilde{z}_n, \\ \text{Im } \tilde{z}_{n+1} \leq \text{Im } w \leq \text{Im } \tilde{z}_n}} \frac{C}{W} \|G(z, w)\|^2 = O_\tau \left( \frac{N^{1+2\varepsilon_0}}{W^2 \text{Im } z} + \frac{N^{2+2\varepsilon_0}}{W^3} \right) \quad (3.10)$$

with high probability for any fixed  $\tau > 0$ .

We now consider interpolation between  $\tilde{z}_n$  and  $\tilde{z}_{n+1}$ :

$$\tilde{z}_{n,m} = \tilde{z}_n - i(\text{Im } \tilde{z}_n - \text{Im } \tilde{z}_{n+1})mN^{-50}, \quad m \in \llbracket 0, N^{50} \rrbracket.$$

We would like to use Lemma 2.1 and induction to prove that (3.8) holds for  $G(z, \tilde{z}_{n,m}) - M(z, \tilde{z}_{n,m})$  for all  $m$ . First, we know (3.8) holds for  $G(z, \tilde{z}_n)$ . Then suppose (3.8) holds for  $G(z, \tilde{z}_{n,j})$  for all  $j \leq m-1$ . We now verify that (2.1) holds for  $G(z, \tilde{z}_{n,m})$  with  $\Phi^2 = N^\tau \Phi_0^2$  for any fixed  $\tau > 0$ , where

$$\Phi_0^2 := \frac{N^{1+2\varepsilon_0}}{W^2 \text{Im } z} + \frac{N^{2+2\varepsilon_0}}{W^3}.$$

To this end, we note that (3.10) already verifies the bound on  $\|T(z, \tilde{z}_{n,m})\|_{\max}$  in (2.1) for all  $m \in \llbracket 0, N^{50} \rrbracket$ . By using  $\|\partial_{\bar{z}} G\|_{\max} \leq N \|G\|_{\max}^2$  (which follows from (3.5)), (1.10),  $|\tilde{z}_{n,m-1} - \tilde{z}_{n,m}| \leq N^{-50}$ , and (3.10) (to bound  $\|G\|_{\max}^2$  by  $\|G\|^2$ ), we note that for sufficiently small constant  $\delta > 0$ ,

$$\|G(z, \tilde{z}_{n,m-1}) - M(z, \tilde{z}_{n,m-1})\|_{\max} \leq N^{-2\delta} \implies \|G(z, \tilde{z}_{n,m}) - M(z, \tilde{z}_{n,m})\|_{\max} \leq N^{-\delta}.$$

This proves the first bound in (2.1) for  $G(z, \tilde{z}_{n,m})$ . Then Lemma 2.1 asserts that (2.2) holds for  $G(z, \tilde{z}_{n,m})$  with  $N^\tau \Phi_0$  for any fixed  $\tau > 0$ . This proves (3.8) (i.e. the  $m = N^{50}$  case) by induction.

**Step 3:** Suppose that for some  $n \in \mathbb{N}$  with  $\text{Im } \tilde{z}_n \geq N^{-10}$ , (1.14) holds for  $G(z, \tilde{z}_n)$  and  $M(z, \tilde{z}_n)$  for any large  $D > 0$ . We have proved that (3.8) and (3.10) hold for  $G(z, \tilde{z}_{n+1})$ . We now apply Lemma 2.8 to prove the weak form of Theorem 1.4. First, the condition (2.10) holds with  $\Phi = \frac{N^{1/2+\varepsilon_0}}{W\sqrt{\text{Im } z}} + \frac{N^{1+\varepsilon_0}}{W^{3/2}}$ . In order for the condition (2.12) to hold, we need

$$\frac{N^{1/2+\varepsilon_0}}{W\sqrt{\text{Im } z}} + \frac{N^{1+\varepsilon_0}}{W^{3/2}} \leq \min \left\{ \frac{W}{N^{1+\varepsilon_*+\varepsilon^*}}, \frac{W^2}{N^{2+\varepsilon^*}} \right\}, \quad (3.11)$$

which is satisfied if

$$W \geq 2 \max \left( N^{\frac{6}{7} + \frac{2}{7}\varepsilon_0 + \frac{2}{7}\varepsilon^*}, N^{\frac{3}{4} + \frac{3}{4}\varepsilon_* + \frac{1}{2}\varepsilon_0 + \frac{1}{2}\varepsilon^*} \right).$$

If we take  $\varepsilon_0 < \varepsilon^*$ , (2.10) implies (2.13) under the condition (1.13). We then apply Lemma 2.8 again, and after at most  $3/\varepsilon^*$  iterations we obtain that

$$\|G(z, \tilde{z}_{n+1}) - M(z, \tilde{z}_{n+1})\|_{\max} \leq N^\tau \left( \frac{1}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2}}{W} \right). \quad (3.12)$$

By induction on  $n$  (with the number of inductions  $\leq 10/\varepsilon_0$ ), the main estimate (3.12) for  $G(z, \tilde{z}_n)$  holds for all  $n$  as long as  $\operatorname{Im} \tilde{z}_n \geq N^{-10}$ .

Similarly, we can apply Lemma 2.14 to prove the strong form of Theorem 1.4. As in the previous argument, (3.8) and (3.10) hold for  $G(z, \tilde{z}_{n+1})$  assuming (1.14) for  $G(z, \tilde{z}_n)$  and  $\operatorname{Im} \tilde{z}_n \geq N^{-10}$ . Therefore, we can choose  $\Phi$  and  $\tilde{\Phi}$  as

$$\Phi = \frac{N^{1/2+\varepsilon_0}}{W\sqrt{\operatorname{Im} z}} + \frac{N^{1+\varepsilon_0}}{W^{3/2}}, \quad \tilde{\Phi} = \frac{N^{\varepsilon_0}}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2+\varepsilon_0}}{W},$$

where the choice of  $\tilde{\Phi}$  follows from using (3.10). It is easy to see that (2.29) holds. In order to apply Lemma 2.14, we need (2.31), i.e.,

$$\left( \frac{N^{\varepsilon_0}}{\sqrt{W \operatorname{Im} z}} + \frac{N^{1/2+\varepsilon_0}}{W} \right)^2 \leq \min \left\{ \frac{W}{N^{1+\varepsilon_*+\varepsilon^*}}, \frac{W^2}{N^{2+\varepsilon^*}} \right\},$$

which is satisfied if

$$W \geq 2 \max \left( N^{\frac{3}{4}+\frac{1}{2}\varepsilon_0+\frac{1}{4}\varepsilon^*}, N^{\frac{1}{2}+\varepsilon_*+\varepsilon_0+\frac{1}{2}\varepsilon^*} \right).$$

Clearly, the assumption (1.15) guarantees this condition if we choose  $\varepsilon_0 < \varepsilon^*/2$ . Again, we can apply Lemma 2.14 iteratively until we get (3.12) for  $G(z, \tilde{z}_{n+1})$ . The rest of the proof for the strong form of Theorem 1.4 is the same as the proof for the weak form.

**Step 4:** We now prove (1.14) for  $G(z, \tilde{z})$  with  $\operatorname{Im} \tilde{z} = 0$  by using continuity from the estimate for  $G(z, \tilde{z})$  with  $\operatorname{Im} \tilde{z} = N^{-10}$  established in **Step 3**. It is easy to see that

$$\partial_{\tilde{z}} \|G(z, \tilde{z})\|_{\max} \leq \|\partial_{\tilde{z}} G(z, \tilde{z})\|_{\max} \leq N \|G(z, \tilde{z})\|_{\max}^2. \quad (3.13)$$

With (3.13) and using (3.12) for  $G(z, \operatorname{Re} \tilde{z} + iN^{-10})$ , we can obtain that

$$\sup_{0 \leq \eta \leq N^{-10}} \|G(z, \operatorname{Re} \tilde{z} + i\eta)\|_{\max} = O(1), \quad w.h.p.$$

Then using (1.10), (3.5) and (3.12) for  $G(z, \operatorname{Re} \tilde{z} + iN^{-10})$ , we obtain that (1.14) holds for  $G(z, \operatorname{Re} \tilde{z})$ .

**Remark 3.1.** *If we use the bound in Remark 2.13 and the condition (2.27) instead of (2.12), then the restriction (3.11) becomes*

$$\left( \frac{N^{1/2+\varepsilon_0}}{W\sqrt{\operatorname{Im} z}} + \frac{N^{1+\varepsilon_0}}{W^{3/2}} \right)^2 \leq \min \left\{ \frac{W}{N^{1+\varepsilon_*+\varepsilon^*}}, \frac{W^2}{N^{2+\varepsilon^*}} \right\}$$

*which gives restriction in (2.28). So we get a result in between the weak and strong forms of Theorem 1.4.*

## 4 PROPERTIES OF $M$

The main goal of this section is to derive some deterministic estimates related to  $(M_\zeta^{\mathbf{g}})_i$ ,  $i \in \mathbb{Z}_N$ . In particular, we will finish the proof of Lemma 1.3 and Lemma 2.7.



**4.1 The stability.** The system of self-consistent equations (1.6) is a perturbation of the standard self-consistent equation

$$m_{\text{sc}}^{-1} = -\tilde{z} - m_{\text{sc}}$$

for  $m_{\text{sc}}(\tilde{z})$ . Thus our basic strategy is to use the standard perturbation theory (see (4.13) below) combined with a stability estimate for the self-consistent equation (i.e. the operator bound (4.4)). We first recall the following elementary properties of  $m_{\text{sc}}$ , which can be proved directly using (1.7).

**Lemma 4.1.** *We have for all  $z = E + i\eta$  with  $\eta > 0$  that*

$$|m_{\text{sc}}(z)| = |m_{\text{sc}}(z) + z|^{-1} \leq 1.$$

Furthermore, there is a constant  $c > 0$  such that for  $E \in [-10, 10]$  and  $\eta \in (0, 10]$  we have

$$c \leq |m_{\text{sc}}(z)| \leq 1 - c\eta, \quad (4.1)$$

$$|\partial_z m_{\text{sc}}(z)| \leq c^{-1}(\kappa + \eta)^{-1/2}, \quad (4.2)$$

$$|1 - m_{\text{sc}}^2(z)| \asymp \sqrt{\kappa + \eta},$$

as well as

$$\text{Im } m_{\text{sc}}(z) \asymp \begin{cases} \sqrt{\kappa + \eta} & \text{if } |E| \leq 2 \\ \frac{\eta}{\sqrt{\kappa + \eta}} & \text{if } |E| \geq 2 \end{cases},$$

where  $\kappa := ||E| - 2|$  denotes the distance of  $E$  to the spectral edges.

The following lemma will be used in the proof of Lemma 1.3 and Lemma 2.7. Recall that  $S_0$  is the matrix with entries  $s_{ij}$ , which is defined in Definition 1.2.

**Lemma 4.2.** *Assume  $|\text{Re } \tilde{z}| \leq 2 - \kappa$  for some constant  $\kappa > 0$  and denote  $m = m_{\text{sc}}(\tilde{z} + i0^+)$ . Then for any fixed  $\tau > 0$ , there exist constants  $c_1, C_1 > 0$  such that*

$$\left\| \left( \frac{m^2 S_0 + \tau}{1 + \tau} \right)^2 \right\|_{L^\infty \rightarrow L^\infty} < 1 - c_1. \quad (4.3)$$

Furthermore,

$$\|(1 - m^2 S_0)^{-1}\|_{L^\infty \rightarrow L^\infty} \leq C_1. \quad (4.4)$$

*Proof.* For some small constant  $\tau > 0$  we write

$$(1 - m^2 S_0)^{-1} = \frac{1}{1 + \tau} \sum_{k=0}^{\infty} \left( \frac{m^2 S_0 + \tau}{1 + \tau} \right)^k. \quad (4.5)$$

Assuming (4.3), we get that

$$\|(1 - m^2 S_0)^{-1}\|_{L^\infty \rightarrow L^\infty} \leq \frac{1}{1 + \tau} \left( 1 + \left\| \frac{m^2 S_0 + \tau}{1 + \tau} \right\|_{L^\infty \rightarrow L^\infty} \right) \sum_{j=0}^{\infty} \left\| \frac{m^2 S_0 + \tau}{1 + \tau} \right\|_{L^\infty \rightarrow L^\infty}^{2j} \leq C_1,$$

which proves (4.4).

We now prove (4.3). Suppose that there is a vector  $\mathbf{v} \in \mathbb{C}^N$  so that  $\|\mathbf{v}\|_\infty = 1$  and

$$\left| \frac{[(m^2 S_0 + \tau)^2 \mathbf{v}]_i}{(1 + \tau)^2} \right| = 1 - \varepsilon$$

for some  $i \in \mathbb{Z}_N$  and  $\varepsilon \equiv \varepsilon_N \rightarrow 0^+$ . Hence

$$(1 + 2\tau + \tau^2)(1 - \varepsilon) = |m^4 b + 2\tau m^2 a + \tau^2 v_i| \leq |b| + 2\tau|a| + \tau^2|v_i| \leq 1 + 2\tau + \tau^2, \quad (4.6)$$

where  $a := (S_0 \mathbf{v})_i$ ,  $b := (S_0^2 \mathbf{v})_i$  and we have used the bounds  $|m| \leq 1$ ,  $|a| \leq 1$  and  $|b| \leq 1$  (since  $\|S_0\|_{L^\infty \rightarrow L^\infty} = 1$ ). It will be clear that the  $|m| = 1$  case is most difficult and we will assume this condition in the following proof. Moreover, we assume with loss of generality that  $v_i > 0$  (by changing the global phase of  $\mathbf{v}$ ). Now  $m$ ,

$a$  and  $b$  are complex numbers, and the inequality (4.6) implies that  $m^4b$ ,  $m^2a$  and  $v_i$  have almost the same phases. Since  $|v_i| \leq 1$ ,  $|b| \leq 1$  and  $|a| \leq 1$ , (4.6) implies that for some constant  $C > 0$  independent of  $\varepsilon$ ,

$$v_i \geq 1 - C\varepsilon, \quad |b - m^{-4}| \leq C\varepsilon, \quad |a - m^{-2}| \leq C\varepsilon. \quad (4.7)$$

Since  $m$  is a unit modulus complex number with imaginary part of order 1, we have that  $\delta := |m^{-2} - m^{-4}|$  is a number of order 1 and

$$|a - b| > \delta/2.$$

Fix the index  $i$  and denote  $c_j := (S_0)_{ij}$ ,  $d_j := (S_0^2)_{ij}$ . Then  $\sum_j c_j = 1 = \sum_j d_j$ . Hence (4.7) implies

$$1 - O(\varepsilon) = \operatorname{Re}(a\bar{a}) = \sum_j c_j \operatorname{Re}(v_j\bar{a}), \quad 1 - O(\varepsilon) = \operatorname{Re}(b\bar{b}) = \sum_j d_j \operatorname{Re}(v_j\bar{b}),$$

where  $O(\varepsilon)$  denote a positive number bounded by  $C\varepsilon$  for some constant  $C > 0$  independent of  $\varepsilon$ . For any  $0 < r < 1$ , denote by  $A_r := \{j : \operatorname{Re}(v_j\bar{a}) \geq 1 - r\}$  and let  $\alpha_r := \sum_{j \in A_r} c_j$ . Then we have

$$\alpha_r \geq \sum_{j \in A_r} c_j \operatorname{Re}(v_j\bar{a}) = 1 - O(\varepsilon) - \sum_{j \notin A_r} c_j \operatorname{Re}(v_j\bar{a}) \geq 1 - O(\varepsilon) - (1 - \alpha_r)(1 - r) = \alpha_r + r - \alpha_r r - O(\varepsilon),$$

which implies that

$$\sum_{j \in A_r} c_j = \alpha_r \geq 1 - O(\varepsilon)r^{-1}, \quad \sum_{j \notin A_r} c_j = O(\varepsilon)r^{-1}. \quad (4.8)$$

Similarly, if we define  $B_r := \{j : \operatorname{Re}(v_j\bar{b}) \geq 1 - r\}$ , then

$$\sum_{j \notin B_r} d_j = O(\varepsilon)r^{-1}. \quad (4.9)$$

We claim that if  $r \geq C\varepsilon$  for some large enough constant  $C > 0$ , then  $A_r \cap B_r \neq \emptyset$ . To see this, we define  $U := \{j : |i - j| \leq W\}$ . By (1.3) and the definition of  $c_j$ , we have  $c_j \geq c_s W^{-1}$  for  $j \in U$ . Clearly, we also have  $d_j \geq \frac{1}{2}c_s W^{-1}$  for  $j \in U$ . Then with (4.8) and (4.9), we have

$$\#\{j \in U \setminus A_r\} = O(\varepsilon)r^{-1}c_s^{-1}W, \quad \#\{j \in U \setminus B_r\} = O(\varepsilon)r^{-1}c_s^{-1}W.$$

If we choose  $r = C\varepsilon$  for some large enough constant  $C > 0$ , then the above two inequalities imply  $A_r \cap B_r \neq \emptyset$ , since  $|U| = W$ . Thus there is an index  $j$  such that

$$\operatorname{Re}(v_j\bar{a}) \geq 1 - r, \quad \operatorname{Re}(v_j\bar{b}) \geq 1 - r. \quad (4.10)$$

Since  $|a| \leq 1$ ,  $|b| \leq 1$ ,  $|v_j| \leq 1$  and  $|a - b| > \delta/2$ , (4.10) is possible only if  $r \gtrsim \delta$ , which contradicts the fact that  $r \rightarrow 0$  when  $\varepsilon \rightarrow 0$ . This proves (4.3).  $\square$

**4.2 Proof of Lemma 1.3.** With Lemma 4.2, we can now give the proof of Lemma 1.3.

*Proof of Lemma 1.3.* We first prove the existence and continuity of the solutions to (1.6). The proof is a standard application of the contraction principle. Denote by  $\mathbf{z} := (z_1, \dots, z_N)$ ,  $\mathbf{x} := (x_1, \dots, x_N)$  and  $\mathbf{M} := ((M_\zeta^\mathbf{g})_1, \dots, (M_\zeta^\mathbf{g})_N)$  with

$$z_i = z\mathbf{1}_{i \in [1, W]} + \tilde{z}\mathbf{1}_{i \notin [1, W]},$$

and

$$x_i \equiv (x_\zeta^\mathbf{g})_i(z, \tilde{z}) := (M_\zeta^\mathbf{g})_i(z, \tilde{z}) - m, \quad \mathbf{M} = \mathbf{x} + m\mathbf{e}_1, \quad m := m_{\text{sc}}(\tilde{z} + i0^+), \quad \mathbf{e}_1 = (1, 1, \dots, 1). \quad (4.11)$$

Using the above notations and recalling Definition 1.2, we can rewrite (1.6) into the following form

$$(m + x_i)^{-1} = M_i^{-1} = -z_i - g_i - (S_0 M)_i + \zeta(\Sigma M)_i = -z_i - g_i - (S_0 \mathbf{x})_i - m(S_0 \mathbf{e}_1)_i + \zeta(\Sigma \mathbf{x})_i + \zeta m(\Sigma \mathbf{e}_1)_i. \quad (4.12)$$

Subtracting  $m^{-1} = -\tilde{z} - m$  from the last equation and using  $S_0 \mathbf{e}_1 = \mathbf{e}_1$ , we get that

$$m^{-1} - (m + x_i)^{-1} = g_i + (z_i - \tilde{z}) + (S_0 \mathbf{x})_i - \zeta m(\Sigma \mathbf{e}_1)_i - \zeta(\Sigma \mathbf{x})_i.$$

Then (4.12) is equivalent to

$$[(1 - m^2 S_0)\mathbf{x}]_i = m^2(g_i + (z_i - \tilde{z})) + m^2 \left( \frac{1}{m + x_i} - \frac{1}{m} + \frac{x_i}{m^2} \right) - \zeta m^3 (\Sigma \mathbf{e}_1)_i - \zeta m^2 (\Sigma \mathbf{x})_i. \quad (4.13)$$

Define iteratively a sequence of vectors  $\mathbf{x}^k \in \mathbb{C}^N$  such that  $\mathbf{x}^0 = \mathbf{0} \in \mathbb{C}^N$  and

$$[(1 - m^2 S_0)\mathbf{x}^{k+1}]_i := m^2(g_i + (z_i - \tilde{z})) + m^2 \left( \frac{1}{m + (\mathbf{x}^k)_i} - \frac{1}{m} + \frac{(\mathbf{x}^k)_i}{m^2} \right) - \zeta m^3 (\Sigma \mathbf{e}_1)_i - \zeta m^2 (\Sigma \mathbf{x}^k)_i. \quad (4.14)$$

In other words, (4.14) defines a mapping  $h : l^\infty(\mathbb{Z}_N) \rightarrow l^\infty(\mathbb{Z}_N)$ :

$$\mathbf{x}^{k+1} = h(\mathbf{x}^k), \quad h_i(\mathbf{x}) := \sum_j (1 - m^2 S_0)_{ij}^{-1} [m^2(g_j + (z_j - \tilde{z})) + q(x_j) - \zeta m^3 (\Sigma \mathbf{e}_1)_j - \zeta m^2 (\Sigma \mathbf{x}^k)_j], \quad (4.15)$$

where

$$q(x) := m^2 \left( \frac{1}{m + x} + \frac{x}{m^2} - \frac{1}{m} \right) = \frac{x^2}{m + x}.$$

Note by the assumptions of Lemma 1.3,  $c_\kappa \leq m \leq 1$  for some constant  $c_\kappa > 0$  depending only on  $\kappa$ . Then with (4.4), it is easy to see that there exists a sufficiently small constant  $0 < \alpha < c_\kappa/2$ , such that  $h$  is a self-mapping

$$h : B_r(l^\infty(\mathbb{Z}_N)) \rightarrow B_r(l^\infty(\mathbb{Z}_N)), \quad B_r(l^\infty(\mathbb{Z}_N)) := \{\mathbf{x} \in l^\infty(\mathbb{Z}_N) : \|\mathbf{x}\|_\infty \leq r\},$$

as long as  $r \leq \alpha$  and

$$\zeta + \|\mathbf{g}\|_\infty + |z - \tilde{z}| \leq c_r \quad (4.16)$$

for some constant  $c_r > 0$  depending on  $r$ . Now it suffices to prove that  $h$  restricted to  $B_r(l^\infty(\mathbb{Z}_N))$  is a contraction, which then implies that  $\mathbf{x} := \lim_{k \rightarrow \infty} \mathbf{x}^k$  exists and is a unique solution to (4.13) subject to the condition  $\|\mathbf{x}\|_\infty \leq r$ .

From the iteration relation (4.15), we obtain that

$$\mathbf{x}^{k+1} - \mathbf{x}^k = \frac{1}{1 - m^2 S_0} [q(\mathbf{x}^k) - q(\mathbf{x}^{k-1})] - \frac{\zeta m^2}{1 - m^2 S_0} \Sigma(\mathbf{x}^k - \mathbf{x}^{k-1}), \quad (4.17)$$

where  $q(\mathbf{x})$  denotes a vector with components  $q(x_i)$ . Using  $|q'(0)| = 0$  and (4.4), we get from (4.17) that

$$\|\mathbf{x}^{k+1} - \mathbf{x}^k\|_\infty \leq C_\kappa (\zeta + \|\mathbf{x}^k\|_\infty + \|\mathbf{x}^{k-1}\|_\infty) \cdot \|\mathbf{x}^k - \mathbf{x}^{k-1}\|_\infty$$

for some constant  $C_\kappa > 0$  depending only on  $\kappa$ . Thus we can first choose a sufficiently small constant  $0 < r < \alpha$  and then the constant  $c_r > 0$  such that  $C_\kappa (c_r + 2r) < 1$ , and  $h$  is a self-mapping on  $B_r(l^\infty(\mathbb{Z}_N))$  under the condition (4.16). In other words,  $h$  is indeed a contraction, which proves the existence and uniqueness of the solution.

Note that with (4.4) and  $\mathbf{x}^0 = \mathbf{0}$ , we get from (4.15) that

$$\|\mathbf{x}^1\|_\infty = O(|z - \tilde{z}| + \zeta + \|\mathbf{g}\|_\infty).$$

With the contraction mapping, we have the bound

$$\|\mathbf{x}\|_\infty \leq \sum_{k=0}^{\infty} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|_\infty \leq \frac{\|\mathbf{x}^1\|_\infty}{1 - C_\kappa (\zeta + 2r)} = O(|z - \tilde{z}| + \zeta + \|\mathbf{g}\|_\infty).$$

This gives the bound (1.9).

We now prove (1.10). We have proved above that both  $(M_\zeta^{\mathbf{g}})_i(z, \tilde{z})$  and  $(M_{\zeta'}^{\mathbf{g}'})_i(z', \tilde{z}')$  exist and satisfy (1.9). Denote by  $m' := m_{\text{sc}}(\tilde{z}' + i0^+)$  and  $x'_i := (M_{\zeta'}^{\mathbf{g}'})_i(z', \tilde{z}') - m'$ . By (4.2), we have

$$|m' - m| = O(|\tilde{z} - \tilde{z}'|). \quad (4.18)$$

Then using (4.13) we can obtain that

$$\begin{aligned} \|\mathbf{x}' - \mathbf{x}\|_\infty &\leq C \|(1 - m^2 S_0)^{-1}\|_{L^\infty \rightarrow L^\infty} \cdot \{ |\tilde{z} - \tilde{z}'| \cdot [\|\mathbf{x}'\|_\infty + \|\mathbf{g}'\|_\infty + |z' - \tilde{z}'| + \|\mathbf{x}'\|_\infty^2 + \zeta'(1 + \|\mathbf{x}'\|_\infty)] \\ &\quad + [\|\mathbf{g} - \mathbf{g}'\|_\infty + |z - z'| + |\tilde{z} - \tilde{z}'| + |\zeta - \zeta'| (1 + \|\mathbf{x}'\|_\infty) + (\zeta + \|\mathbf{x}\|_\infty + \|\mathbf{x}'\|_\infty) \cdot \|\mathbf{x}' - \mathbf{x}\|_\infty] \} \\ &\leq C (\zeta + \|\mathbf{x}\|_\infty + \|\mathbf{x}'\|_\infty) \cdot \|\mathbf{x}' - \mathbf{x}\|_\infty + C (\|\mathbf{g} - \mathbf{g}'\|_\infty + |z - z'| + |\tilde{z} - \tilde{z}'| + |\zeta - \zeta'|). \end{aligned}$$

Applying (1.9) to both  $(M_\zeta^{\mathbf{g}})_i(z, \tilde{z})$  and  $(M_{\zeta'}^{\mathbf{g}'})_i(z', \tilde{z}')$ , we see that for small enough  $c$ ,

$$\|\mathbf{x}' - \mathbf{x}\|_\infty \leq C (\|\mathbf{g} - \mathbf{g}'\|_\infty + |z - z'| + |\tilde{z} - \tilde{z}'| + |\zeta - \zeta'|).$$

Together with (4.18), we obtain (1.10) as desired.  $\square$

**4.3 Proof of Lemma 2.7.** To prove Lemma 2.7, it suffices to prove the result for the case  $\mathbf{g} = 0$ , and we will describe how to relax to the condition  $\mathbf{g} = O(W^{-3/4})$  by using the Lipschitz continuity estimate (1.10) at the end of the proof. In preparation for the proof, we first prove the following lemma.

**Lemma 4.3.** *Suppose that  $\mathbf{g} = 0$  and the assumptions (1.11), (1.12) and (1.15) hold. Then there exist constants  $c > 0$  and  $C > 0$  such that*

$$|(M_\zeta^{\mathbf{0}})_n|^2 - |m|^2 \leq C (|z - \tilde{z}| + \zeta) e^{-c \frac{|n|}{W}}, \quad n \in \mathbb{Z}_N, \quad (4.19)$$

and

$$\frac{1}{W} \sum_{n \in \mathbb{Z}_N} (|m|^2 |(M_\zeta^{\mathbf{0}})_n|^{-2} - 1) \geq c (\operatorname{Im} z - \operatorname{Im} \tilde{z}) - \zeta + O\left(N^{-\frac{3}{2}\varepsilon_*} + N^{-\varepsilon_*} \operatorname{Im} \tilde{z}\right), \quad (4.20)$$

where  $m := m_{\text{sc}}(\tilde{z} + i0^+)$ .

*Proof of Lemma 4.3.* First with (4.5) and the fact that  $(S_0)_{ij} = 0$  if  $|i - j| \geq C_s W$ , we get that

$$[(1 - m^2 S_0)^{-1}]_{ij} - \delta_{ij} = [m^2 (1 - m^2 S_0)^{-1} S_0]_{ij} = O(W^{-1}) \sum_{k \geq \frac{|i-j|}{C_s W}} \left\| \frac{m^2 S_0 + \tau}{1 + \tau} \right\|_{L^\infty \rightarrow L^\infty}^k$$

Therefore with (4.3), we obtain immediately that

$$|[(1 - m^2 S_0)^{-1}]_{ij} - \delta_{ij}| \leq C W^{-1} e^{-c \frac{|i-j|}{W}} \quad (4.21)$$

for some constants  $c, C > 0$ . As in the proof of Lemma 1.3, with  $\mathbf{x}^k$  defined in (4.14), we know that

$$x_n = M_n - m = x_n^1 + \sum_{k \geq 1} (x_n^{k+1} - x_n^k), \quad M_n := (M_\zeta^{\mathbf{0}})_n. \quad (4.22)$$

(Recall that we have proved that  $x_n = \lim_{k \rightarrow \infty} x_n^k$  in the proof of Lemma 1.3 above.) In particular, according to (4.14),  $\mathbf{x}^1$  is given by

$$[(1 - m^2 S_0) \mathbf{x}^1]_i = m^2 (z_i - \tilde{z}) - \zeta m^3 (\Sigma \mathbf{e}_1)_i. \quad (4.23)$$

Then with (4.21) and (4.23), one can show that

$$|x_n^1| \leq C e^{-c \frac{|n|}{W}} (|z - \tilde{z}| + \zeta), \quad n \in \mathbb{Z}_N. \quad (4.24)$$

By (4.17) and (4.21), we have

$$|x_i^{k+1} - x_i^k| \leq C \sum_j \left( W^{-1} e^{-c \frac{|i-j|}{W}} + \delta_{ij} \right) \left[ (|x_j^k| + |x_j^{k-1}|) |x_j^k - x_j^{k-1}| + \zeta \mathbf{1}_{j \in [1, W]} \max_{j' \in [1, W]} |x_{j'}^k - x_{j'}^{k-1}| \right].$$

By induction, it is easy to prove that there are constants  $c, C > 0$  such that

$$|x_n^{k+1} - x_n^k| \leq C e^{-c \frac{|n|}{W}} (|z - \tilde{z}| + \zeta)^{k+1}. \quad (4.25)$$

Together with (4.24) and (4.22), this implies

$$|x_n| = |M_n - m| \leq C(|z - \tilde{z}| + \zeta) e^{-c \frac{|n|}{W}}, \quad n \in \mathbb{Z}_N. \quad (4.26)$$

This proves (4.19) since  $||M_n|^2 - |m|^2| \leq |M_n^2 - m^2|$ .

We now prove (4.20). Using (4.19), we have

$$\frac{1}{W} \sum_{n \in \mathbb{Z}_N} (|m|^2 |M_n|^{-2} - 1) = \frac{1}{W|m|^2} \sum_{n \in \mathbb{Z}_N} (|m|^2 - |M_n|^2) + O(|z - \tilde{z}|^2 + \zeta^2). \quad (4.27)$$

By definition (4.11),

$$|M_n|^2 = |m|^2 + 2 \operatorname{Re}(\bar{m}x_n) + |x_n|^2.$$

Then with (4.26) we get that

$$\frac{1}{W} \sum_n (|M_n|^2 - |m|^2) = \frac{1}{W} \sum_n [2 \operatorname{Re}(\bar{m}x_n) + |x_n|^2] = \frac{2}{W} \sum_n \operatorname{Re}(\bar{m}x_n) + O(|z - \tilde{z}|^2 + \zeta^2).$$

By (1.11) and (1.12), we have

$$\zeta^2 + |\operatorname{Re}(z - \tilde{z})|^2 \leq T^2 + r^2 \leq N^{-3\varepsilon_*/2}, \quad 0 \leq \operatorname{Im} \tilde{z} \leq \operatorname{Im} z \leq N^{-\varepsilon^*},$$

which implies that

$$\zeta^2 + |z - \tilde{z}|^2 \leq \zeta^2 + |\operatorname{Re}(z - \tilde{z})|^2 + \operatorname{Im}(z - \tilde{z})^2 \leq N^{-3\varepsilon_*/2} + N^{-\varepsilon^*} \operatorname{Im}(z - \tilde{z}).$$

Then using (4.22) and (4.25), we obtain that

$$\begin{aligned} \frac{1}{W} \sum_n (|M_n|^2 - |m|^2) &= \frac{2}{W} \sum_n \operatorname{Re}(\bar{m}x_n) + O\left(N^{-\frac{3}{2}\varepsilon_*} + N^{-\varepsilon^*} \operatorname{Im}(z - \tilde{z})\right) \\ &= \frac{2}{W} \sum_n \operatorname{Re}(\bar{m}x_n^1) + O\left(N^{-\frac{3}{2}\varepsilon_*} + N^{-\varepsilon^*} \operatorname{Im}(z - \tilde{z})\right). \end{aligned} \quad (4.28)$$

Summing (4.23) over  $i$ , we get that (recall that we take  $\mathbf{g} = 0$ )

$$(1 - m^2) \sum_i x_i^1 := m^2 \sum_i (z_i - \tilde{z}) - \zeta m^3 (W + 1) = m^2 W (z - \tilde{z}) - \zeta m^3 W + O(1),$$

where we used that  $\sum_i (S_0)_{ij} = 1$  and  $(\Sigma \mathbf{e}_1)_i = 1 + W^{-1}$  for  $i \in \llbracket 1, W \rrbracket$ . Thus for the second term in the second line of (4.28), we have

$$\begin{aligned} \sum_n \operatorname{Re}(\bar{m}x_n^1) &= |m|^2 W \operatorname{Re} \left( \frac{(z - \tilde{z})m - \zeta m^2}{1 - m^2} \right) + O(1) \\ &= |m|^2 W \left( \frac{\zeta}{2} - \frac{\operatorname{Im} z - \operatorname{Im} \tilde{z}}{\sqrt{4 - |\operatorname{Re} \tilde{z}|^2}} + O\left(N^{-\varepsilon^*} \operatorname{Im} \tilde{z}\right) \right) + O(1), \end{aligned} \quad (4.29)$$

where we have used the following special properties of  $m(\tilde{z} + i0^+)$  when  $\tilde{z}$  is a real number, in which case  $m(\tilde{z} + i0^+)$  has unit modulus:

$$\operatorname{Re} \frac{m(a^+)}{1 - m^2(a^+)} = 0, \quad \operatorname{Im} \frac{m(a^+)}{1 - m^2(a^+)} = \frac{1}{\sqrt{4 - a^2}}, \quad \operatorname{Re} \frac{m^2(a^+)}{1 - m^2(a^+)} = -\frac{1}{2}, \quad |a| < 2, \quad a^+ := a + i0^+. \quad (4.30)$$

Here the error  $O(N^{-\varepsilon^*} \operatorname{Im} \tilde{z})$  in (4.29) is due to  $|m(\tilde{z}) - m(\operatorname{Re} \tilde{z} + i0^+)| \leq C \operatorname{Im} \tilde{z}$ . Inserting (4.29) into (4.28), we obtain that for some constant  $c > 0$ ,

$$\frac{1}{W} \sum_n (|M_n|^2 - |m|^2) \leq -c(\operatorname{Im} z - \operatorname{Im} \tilde{z}) + \zeta |m|^2 + O\left(N^{-\frac{3}{2}\varepsilon_*} + N^{-\varepsilon^*} \operatorname{Im} \tilde{z}\right),$$

which, together with (4.27), proves (4.20).  $\square$

With Lemma 4.3, we now finish the proof of Lemma 2.7.

*Proof of Lemma 2.7.* We first assume that  $\mathbf{g} = 0$ . With (4.3) and a perturbation argument, we can show that

$$\left\| \left( \frac{M^2 S + \tau}{1 + \tau} \right)^2 \right\|_{L^\infty \rightarrow L^\infty} < 1 - c$$

for some constant  $c > 0$ . Then (2.6) can be proved as in (4.21). Our main task is to prove (2.7). Assume that

$$(1 - |M|^2 S) \mathbf{u}^0 = \mathbf{v}^0 \quad (4.31)$$

for some vectors  $\mathbf{u}^0, \mathbf{v}^0 \in \mathbb{R}^N$ . Multiplying (4.31) with  $\mathbf{u}^0 |M|^{-2}$  from the left and using the definition of  $S$ , we obtain that

$$\sum_i (|M_i|^{-2} - 1) |\mathbf{u}_i^0|^2 + \sum_{1 \leq i \leq W} \zeta (1 + W^{-1}) |\mathbf{u}_i^0|^2 + \frac{1}{2} \sum_{i,j} S_{ij} (\mathbf{u}_i^0 - \mathbf{u}_j^0)^2 = (\mathbf{u}^0, |M|^{-2} \mathbf{v}^0). \quad (4.32)$$

We define a symmetric operator  $H : L^2(\mathbb{T}) \mapsto L^2(\mathbb{T})$ , where  $\mathbb{T} := \llbracket -(\log N)^4 W, (\log N)^4 W \rrbracket$  and

$$H := H_0 + H_1,$$

with

$$H_0 : (\mathbf{u}, H_0 \mathbf{v}) = \frac{1}{4} \sum_{i,j \in \mathbb{T}} S_{ij} (\mathbf{u}_i - \mathbf{u}_j) (\mathbf{v}_i - \mathbf{v}_j), \quad \mathbf{u}, \mathbf{v} \in L^2(\mathbb{T}),$$

and

$$(H_1)_{ij} := \delta_{ij} [(|M_i|^{-2} - 1) + \zeta \mathbf{1}_{1 \leq i \leq W} (1 + W^{-1})].$$

For any vector  $\mathbf{u}$ , we denote by  $\mathbf{u}|_{\mathbb{T}}$  the restriction of  $\mathbf{u}$  to  $L^2(\mathbb{T})$ . Then with (4.19) and the fact that  $|m| \leq 1$ , we can rewrite (4.32) as

$$(\mathbf{u}^0|_{\mathbb{T}}, H \mathbf{u}^0|_{\mathbb{T}}) + \frac{1}{4} \sum_{i,j} S_{ij} (\mathbf{u}_i^0 - \mathbf{u}_j^0)^2 \leq (\mathbf{u}^0, |M|^{-2} \mathbf{v}^0) + O(N^{-10}) \|\mathbf{u}^0\|_2^2. \quad (4.33)$$

First we claim that

$$H \geq c \operatorname{Im} z (\log N)^{-4} \quad (4.34)$$

for some constant  $c > 0$ . With Temple's inequality, we have the following estimate on the ground state energy of  $H$ :

$$H \geq E_0(H) \geq \langle H \rangle_\phi - \frac{\langle (H)^2 \rangle_\phi - \langle H \rangle_\phi^2}{E_1(H) - \langle H \rangle_\phi}, \quad (4.35)$$

for any  $\phi \in L^2(\mathbb{T})$  such that  $\|\phi\|_2 = 1$  and  $\langle H \rangle_\phi < E_1(H)$ , where  $E_0(H)$  and  $E_1(H)$  are the lowest two eigenvalues of  $H$ . Applying min-max principle to  $H \geq H_0 - \|H_1\|_{L^2 \rightarrow L^2}$ , we obtain that

$$E_1(H) \geq E_1(H_0) - \|H_1\|_{L^2 \rightarrow L^2}. \quad (4.36)$$

By (4.19), we have  $\|H_1\|_{L^2 \rightarrow L^2} = O(|z - \tilde{z}| + \zeta + \operatorname{Im} \tilde{z})$ . We then claim that

$$E_1(H_0) \geq c (\log N)^{-13} \quad (4.37)$$

for some constant  $c > 0$ . Recall that  $S \equiv S^\zeta = S_0 - \zeta \Sigma$  with

$$\frac{1}{4} \sum_{i,j \in \mathbb{T}} \Sigma_{ij} (\mathbf{u}_i - \mathbf{u}_j)^2 \leq 1, \quad \forall \mathbf{u} \in L^2(\mathbb{T}), \quad \|\mathbf{u}\|_2 = 1.$$

Then again by min-max principle, it suffices to prove the following lemma.

**Lemma 4.4.** *For  $s_{ij}$  satisfying (1.1)-(1.3), there exists a constant  $c > 0$  such that*

$$\frac{1}{4} \sum_{i,j \in \mathbb{T}} s_{ij} (\mathbf{u}_i - \mathbf{u}_j)^2 \geq c (\log N)^{-13}, \quad \forall \mathbf{u} \in L^2(\mathbb{T}), \quad \|\mathbf{u}\|_2 = 1, \quad \mathbf{u} \perp (1, 1, \dots, 1).$$

We postpone its proof until we finish the proof of Lemma 2.7. We now choose the trial state  $\phi$  as a constant vector in (4.35), i.e.,

$$\phi_0 = \frac{1}{\sqrt{|\mathbb{T}|}}(1, 1, \dots, 1).$$

Then by definition,  $H_0\phi_0 = \mathbf{0}$  and  $\langle H \rangle_{\phi_0} \leq \|H_1\|_{L^2 \rightarrow L^2} \ll E_1(H)$  by (4.36) and (4.37). Then by (4.35) and (4.36), we have

$$H \geq \langle H_1 \rangle_{\phi_0} - \frac{\|H_1\|_{L^2 \rightarrow L^2}^2}{E_1(H) - \langle H_1 \rangle_{\phi_0}} \geq \langle H_1 \rangle_{\phi_0} - \frac{\|H_1\|_{L^2 \rightarrow L^2}^2}{E_1(H_0) - 2\|H_1\|_{L^2 \rightarrow L^2}}. \quad (4.38)$$

By the definition of  $H_1$ , we have

$$\begin{aligned} \langle H_1 \rangle_{\phi_0} &= \frac{1}{|\mathbb{T}|} \sum_{n \in \mathbb{T}} [(|M_n|^{-2} - 1) + \zeta \mathbf{1}_{n \in [1, W]}(1 + W^{-1})] \\ &= \frac{1}{|\mathbb{T}|} \sum_{n \in \mathbb{T}} (1 - |m|^2)|M_n|^{-2} + \frac{1}{|\mathbb{T}|} \sum_{n \in \mathbb{T}} (|m|^2|M_n|^{-2} - 1) + \frac{\zeta(W+1)}{|\mathbb{T}|} \\ &\geq c \operatorname{Im} \tilde{z} + O(N^{-10}) + \frac{1}{|\mathbb{T}|} \sum_{n \in \mathbb{Z}_N} (|m|^2|M_n|^{-2} - 1) + \frac{\zeta(W+1)}{|\mathbb{T}|} \\ &\geq c \operatorname{Im} z (\log N)^{-4} + O\left(N^{-\frac{3}{2}\varepsilon^*} + N^{-\varepsilon^*} \operatorname{Im} z\right), \end{aligned}$$

where we used (4.19) and  $|m|^2 \leq 1 - c \operatorname{Im} \tilde{z}$  (by (4.1)) in the third step, and (4.20) in the last step. Together with (4.38),  $\|H_1\|_{L^2 \rightarrow L^2}^2 = O(N^{-\frac{3}{2}\varepsilon^*} + N^{-\varepsilon^*} \operatorname{Im} z)$  and (4.37), this proves (4.34).

With (4.34), (4.33) gives that for some  $c > 0$ ,

$$c \operatorname{Im} z (\log N)^{-4} \sum_{i \in \mathbb{T}} |\mathbf{u}_i^0|^2 + \frac{1}{4} \sum_{i, j} S_{ij} (\mathbf{u}_i^0 - \mathbf{u}_j^0)^2 \leq (\mathbf{u}^0, |M|^{-2} \mathbf{v}^0) + O(N^{-10}) \|\mathbf{u}^0\|_2^2.$$

Now for some fixed  $i_0 \in \mathbb{Z}_N$ , we choose  $\mathbf{v}^0 = S \mathbf{e}_{i_0}$ . Then the above inequality becomes

$$c \operatorname{Im} z (\log N)^{-4} \sum_{i \in \mathbb{T}} |\mathbf{u}_i^0|^2 + \frac{1}{4} \sum_{i, j} S_{ij} (\mathbf{u}_i^0 - \mathbf{u}_j^0)^2 \leq (S|M|^{-2} \mathbf{u}^0)_{i_0} + O(N^{-10}) \|\mathbf{u}^0\|_2^2. \quad (4.39)$$

In the following, we suppose  $\|\mathbf{u}^0\|_\infty \gg W^{-1}$ , otherwise the proof is done. Since for any  $i \in \mathbb{Z}_N$ ,

$$(\mathbf{u}^0 - |M|^2 S \mathbf{u}^0)_i = (S \mathbf{e}_{i_0})_i = O(W^{-1}), \quad (4.40)$$

we must have

$$\|\mathbf{u}^0\|_\infty \asymp \|S \mathbf{u}^0\|_\infty.$$

Now we decompose  $\mathbf{u}^0$  as follows:

$$\mathbf{u}_i^0 = u + \tilde{\mathbf{u}}_i, \quad \text{with } u = \frac{1}{N} \sum_{i \in \mathbb{Z}_N} \mathbf{u}_i^0, \quad \sum_i \tilde{\mathbf{u}}_i = 0.$$

Suppose  $|u| \geq 10 \|\tilde{\mathbf{u}}\|_\infty$ , then we have

$$\max_i |\mathbf{u}_i^0| \leq 2 \min_i |\mathbf{u}_i^0|.$$

Together with (4.39), it implies that if  $|u| \geq 10 \|\tilde{\mathbf{u}}\|_\infty$ , then

$$\|\mathbf{u}^0\|_\infty \leq 2|u| \leq C(W \operatorname{Im} z)^{-1}. \quad (4.41)$$

On the other hand, if  $|u| \leq 10 \|\tilde{\mathbf{u}}\|_\infty$ , with (4.31), (4.19) and the definition of  $S$  in Definition 1.2, we get that

$$\tilde{\mathbf{u}} - |M|^2 S \tilde{\mathbf{u}} = O(W^{-1} + (\zeta + |z - \tilde{z}|)|u|). \quad (4.42)$$

Then in this case, with (4.40) and (4.42) it is easy to see that

$$\|\mathbf{u}^0\|_\infty \asymp \|S \mathbf{u}^0\|_\infty \asymp \|\tilde{\mathbf{u}}\|_\infty \asymp \|S \tilde{\mathbf{u}}\|_\infty \asymp \|S_0 \tilde{\mathbf{u}}\|_\infty. \quad (4.43)$$

By (1.2), we have

$$\sum_j (S_0 \tilde{\mathbf{u}})_j = 0,$$

which implies

$$\|S_0 \tilde{\mathbf{u}}\|_\infty \leq \max_{i,j} |(S_0 \tilde{\mathbf{u}})_j - (S_0 \tilde{\mathbf{u}})_i|. \quad (4.44)$$

Using (1.2), for fixed  $i \leq j \in \mathbb{Z}_N$  we have

$$|(S_0 \tilde{\mathbf{u}})_j - (S_0 \tilde{\mathbf{u}})_i|^2 = \left| \sum_{x,y} (S_0)_{ix} (S_0)_{jy} (\tilde{\mathbf{u}}_x - \tilde{\mathbf{u}}_y) \right|^2 \leq \sum_{x,y} (S_0)_{ix} (S_0)_{jy} |\tilde{\mathbf{u}}_x - \tilde{\mathbf{u}}_y|^2. \quad (4.45)$$

The lower bound in (1.3) shows that  $S_0$  has a core, i.e., there is a constant  $c_s > 0$  such that  $(S_0)_{xy} \geq c_s W^{-1}$  if  $|x - y| \leq W$ . Then for any fixed  $i \leq j \in \mathbb{Z}_N$ , we choose  $x_0, x_1, x_2, \dots, x_n$  for some  $n = O(N/W)$  such that

$$i = x_0 \leq x_1 \leq x_2 \leq \dots \leq x_{n-1} \leq x_n = j, \quad \text{with } W/3 \leq |x_k - x_{k+1}| \leq W/2, \quad \forall k.$$

Furthermore, set  $x'_0 = x$  and  $x'_n = y$ . Clearly for any choices of  $x'_k$ ,  $1 \leq k \leq n-1$ , we have

$$\tilde{\mathbf{u}}_y - \tilde{\mathbf{u}}_x = \sum_{k=1}^n (\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}) \Rightarrow |\tilde{\mathbf{u}}_y - \tilde{\mathbf{u}}_x|^2 \leq \frac{CN}{W} \sum_{k=1}^n |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2.$$

For our goal, we will choose  $x'_k$ 's such that

$$x'_k \in \llbracket x_k - W/4, x_k + W/4 \rrbracket, \quad 1 \leq k \leq n-1.$$

Taking averaging over all  $x'_k$ ,  $1 \leq k \leq n-1$ , in the above regions, we get that

$$|\tilde{\mathbf{u}}_y - \tilde{\mathbf{u}}_x|^2 \leq \frac{N}{W} \left( \text{Average}_{x'_1, x'_2, \dots, x'_{n-1}} \right) \sum_{k=1}^n |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2.$$

Note that by our choices, we always have  $|x'_k - x'_{k-1}| \leq W$  and  $S_{x'_k x'_{k-1}} \geq \frac{1}{2} c_s W^{-1}$  for  $2 \leq k \leq n-1$ , which gives that

$$\begin{aligned} \text{Average}_{x'_{k-1}, x'_k} |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2 &\leq \frac{4}{W^2} \sum_{x'_k, x'_{k-1} \in \llbracket x_{k-1} - W/4, x_k + W/4 \rrbracket} |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2 \\ &\leq \frac{8c_s^{-1}}{W} \sum_{x'_k, x'_{k-1} \in \llbracket x_{k-1} - W/4, x_k + W/4 \rrbracket} S_{x'_k x'_{k-1}} |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2. \end{aligned}$$

Together with (4.45), we get that for some constant  $C > 0$ ,

$$\begin{aligned} |(S_0 \tilde{\mathbf{u}})_j - (S_0 \tilde{\mathbf{u}})_i|^2 &\leq \sum_{x,y} (S_0)_{ix} (S_0)_{jy} \left[ \frac{CN}{W^2} \sum_{k=2}^{n-1} \sum_{x'_k, x'_{k-1} \in \llbracket x_{k-1} - W/4, x_k + W/4 \rrbracket} S_{x'_k x'_{k-1}} |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2 \right] \\ &\quad + \sum_{x,y} (S_0)_{ix} (S_0)_{jy} \frac{CN}{W} \left[ \frac{2}{W} \sum_{x': |x' - x_1| \leq W/4} |\tilde{\mathbf{u}}_{x'} - \tilde{\mathbf{u}}_x|^2 + \frac{2}{W} \sum_{y': |y' - x_{n-1}| \leq W/4} |\tilde{\mathbf{u}}_y - \tilde{\mathbf{u}}_{y'}|^2 \right]. \end{aligned}$$

For the first term on the right-hand side, we have

$$\sum_{k=2}^{n-1} \sum_{x'_k, x'_{k-1} \in \llbracket x_{k-1} - W/4, x_k + W/4 \rrbracket} S_{x'_k x'_{k-1}} |\tilde{\mathbf{u}}_{x'_k} - \tilde{\mathbf{u}}_{x'_{k-1}}|^2 \leq C \sum_{k,l \in \mathbb{Z}_N} S_{kl} (\tilde{\mathbf{u}}_k - \tilde{\mathbf{u}}_l)^2.$$

For the terms in the second line, we notice that

$$|x' - x| \leq |x' - x_1| + |x_1 - i| + |i - x| \leq C_s W + W$$



for all  $x'$  such that  $|x' - x_1| \leq W/4$ , where  $C_s$  is the constant appeared in (1.3). Then we can subdivide the interval  $\llbracket x, x' \rrbracket$  or  $\llbracket x', x \rrbracket$  into subintervals with lengths  $\leq W/2$ , and proceed as above to get

$$\sum_x (S_0)_{ix} \frac{2}{W} \sum_{|x' - x_1| \leq W/4} |\tilde{\mathbf{u}}_{x'} - \tilde{\mathbf{u}}_x|^2 \leq \frac{C}{W} \sum_{1 \leq k, l \leq N} S_{kl} (\tilde{\mathbf{u}}_k - \tilde{\mathbf{u}}_l)^2$$

for some constant  $C > 0$  that is independent of the choice of  $x'$ . In sum, we have obtained that

$$|(S_0 \tilde{\mathbf{u}})_j - (S_0 \tilde{\mathbf{u}})_i|^2 \leq \frac{CN}{W^2} \sum_{1 \leq k, l \leq N} S_{kl} (\tilde{\mathbf{u}}_k - \tilde{\mathbf{u}}_l)^2 = \frac{CN}{W^2} \sum_{1 \leq k, l \leq N} S_{kl} (\mathbf{u}_k^0 - \mathbf{u}_l^0)^2.$$

Then from (4.43) and (4.44), we obtain that

$$\|\mathbf{u}^0\|_\infty^2 \leq \frac{CN}{W^2} \sum_{1 \leq k, l \leq N} S_{kl} (\mathbf{u}_k^0 - \mathbf{u}_l^0)^2.$$

Plugging it into (4.39), we get that if  $|u| \leq 10\|\tilde{\mathbf{u}}\|_\infty$ , then

$$\frac{W^2}{N} \|\mathbf{u}^0\|_\infty^2 \leq C \sum_{1 \leq k, l \leq N} S_{kl} (\mathbf{u}_k^0 - \mathbf{u}_l^0)^2 \leq C \|\mathbf{u}^0\|_\infty + O(N^{-10}) \|\mathbf{u}^0\|_2^2 \Rightarrow \|\mathbf{u}^0\|_\infty \leq \frac{CN}{W^2}. \quad (4.46)$$

In sum, by our choice of  $\mathbf{v}^0 = S\mathbf{e}_{i_0}$  and (4.31), we obtain from (4.41) and (4.46) that

$$\left\| (1 - S|M|^2)^{-1} S \right\|_{\max} \leq C \left( \frac{1}{W \operatorname{Im} z} + \frac{N}{W^2} \right),$$

which completes the proof of (2.7) in the case with  $\mathbf{g} = 0$ .

Given any  $\mathbf{g} \in \mathbb{R}^N$  such that  $\|\mathbf{g}\|_\infty \leq W^{-3/4}$ , we can write

$$M_\zeta^{\mathbf{g}} = M_\zeta^0 + \mathcal{E},$$

where  $\mathcal{E}$  is a diagonal matrix with  $\max_i |\mathcal{E}_{ii}| = O(\|\mathbf{g}\|_\infty) = O(W^{-3/4})$  by the Lipschitz continuity estimate (1.10). Then (2.6) can be obtained by combing (2.6) in the case  $\mathbf{g} = 0$  with a standard perturbation argument. For (2.7), we write

$$\left(1 - S|M_\zeta^{\mathbf{g}}|^2\right)^{-1} S = \left(1 - S|M_\zeta^0|^2\right)^{-1} S + \left(1 - S|M_\zeta^0|^2\right)^{-1} S \left(|M_\zeta^{\mathbf{g}}|^2 - |M_\zeta^0|^2\right) \left(1 - S|M_\zeta^{\mathbf{g}}|^2\right)^{-1} S. \quad (4.47)$$

Using (2.7) in the case  $\mathbf{g} = 0$  and the bound

$$\left\| \left(1 - S|M_\zeta^0|^2\right)^{-1} S \right\|_{L^\infty \rightarrow L^\infty} \leq N \left\| \left(1 - S|M_\zeta^0|^2\right)^{-1} S \right\|_{\max},$$

we get from (4.47) that

$$\left\| \left(1 - S|M_\zeta^{\mathbf{g}}|^2\right)^{-1} S \right\|_{\max} \leq \left\| \left(1 - S|M_\zeta^0|^2\right)^{-1} S \right\|_{\max} + O\left(\left(\frac{N}{W \operatorname{Im} z} + \frac{N^2}{W^2}\right) W^{-3/4}\right) \cdot \left\| \left(1 - S|M_\zeta^{\mathbf{g}}|^2\right)^{-1} S \right\|_{\max}.$$

Together with (1.15), this implies (2.7) for any  $\mathbf{g}$  such that  $\|\mathbf{g}\|_\infty \leq W^{-3/4}$ .  $\square$

*Proof of Lemma 4.4.* Since the matrix  $S_0 = (s_{ij})$  has a core by (1.3), it suffices to prove that

$$\sum_{i, j \in \mathbb{T}} \hat{s}_{ij} (\mathbf{u}_i - \mathbf{u}_j)^2 \geq c(\log N)^{-13}, \quad \forall \mathbf{u} \in L^2(\mathbb{T}), \quad \|\mathbf{u}\|_2 = 1, \quad \mathbf{u} \perp (1, 1, \dots, 1), \quad (4.48)$$

where

$$\hat{s}_{ij} := \frac{1}{W} \mathbf{1}_{|i-j| \leq W}.$$

Then we define the following two symmetric operators  $F_{0,1} : L^2(\mathbb{T}) \mapsto L^2(\mathbb{T})$  such that for any  $\mathbf{u}, \mathbf{v} \in L^2(\mathbb{T})$ ,

$$(\mathbf{u}, F_0 \mathbf{v}) = \frac{1}{W(\log N)^5} \sum_{i,j \in \mathbb{T}, |i-j|_{\mathbb{T}} \leq W} (\mathbf{u}_i - \mathbf{u}_j) (\mathbf{v}_i - \mathbf{v}_j),$$

where  $|\cdot|_{\mathbb{T}}$  denotes the periodic distance on  $\mathbb{T}$ , and

$$(\mathbf{u}, F_1 \mathbf{v}) = \sum_{i,j \in \mathbb{T}} \tilde{s}_{ij} (\mathbf{u}_i - \mathbf{u}_j) (\mathbf{v}_i - \mathbf{v}_j), \quad \tilde{s}_{ij} := \hat{s}_{ij} - \frac{1}{W(\log N)^5} \mathbb{1}_{|i-j|_{\mathbb{T}} \leq W}.$$

We first show that for some constant  $c > 0$ ,

$$E_1(F_0) \geq c(\log N)^{-13}, \quad (4.49)$$

where  $E_1(F_0)$  denotes the second lowest eigenvalue of  $F_0$ . Without loss of generality, we can regard  $F_0$  as an operator on  $L^2(\mathbb{T}, \mathbb{C})$  consisting of *complex*  $L^2$  vectors. Since  $F_0$  is a periodic operator on  $L^2(\mathbb{T}, \mathbb{C})$ , its eigenvectors are the unit complex vectors with Fourier components:

$$\mathbf{w}_p : (\mathbf{w}_p)_k := \frac{1}{\sqrt{|\mathbb{T}|}} e^{ipk}, \quad k \in \mathbb{T}, \quad \text{with } p = \frac{2\pi n}{|\mathbb{T}|}, \quad n \in \mathbb{T}.$$

Then for any  $p \neq 0$ , we have

$$\begin{aligned} (\mathbf{w}_p, F_0 \mathbf{w}_p) &= \frac{1}{W(\log N)^5} \sum_{|k-l|_{\mathbb{T}} \leq W} |(\mathbf{w}_p)_k - (\mathbf{w}_p)_l|^2 = \frac{1}{|\mathbb{T}|W(\log N)^5} \sum_{|k-l|_{\mathbb{T}} \leq W} [2 - 2 \cos(p(k-l))] \\ &= \frac{1}{W(\log N)^5} \sum_{|n| \leq W} [2 - 2 \cos(pn)] \geq \frac{c}{W(\log N)^5} \frac{W^3}{|\mathbb{T}|^2} \geq c(\log N)^{-13}. \end{aligned}$$

This proves (4.49).

We now show that  $F_1$  defines a positive operator. For simplicity of notations, we let  $L = |\mathbb{T}|$  and shift  $\mathbb{T}$  to  $\mathbb{T} := \llbracket 1, L \rrbracket$ . Then  $\tilde{s}_{ij}$  can be written as

$$\tilde{s}_{ij} = (1 - (\log N)^{-5}) \hat{s}_{ij} - \frac{1}{W(\log N)^5} (\mathbb{1}_{1 \leq i \leq W, L-W+i \leq j \leq L} + \mathbb{1}_{1 \leq j \leq W, L-W+j \leq i \leq L}). \quad (4.50)$$

Fix any  $\mathbf{u} \in L^2(\mathbb{T})$ . The following proof is very similar to the one below (4.45), so we shall omit some details. For any fixed  $1 \leq i \leq W$  and  $L-W \leq j \leq L$ , we choose  $x_0, x_1, \dots, x_n$  for some  $n = O((\log N)^4)$  such that

$$i = x_0 \leq x_1 \leq x_2 \leq \dots \leq x_{n-1} \leq x_n = j, \quad \text{with } W/3 \leq |x_k - x_{k+1}| \leq W/2, \quad \forall k.$$

Moreover, we set  $x'_0 = i$  and  $x'_n = j$ . Then we can get as before that

$$|\mathbf{u}_i - \mathbf{u}_j|^2 \leq C(\log N)^4 \left( \text{Average}_{x'_1, x'_2, \dots, x'_{n-1}} \right) \sum_{k=1}^n \left| \mathbf{u}_{x'_k} - \mathbf{u}_{x'_{k-1}} \right|^2,$$

where we took average over all  $x'_k \in \llbracket x_k - W/4, x_k + W/4 \rrbracket$ ,  $1 \leq k \leq n-1$ . Note that by our choices, we always have  $|x'_k - x'_{k-1}| \leq W$  and  $\hat{s}_{x'_k, x'_{k-1}} = W^{-1}$  for  $1 \leq k \leq n$ , which gives that

$$\begin{aligned} & \frac{1}{W(\log N)^5} \sum_{1 \leq i \leq W, L-W \leq j \leq L} |\mathbf{u}_i - \mathbf{u}_j|^2 \\ & \leq \frac{1}{W(\log N)^5} \sum_{1 \leq i \leq W, L-W \leq j \leq L} \left[ \frac{C(\log N)^4}{W} \sum_{k=2}^{n-1} \sum_{x'_k, x'_{k-1} \in \llbracket x_{k-1} - W/4, x_k + W/4 \rrbracket} \hat{s}_{x'_k, x'_{k-1}} \left| \mathbf{u}_{x'_k} - \mathbf{u}_{x'_{k-1}} \right|^2 \right] \\ & + \frac{1}{W(\log N)^5} \sum_{1 \leq i \leq W, L-W \leq j \leq L} C(\log N)^4 \left[ \sum_{x: |x-x_1| \leq W/4} \hat{s}_{xi} |\mathbf{u}_x - \mathbf{u}_i|^2 + \sum_{y: |y-x_{n-1}| \leq W/4} \hat{s}_{jy} |\mathbf{u}_y - \mathbf{u}_j|^2 \right] \\ & \leq C(\log N)^{-1} \sum_{k,l \in \mathbb{T}} \hat{s}_{kl} (\mathbf{u}_k - \mathbf{u}_l)^2. \end{aligned}$$

Then by (4.50), it is easy to see that  $F_1$  is a positive operator. Thus by min-max principle we have

$$E_1(F_0 + F_1) \geq E_1(F_0),$$

which proves (4.48) together with (4.49). □

## REFERENCES

- [1] P. Bourgade, L. Erdős, H.-T. Yau, and J. Yin, *Universality for a class of random band matrices*, Advances in Theoretical and Mathematical Physics **21** (2017), no. 3, 739–800.
- [2] P. Bourgade, H.-T. Yau, and J. Yin, *Random band matrices in the delocalized phase, I: Quantum unique ergodicity and universality*, in preparation (2018).
- [3] L. Erdős, A. Knowles, and H.-T. Yau, *Averaging fluctuations in resolvents of random band matrices*, Ann. Henri Poincaré **14** (2013), 1837–1926.
- [4] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin, *Delocalization and diffusion profile for random band matrices*, Comm. Math. Phys. **323** (2013), no. 1, 367–416.
- [5] ———, *The local semicircle law for a general class of random matrices*, Elect. J. Probab. **18** (2013), no. 59, 1–58.
- [6] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin, *Spectral statistics of Erdős-Rényi graphs II: Eigenvalue spacing and the extreme eigenvalues*, Comm. Math. Phys. **314** (2012), 587–640.
- [7] L. Erdős, B. Schlein, and H.-T. Yau, *Local semicircle law and complete delocalization for Wigner random matrices*, Commun. Math. Phys. **287** (2008), no. 2, 641–655.
- [8] L. Erdős, H.-T. Yau, and J. Yin, *Bulk universality for generalized Wigner matrices*, Probab. Theory Related Fields **154** (2012), no. 1-2, 341–407.
- [9] A. Knowles and J. Yin, *The isotropic semicircle law and deformation of Wigner matrices*, Comm. Pure Appl. Math. **66** (2013), 1663–1749.
- [10] F. Yang and J. Yin, *Random band matrices in the delocalized phase, III: Averaging fluctuations*, in preparation (2018).