

## Local single ring theorem on optimal scale

Zhigang Bao\*  
HKUST  
mazgbao@ust.hk

László Erdős\*  
IST Austria  
lerdos@ist.ac.at

Kevin Schnelli\*  
KTH Royal Institute of Technology  
schnelli@kth.se

Let  $U$  and  $V$  be two independent  $N$  by  $N$  random matrices that are distributed according to Haar measure on  $U(N)$ . Let  $\Sigma$  be a non-negative deterministic  $N$  by  $N$  matrix. The *single ring theorem* [26] asserts that the empirical eigenvalue distribution of the matrix  $X := U\Sigma V^*$  converges weakly, in the limit of large  $N$ , to a deterministic measure which is supported on a single ring centered at the origin in  $\mathbb{C}$ . Within the bulk regime, *i.e.*, in the interior of the single ring, we establish the convergence of the empirical eigenvalue distribution on the optimal local scale of order  $N^{-1/2+\varepsilon}$  and establish the optimal convergence rate. The same results hold true when  $U$  and  $V$  are Haar distributed on  $O(N)$ .

*Date:* March 1, 2019

*Keywords:* Non-hermitian random matrices, local eigenvalue density, single ring theorem, free convolution

*AMS Subject Classification (2010):* 46L54, 60B20

### 1. INTRODUCTION AND MAIN RESULT

Consider the  $N \times N$  random matrix of the form

$$X \equiv X_N = U\Sigma V^*, \quad (1.1)$$

where  $U \equiv U_N$  and  $V \equiv V_N$  are two independent sequences of random matrices, which are both Haar distributed on either the unitary group,  $U(N)$ , of degree  $N$ , or on the orthogonal group,  $O(N)$ , of degree  $N$ . Moreover, let  $\Sigma \equiv \Sigma_N$  be a sequence of  $N \times N$  deterministic non-negative definite diagonal matrices. Note that in general  $X$  is not hermitian and most of its eigenvalues are genuinely complex numbers. In fact, almost surely the matrix  $X$  is not normal. Let  $\lambda_j(X)$ ,  $j = 1, 2, \dots, N$ , be the eigenvalues of  $X$  and let

$$\mu_X := \frac{1}{N} \sum_{j=1}^N \delta_{\lambda_j(X)} \quad (1.2)$$

be the (normalized) empirical spectral distribution of  $X$ . We define  $\mu_\Sigma$  analogously.

**Assumption 1.1.** *We assume that the sequence  $(\Sigma_N)$  is uniformly bounded, *i.e.*, there exists a finite constant  $S_+$  such that*

$$0 \leq \Sigma_N \leq S_+. \quad (1.3)$$

From this assumption it follows that there is a constant  $0 < s_+ < \infty$  such that, for all  $N \in \mathbb{N}$ ,

$$\text{supp } \mu_\Sigma \subset [0, s_+]. \quad (1.4)$$

---

\*Partially supported by ERC Advanced Grant RANMAT No. 338804.

We first consider the situation where there exists a limiting measure  $\mu_\sigma^*$  of  $\mu_\Sigma$ , *i.e.*,

$$d_L(\mu_\Sigma, \mu_\sigma) \rightarrow 0, \quad (1.5)$$

as  $N \rightarrow \infty$ , where  $d_L$  denotes the Lévy distance. Given such a  $\mu_\sigma$  on  $[0, \infty)$ , we define

$$r_- := \left( \int_{\mathbb{R}^+} x^{-2} d\mu_\sigma(x) \right)^{-\frac{1}{2}}, \quad r_+ := \left( \int_{\mathbb{R}^+} x^2 d\mu_\sigma(x) \right)^{\frac{1}{2}}, \quad (1.6)$$

where we set  $r_- = 0$  in case the integral in its definition diverges. Note that if  $\mu_\sigma$  is supported more than one point, we have  $r_- < r_+$  as follows from Schwarz inequality. We let

$$\mathcal{R}_\sigma \equiv \mathcal{R}(\mu_\sigma) := \{w \in \mathbb{C} : r_- < |w| < r_+\} \quad (1.7)$$

be the ring in  $\mathbb{C}$  with radii  $r_-$  and  $r_+$ . In case  $r_- = 0$ ,  $\mathcal{R}_\sigma$  is the punctuated disc of radius  $r_+$ .

For a probability measure  $\mu$  on  $\mathbb{R}$  we denote by  $\mu^{\text{sym}}$  its symmetrization, *i.e.*,  $\mu^{\text{sym}}(A) := \frac{1}{2}[\mu(A) + \mu(-A)]$  for any Borel set  $A \subset \mathbb{R}$ . For  $r \in \mathbb{R}^+$ , set

$$\mu_{\sigma,r} := \mu_\sigma^{\text{sym}} \boxplus \delta_r^{\text{sym}}, \quad (1.8)$$

where  $\boxplus$  denotes the free additive convolution of probability measures on  $\mathbb{R}$ ; see Subsection 2.1.

Given a probability measure  $\mu$  on  $\mathbb{R}$ , its *Stieltjes transform*,  $m_\mu$ , on the complex upper half-plane  $\mathbb{C}^+ := \{z \in \mathbb{C} : \text{Im } z > 0\}$  is defined by

$$m_\mu(z) := \int_{\mathbb{R}} \frac{d\mu(x)}{x - z}, \quad z \in \mathbb{C}^+. \quad (1.9)$$

**Theorem 1.2** (Single ring theorem, [26]). *Assume that Assumption 1.1 holds and that there is a compactly supported probability measure  $\mu_\sigma$  on  $[0, \infty)$ , which is supported at more than one point, such that (1.5) holds. Assume in addition that there are constants  $k, k_1 > 0$  such that*

$$\text{Im } m_{\mu_\Sigma}(z) \leq k_1 \quad (1.10)$$

*on  $\{z \in \mathbb{C}^+ : \text{Im } z > N^{-k}\}$ . Then the empirical spectral distribution  $\mu_X$  converges weakly (in probability) to a deterministic probability measure  $\rho_\sigma$  supported on  $\overline{\mathcal{R}_\sigma}$ . The limiting measure is absolutely continuous with respect to Lebesgue measure and given by*

$$\rho_\sigma(w) d^2w = \frac{1}{2\pi} \Delta_w \left( \int_{\mathbb{R}} \log |s| \mu_{\sigma,|w|}(ds) \right) d^2w, \quad w \in \mathcal{R}_\sigma, \quad (1.11)$$

where  $\Delta_w = 4\partial_w \partial_{\bar{w}}$  is the Laplacian on  $\mathbb{C}$  and  $d^2w \equiv dw \wedge d\bar{w}$  is Lebesgue measure on  $\mathbb{C}$ .

*Remark 1.3.* In Theorem 1.2,  $U$  and  $V$  may be both Haar distributed on  $U(N)$  or on  $O(N)$ .

*Remark 1.4.* In its original form Theorem 1.2 was proved by Guionnet, Krishnapur and Zeitouni in [26] under a further assumption on the smallest singular value of the matrix  $X - z$ ,  $z \in \mathbb{C}$ . This hard-to-check condition was removed by Rudelson and Vershynin in [34] (*c.f.*, Theorem 2.6 below), which yields Theorem 1.2.

*Remark 1.5.* The measure  $\rho_\sigma$  was first computed in [27]. It has a direct interpretation in free probability theory, in fact it is the Brown measure of the free product of a Haar unitary and an element  $\sigma$  on a noncommutative probability space; see [27] for more details.

---

\*We will often use the convention that capital letters indicate random matrices and the corresponding small letters indicate their limiting objects.

**1.1. Local single ring law.** To state our results, we use the following definition on high-probability estimates from [20]. In Appendix A we collect some of its properties.

**Definition 1.6.** Let  $\mathcal{X} \equiv \mathcal{X}^{(N)}$ ,  $\mathcal{Y} \equiv \mathcal{Y}^{(N)}$  be two sequences of nonnegative random variables. We say that  $\mathcal{Y}$  stochastically dominates  $\mathcal{X}$  if, for all (small)  $\epsilon > 0$  and (large)  $D > 0$ ,

$$\mathbb{P}(\mathcal{X}^{(N)} > N^\epsilon \mathcal{Y}^{(N)}) \leq N^{-D}, \quad (1.12)$$

for sufficiently large  $N \geq N_0(\epsilon, D)$ , and we write  $\mathcal{X} \prec \mathcal{Y}$ . When  $\mathcal{X}^{(N)}$  and  $\mathcal{Y}^{(N)}$  depend on a parameter  $v \in \mathcal{V}$  (typically an index label or a spectral parameter), then  $\mathcal{X}(v) \prec \mathcal{Y}(v)$ , uniformly in  $v \in \mathcal{V}$ , means that the threshold  $N_0(\epsilon, D)$  can be chosen independently of  $v$ .

Motivated by (1.11) we introduce a probability measure  $\rho_\Sigma$  on  $\mathbb{C}$  by requiring

$$d\rho_\Sigma(w) = \frac{1}{2\pi} \Delta_w \left( \int_{\mathbb{R}} \log |s| d\mu_{\Sigma, |w|}(s) \right) d^2w, \quad w \in \mathbb{C}, \quad (1.13)$$

where

$$\mu_{\Sigma, r} := \mu_\Sigma^{\text{sym}} \boxplus \delta_r^{\text{sym}}, \quad r \geq 0, \quad (1.14)$$

and  $\Delta_w$  is the Laplacian on  $\mathbb{C}$  in the sense of distributions.

*Remark 1.7.* The fact that formula (1.13) defines a probability measure follows from previous work on the subject which we shortly summarize here.

Consider a non-commutative  $W^*$ -probability space  $(\mathcal{M}, \tau)$ , with  $\tau$  a trace. Let  $u$  be a Haar unitary element and let  $t = t^*$  be  $*$ -free from  $u$  and such that the distribution of  $t$ , i.e., its spectral measure, is given by  $\mu_\Sigma$ . Let  $\tilde{\mu}_{\Sigma, w}$  be the spectral measure of  $|ut - uid|$ , with  $id$  the unit in  $\mathcal{M}$  and  $w \in \mathbb{C}$ . Then the Brown measure for the product  $ut$  is given by the Riesz measure associated to the subharmonic function

$$\mathbb{C} \ni w \mapsto \int_{\mathbb{R}} \log |s| d\tilde{\mu}_{\Sigma, w}(s), \quad (1.15)$$

c.f., Section 2 of [27]. Haagerup and Larsen showed in Proposition 3.5 in [27] that  $\tilde{\mu}_{\Sigma, w} = \mu_{\Sigma, |w|}$ . Hence  $\rho_\Sigma$  in (1.13) can be characterized as the Brown measure of  $ut$  which by construction is a probability measure.

The main result of this paper is the following local single theorem in the bulk. Notice that (1.5) is not assumed, we only require that  $d_L(\mu_\Sigma, \mu_\sigma) \leq b$ , for some small constant  $b > 0$ , for  $N$  sufficiently large.

**Theorem 1.8.** Suppose that Assumption 1.1 holds. Let  $\mu_\sigma$  be a compactly supported probability measure on  $[0, \infty)$  which is supported at more than one point. Fix any (small)  $\tau > 0$  and define

$$\mathcal{R}_\sigma^\tau := \{w \in \mathbb{C} : r_- + \tau \leq |w| \leq r_+ - \tau\} \subset \mathcal{R}_\sigma, \quad (1.16)$$

where  $r_\pm \equiv r_\pm(\mu_\sigma)$  are given in (1.10). Then there exists a (small) constant  $b_0 > 0$  and  $N_0 \in \mathbb{N}$ , depending only on  $\mu_\sigma$  and  $S_+$ , such that whenever the Lévy distance  $d_L(\mu_\Sigma, \mu_\sigma)$  satisfies

$$\sup_{N \geq N_0} d_L(\mu_\Sigma, \mu_\sigma) \leq b, \quad (1.17)$$

for some  $b \leq b_0$ , then the following holds. Choose any  $w_0 \in \mathcal{R}_\sigma^\tau$ . Let  $f : \mathbb{C} \rightarrow \mathbb{R}$  be a smooth function such that  $\|f\|_\infty \leq C_0$  and  $f(z) = 0$  for all  $|z| \geq C_0$ , for some positive constant  $C_0$ . For  $\alpha \in (0, 1/2)$  set

$$f_{w_0}(w) := N^{2\alpha} f(N^\alpha(w - w_0)). \quad (1.18)$$

Then we have for any  $\alpha \in (0, 1/2)$  that the estimate

$$\left| \frac{1}{N} \sum_{i=1}^N f_{w_0}(\lambda_i(X)) - \int_{\mathcal{R}_\sigma} f_{w_0}(w) d\rho_\Sigma(w) \right| \prec N^{-1+2\alpha} \|\Delta f\|_{L^1(\mathbb{C})} \quad (1.19)$$

holds uniformly in  $f$  and in  $w_0 \in \mathcal{R}_\sigma^\tau$ , for  $N$  sufficiently large, depending on  $\tau$ ,  $S_+$ ,  $\mu_\sigma$  and  $C_0$ .

*Remark 1.9.* Note that we can choose  $\alpha$  in (1.19), almost as large as  $1/2$  in order to have an effective bound on the error term. Since the typical distance between the eigenvalues in the bulk of the ring  $\mathcal{R}_\sigma$  is of order  $N^{-1/2}$ , our result is optimal, both in terms of range of the exponent  $\alpha$  and the error term on the right side of (1.19). In particular, this improves the recent local single ring theorem of Benaych-Georges in [10] from scale  $(\log N)^{-1/4}$  to the optimal scale  $N^{-1/2+\epsilon}$ , for any small  $\epsilon > 0$ .

*Remark 1.10.* Theorem 1.8 holds with  $U, V$  being Haar distributed on either  $U(N)$  or on  $O(N)$ .

*Remark 1.11.* Note that  $w_0$  in Theorem 1.8 is chosen to be the (open) single ring  $\mathcal{R}_\sigma$ , in particular  $w_0$  stays away from the boundary of  $\mathcal{R}_\sigma$ . In case  $r_- = 0$ ,  $\mathcal{R}_\sigma$  is a punctuated disc. It has been proved in [25, 9] that there are no outliers at an order one distance from  $\mathcal{R}_\sigma$ .

Let  $f : \mathbb{C} \rightarrow \mathbb{R}$  be smooth and supported on  $\mathcal{R}_\sigma^\tau$ , for some (small)  $\tau > 0$ . Following the proof of Theorem 1.8 it is straightforward to verify that (1.19) also holds with  $\alpha = 0$  and  $f_{w_0}$  replaced with  $f$ , provided that the support of the function  $f$  stays away from the spectral edges, *i.e.*, is contained in  $\mathcal{R}_\sigma^\tau$ .

The following corollary of Theorem 1.8 expresses the speed of convergence in the single ring theorem on the macroscopic scale.

**Corollary 1.12.** *Under the conditions and with the notations of Theorem 1.8, we have that*

$$\left| \frac{1}{N} \sum_{i=1}^N f(\lambda_i(X)) - \int_{\mathcal{R}_\sigma} f(w) d\rho_\sigma(w) \right| \prec \|\Delta f\|_{L^1(\mathbb{C})} \left( \frac{1}{N} + b \right), \quad (1.20)$$

uniformly for any function  $f$  supported in  $\mathcal{R}_\sigma^\tau$  with a bound  $\|f\|_\infty \leq C_0$ , for  $N$  sufficiently large, depending on  $\tau$ ,  $S_+$ ,  $\mu_\sigma$  and  $C_0$ .

*Remark 1.13.* In (1.20) the measure  $\rho_\sigma$  is given by (1.11). By Theorem 4.4 and Corollary 4.5 of [27], the measure  $\rho_\sigma$  is absolutely continuous on  $\mathbb{C} \setminus \{0\}$  with respect to Lebesgue measure. Moreover, it satisfies  $\rho_\sigma(\{0\}) = \mu_\sigma(\{0\})$ . (In case  $\mu_\sigma(\{0\}) > 0$ , we have  $r_- = 0$ .) Note however that we have to exclude the point  $w = 0$  in our results since it is outside  $\mathcal{R}_\sigma$ .

*Remark 1.14.* Note that in Theorem 1.8 and Corollary 1.12 we do not require any regularity assumption on the measure  $\mu_\sigma$ , we even allow for atoms in  $\mu_\sigma$ . In particular, sending  $b \rightarrow 0$ , as  $N \rightarrow \infty$ , Corollary 1.12 also implies that Assumption 1.1 and (1.5) together imply  $d_L(\rho_\Sigma, \rho_\sigma) \rightarrow 0$ , as  $N \rightarrow \infty$ , thus removing the regularity condition (1.10) in the bulk from the single ring theorem, this answers a question in [26, Remark 2].

**1.2. Summary of previous results.** The first single ring theorem was established by Feinberg and Zee for a class of unitary invariant ensemble in [23], but without full rigor. The complete mathematical proof was given by Guionnet, Krishnapur and Zeitouni [26]; see also Remarks 1.4 and 1.14 for relaxing some conditions.

In spirit of the Wigner ensemble for the Hermitian case, the Ginibre ensemble can also be naturally extended by considering arbitrary i.i.d. entries; however, the unitary invariance property is lost in this generalization. Starting from the work of Girko [24], until the final result of Tao and Vu [35] with the least moment assumption, there have been many works devoted in proving circular law for general distribution. We refer to the survey [13] for more references in this direction. A prominent idea called *Hermitization* was introduced by Girko in [24]. This method translates spectral distribution problems of a non-Hermitian matrix to those of a Hermitian matrix (of double dimension), whose spectral properties can be studied with more established techniques.

Similarly to Wigner's original semicircle law, the single ring theorem establishes weak convergence of the spectral distribution, *i.e.*, it captures the density of eigenvalues on the global

scale. Since the typical distance between nearby eigenvalues is very small, of order  $N^{-1/2}$ , it is natural to ask whether the empirical density can also be approximated by the deterministic limit density on some local scale. Ideally, such *local law* should hold on the smallest possible scale, *i.e.*, just above the scale  $N^{-1/2}$ . In the Hermitian case, the local laws for Wigner and related ensembles have been extensively studied in the recent years, see *e.g.* [18] for a survey and references therein; the optimal local scale has been first achieved in [21].

With the aid of Girko's Hermitization, local laws for non-Hermitian matrices can be obtained via studying the local law for certain Hermitian matrices. With this strategy, the local circular law on optimal scale was established in the series of works Bourgade, Yau and Yin [14, 15] and Yin [38]. The first local single ring theorem was obtained by Benaych-Georges in [10], down to the scale  $(\log N)^{-\frac{1}{4}}$ , by proving the matrix subordination for Girko's Hermitization of  $X$  in (1.1), *c.f.*, (2.14). The strategy of matrix subordination was originally introduced by Kargin in [28] for proving a local law in the additive matrix model  $A + UBU^*$ , where  $A$  and  $B$  are deterministic Hermitian matrices and  $U$  is a Haar unitary. This additive model shares certain similarities with the Hermitization of the model  $X = U\Sigma V^*$ , but the latter has a block structure and thus we call it *block additive model* (*c.f.*, (4.2)). Recently, in [3, 4, 5], we obtained the local law of the additive model  $A + UBU^*$  on the optimal scale. The approach developed in these works opens up a path to treat the optimal local law in the block additive model, hence also sheds light on the optimal local single ring theorem. The key difference is that in the block additive model the Haar unitary matrices provide only a randomized  $U(N) \times U(N)$  symmetry instead of the full  $U(2N)$  symmetry. In particular, the coupling between the blocks is deterministic, so the mixing mechanism is much weaker. A more detailed overview of the proof strategy and the difficulties will be given in Section 4.2.

**1.3. Notational conventions.** We use the symbols  $O(\cdot)$  and  $o(\cdot)$  for the standard big-O and little-o notation. We use  $c$  and  $C$  to denote strictly positive constants that do not depend on  $N$ . Their values may change from line to line.

We denote by  $M_N(\mathbb{C})$  the set of  $N \times N$  matrices over  $\mathbb{C}$ . For  $A \in M_N(\mathbb{C})$ , we denote by  $\|A\|$  its operator norm and by  $\|A\|_2$  its Hilbert-Schmidt norm. The matrix entries of  $A$  are denoted by  $A_{ij}$ .

Let  $\mathbf{g} = (g_1, \dots, g_N)$  be a real or complex Gaussian vector. We write  $\mathbf{g} \sim \mathcal{N}_{\mathbb{R}}(0, \sigma^2 I_N)$  if  $g_1, \dots, g_N$  are independent and identically distributed (i.i.d.)  $N(0, \sigma^2)$  normal variables; and we write  $\mathbf{g} \sim \mathcal{N}_{\mathbb{C}}(0, \sigma^2 I_N)$  if  $g_1, \dots, g_N$  are i.i.d.  $N_{\mathbb{C}}(0, \sigma^2)$  variables, where  $g_i \sim N_{\mathbb{C}}(0, \sigma^2)$  means that  $\operatorname{Re} g_i$  and  $\operatorname{Im} g_i$  are independent  $N(0, \frac{\sigma^2}{2})$  normal variables.

We use double brackets to denote index sets, *i.e.*, for  $n_1, n_2 \in \mathbb{R}$ ,  $[[n_1, n_2]] := [n_1, n_2] \cap \mathbb{Z}$ .

*Acknowledgment:* Part of this work was accomplished when Z.-G. B. and K. S. were working at IST Austria with the support of ERC Advanced Grant RANMAT No. 338804. Support and hospitality are gratefully acknowledged. We thank an anonymous referee for very useful comments and suggestions.

## 2. PRELIMINARIES AND MAIN TECHNICAL TASK

**2.1. Free additive convolution.** We recall some basic notions and results for the free additive convolution. We follow the notational conventions in our previous paper [2].

Let  $\mu$  be a Borel probability measure on  $\mathbb{R}$  and recall its Stieltjes transform  $m_\mu$  defined in (1.9). Note that  $m_\mu : \mathbb{C}^+ \rightarrow \mathbb{C}^+$  is an analytic function such that

$$\lim_{\eta \nearrow \infty} i\eta m_\mu(i\eta) = -1. \quad (2.1)$$

Conversely, if  $m : \mathbb{C}^+ \rightarrow \mathbb{C}^+$  is an analytic function such that  $\lim_{\eta \nearrow \infty} i\eta m(i\eta) = -1$ , then  $m$  is the Stieltjes transform of a probability measure  $\mu$ .

Given a Borel probability measure  $\mu$  on  $\mathbb{R}$ , let  $F_\mu$  be the *negative reciprocal Stieltjes transform* of  $\mu$ ,

$$F_\mu(z) := -\frac{1}{m_\mu(z)}, \quad z \in \mathbb{C}^+. \quad (2.2)$$

Observe that

$$\lim_{\eta \nearrow \infty} \frac{F_\mu(i\eta)}{i\eta} = 1, \quad (2.3)$$

as follows from (2.1). Note that  $F_\mu$  is analytic on  $\mathbb{C}^+$  with nonnegative imaginary part.

The *free additive convolution* is the symmetric binary operation on Borel probability measures on  $\mathbb{R}$  characterized by the following result.

**Theorem 2.1** (Theorem 4.1 in [8], Theorem 2.1 in [16]). *Given two Borel probability measures,  $\mu_1$  and  $\mu_2$ , on  $\mathbb{R}$ , there exist unique analytic functions,  $\omega_1, \omega_2 : \mathbb{C}^+ \rightarrow \mathbb{C}^+$ , such that,*

(i) *for all  $z \in \mathbb{C}^+$ ,  $\text{Im } \omega_1(z), \text{Im } \omega_2(z) \geq \text{Im } z$ , and*

$$\lim_{\eta \nearrow \infty} \frac{\omega_1(i\eta)}{i\eta} = \lim_{\eta \nearrow \infty} \frac{\omega_2(i\eta)}{i\eta} = 1; \quad (2.4)$$

(ii) *for all  $z \in \mathbb{C}^+$ ,*

$$F_{\mu_1}(\omega_2(z)) = F_{\mu_2}(\omega_1(z)), \quad \omega_1(z) + \omega_2(z) - z = F_{\mu_1}(\omega_2(z)). \quad (2.5)$$

It follows from (2.4) that the analytic function  $F : \mathbb{C}^+ \rightarrow \mathbb{C}^+$  defined by

$$F(z) := F_{\mu_1}(\omega_2(z)) = F_{\mu_2}(\omega_1(z)), \quad (2.6)$$

satisfies the analogue of (2.3). Thus  $F$  is the negative reciprocal Stieltjes transform of a probability measure  $\mu$ , called the free additive convolution of  $\mu_1$  and  $\mu_2$ , denoted by  $\mu \equiv \mu_1 \boxplus \mu_2$ . The functions  $\omega_1$  and  $\omega_2$  are referred to as the *subordination functions* and  $F$  is said to be subordinated to  $F_{\mu_1}$ , respectively to  $F_{\mu_2}$ . The subordination phenomenon was first noted by Voiculescu [37] in a generic situation and extended to full generality by Biane [12]. To exclude trivial shifts of measures, we henceforth assume that both,  $\mu_1$  and  $\mu_2$ , are supported at more than one point. Then the analytic functions  $F$ ,  $\omega_1$  and  $\omega_2$  extend continuously to the real line; see Theorem 2.3 [6] or Theorem 3.3 [7]. We use the same notation for their extensions to  $\mathbb{C}^+ \cup \mathbb{R}$ .

**2.2. The limiting measure  $\mu_{\sigma,r}$ .** Recall the definitions  $\mu_{\Sigma,r} := \mu_{\Sigma}^{\text{sym}} \boxplus \delta_r^{\text{sym}}$  and  $\mu_{\sigma,r} := \mu_{\sigma}^{\text{sym}} \boxplus \delta_r^{\text{sym}}$  from (1.8). In this subsection, we will always assume that  $\mu_{\Sigma}$  and  $\mu_{\sigma}$  satisfy Assumption 1.1. For sake of simplicity of notation, we abbreviate in this subsection

$$\mu_1 \equiv \mu_{\sigma}^{\text{sym}}, \quad \mu_2 \equiv \delta_r^{\text{sym}}. \quad (2.7)$$

The negative reciprocal Stieltjes transform of  $\mu_2 = \delta_r^{\text{sym}}$  is found to be

$$F_{\mu_2}(z) = z - \frac{r^2}{z}, \quad z \in \mathbb{C}^+. \quad (2.8)$$

Substituting (2.8) into (2.5), we obtain

$$F_{\mu_1}(\omega_2(z)) = F_{\mu_2}(\omega_1(z)) = F_{\mu_1}(\omega_2(z)) - \omega_2(z) + z - \frac{r^2}{F_{\mu_1}(\omega_2(z)) - \omega_2(z) + z}.$$

Solving the above equation for  $F_{\mu_1}(\omega_2(z))$  we conclude that the subordination function  $\omega_2(z)$  is the unique solution to

$$F_{\mu_1}(\omega_2(z)) - \omega_2(z) = -z - \frac{r^2}{\omega_2(z) - z}, \quad z \in \mathbb{C}^+, \quad (2.9)$$

subject to the condition  $\text{Im } \omega_2(z) \geq \text{Im } z$ . Comparing once more with (2.5) we immediately find that the other subordination function is given by

$$\omega_1(z) = -\frac{r^2}{\omega_2(z) - z}, \quad z \in \mathbb{C}^+. \quad (2.10)$$

The analysis of the measure  $\mu_{\sigma,r} = \mu_1 \boxplus \mu_2$  thus reduces to the analysis of (2.9) for  $\omega_2$ . We first derive upper and lower bound on  $\omega_2(z)$ . For the purpose of proving Theorem 1.8 it will suffice to consider  $z \in \{i\eta : \eta \geq 0\}$ . Since  $\mu_1$  and  $\mu_2$  are symmetric, we have  $\omega_2(i\eta) = -\overline{\omega_2(i\eta)}$ , i.e.,  $\omega_2(i\eta)$  and  $\omega_1(i\eta)$  are both fully imaginary. This simplifies our analysis; while detailed quantitative properties of the full measure  $\mu_{\sigma,r}$  are still poorly understood, we now have a good control on it near zero, hence on its Stieltjes transform along the imaginary axis. The main result, formulated in Theorem 2.2 below, is that the subordination functions are bounded from below and above on the imaginary axis without any condition on  $\mu_\sigma$ . This theorem is the key input that enables us to dispense with the regularity condition in the single ring theorem; see Remark 1.14.

**Theorem 2.2** (Bounds on subordination functions). *We assume that the support of  $\mu_\sigma$  contains more than one point, equivalently, that  $r_- < r_+$ . Let  $\mu_1 = \mu_\sigma^{\text{sym}}$  and  $\mu_2 = \delta_r^{\text{sym}}$  for some  $r > 0$ . Fix  $\eta_M < \infty$  and a (small)  $\tau > 0$ . Set*

$$J := [r_- + \tau, r_+ - \tau].$$

*There exist constants  $c \equiv c(\mu_1, \tau, \eta_M) > 0$  and  $C \equiv C(\mu_1, \tau, \eta_M) < \infty$  such that*

$$\sup_{r \in J} \sup_{\eta \in [0, \eta_M]} |\omega_1(i\eta)| \leq C, \quad \sup_{r \in J} \sup_{\eta \in [0, \eta_M]} |\omega_2(i\eta)| \leq C, \quad (2.11)$$

$$\inf_{r \in J} \inf_{\eta \in [0, \eta_M]} \text{Im } \omega_1(i\eta) \geq c, \quad \inf_{r \in J} \inf_{\eta \in [0, \eta_M]} \text{Im } \omega_2(i\eta) \geq c, \quad (2.12)$$

and

$$\inf_{r \in J} \inf_{\eta \in [0, \eta_M]} |m_{\mu_1 \boxplus \mu_2}(i\eta)| \geq c, \quad \sup_{r \in J} \sup_{\eta \in [0, \eta_M]} |m_{\mu_1 \boxplus \mu_2}(i\eta)| \leq C. \quad (2.13)$$

*Remark 2.3.* By (2.13), the measure  $\mu_1 \boxplus \mu_2$  has a positive and bounded density at  $E = 0$ . In particular,  $E = 0$  is in the bulk of the measure  $\mu_1 \boxplus \mu_2$ , as defined in Definition 4.2 below.

The proof of Theorem 2.2 is quite technical and independent of the main line of the argument, so we give it in Section 7. In the subsequent sections, we will mainly rely on the following corollary of Theorem 2.2. Let  $m_{\Sigma,r}(z)$  be the Stieltjes transform of  $\mu_{\Sigma,r}$ ; see (1.8).

**Corollary 2.4.** *Fix  $\eta_M < \infty$  and a (small)  $\tau > 0$ . Then there are constants  $C \equiv C(\mu_\sigma^{\text{sym}}, \tau, \eta_M)$ ,  $c \equiv c(\mu_\sigma^{\text{sym}}, \tau, \eta_M)$  and a threshold  $N_0 \equiv N_0(\mu_\sigma^{\text{sym}}, \tau, \eta_M)$  such that the conclusions in Theorem 2.2 hold with  $\mu_1 = \mu_\Sigma^{\text{sym}}$  and  $\mu_2 = \delta_r^{\text{sym}}$ , for  $N \geq N_0$ .*

*Proof.* This follows directly from the continuity of the subordination functions with respect to the Lévy distance (see Lemma 5.1 of [2]), from Theorem 2.2 and from (1.17).  $\square$

**2.3. Key technical inputs.** Following Girko's hermitization technique [24], we introduce for any  $w \in \mathbb{C}$  the  $2N \times 2N$  Hermitian matrix

$$H^w := \begin{pmatrix} 0 & X - w \\ X^* - w^* & 0 \end{pmatrix}. \quad (2.14)$$

The main advantage of working with  $H^w$  is that it is self-adjoint and we thus have a functional calculus at disposal. For any function  $g \in C^2(\mathbb{C})$ , an application of Green's theorem reveals that

$$\frac{1}{N} \sum_{i=1}^N g(\lambda_i(X)) = \frac{1}{2\pi} \int_{\mathbb{C}} (\Delta g)(w) \left( \frac{1}{2N} \text{Tr } \log |H^w| \right) d^2 w, \quad (2.15)$$

which is a manifestation of  $\log |\cdot|$  being the Coulomb potential in two dimensions. The following identity, first used in this context by [36], allows us to efficiently deal with the right side of (2.15). For any (large)  $K > 0$ ,

$$\frac{1}{2N} \text{Tr } \log |H^w| = \frac{1}{2N} \text{Tr } \log |(H^w - iK)| - \text{Im} \int_0^K m^w(i\eta) d\eta, \quad (2.16)$$



with  $|w| > 0$ , where  $m^w(z)$ ,  $z \in \mathbb{C}^+$ , is the Stieltjes transform of the spectral distribution of  $H^w$ . For very large  $K$  the first term on the right side of (2.16) is elementary to control, we hence focus on the second term. Due to the block structure of  $H^w$ , the eigenvalues come in pairs  $\pm\lambda_i^w$ ,  $i \in \llbracket 1, N \rrbracket$ , where  $0 \leq \lambda_1^w \leq \dots \leq \lambda_N^w$  are the non-negative eigenvalues. With these notations  $m^w$  is given by

$$m^w(z) := \frac{1}{2N} \sum_{i=1}^N \left( \frac{1}{\lambda_i^w - z} + \frac{1}{-\lambda_i^w - z} \right) = \frac{1}{N} \sum_{i=1}^N \frac{\lambda_i^w}{(\lambda_i^w)^2 - z^2}, \quad z \in \mathbb{C}^+.$$

Recall the notation  $m_{\Sigma, |w|}$  for the Stieltjes transform of  $\mu_{\Sigma, |w|}$ ; *c.f.*, (1.8). The following result is the main technical input for the proof of Theorem 1.8. Recall  $\mathcal{R}_\sigma^\tau$  from (1.16).

**Theorem 2.5** (Local law for  $H^w$ ). *Under the conditions and with the notations of Theorem 1.8, the estimate*

$$\sup_{w \in \mathcal{R}_\sigma^\tau} |m^w(i\eta) - m_{\Sigma, |w|}(i\eta)| \prec \frac{1}{N\eta}, \quad (2.17)$$

holds uniformly in  $\eta > 0$ , for  $N$  sufficiently large, depending on  $\tau$ ,  $S_+$  and  $\mu_\sigma$ .

This result controls  $|m^w(i\eta) - m_{\Sigma, |w|}(i\eta)|$  along the positive imaginary axis. Note that the error estimate on the right side of (2.17) is effective when  $\eta$  is chosen just above the local scale, *i.e.*, when  $\eta > N^{-1+\gamma}$ , for any small  $\gamma > 0$ . For even smaller  $\eta > 0$ , (2.17) yields the upper bound  $|m^w(i\eta)| \prec (N\eta)^{-1}$  which improves the trivial deterministic bound  $|m^w(i\eta)| \leq \eta^{-1}$  by a factor  $N^{-1}$ . Theorem 2.5 is used to control the integrand in the second term on the right side of (2.16) for  $\eta \gtrsim N^{-1}$ . On very short scales, the behavior of  $m^w(i\eta)$ ,  $\eta \lesssim N^{-1}$ , is essentially random and determined by the smallest (in absolute value) eigenvalues of  $H^w$ . The following estimate on  $\lambda_1^w$ , proved by Rudelson and Vershynin in [34], is then used to control the integrand of the second term on the right side of (2.16) for very small  $\eta \lesssim N^{-1}$ .

**Theorem 2.6** (Theorem 1.1 and Theorem 1.2 in [34]). *There exist positive numerical constants  $c > 0$  and  $C < \infty$ , such that*

$$\mathbb{P}\left(\lambda_1^w \leq \frac{t}{|w|}\right) \leq \left(\frac{t}{|w|}\right)^c N^C, \quad (2.18)$$

uniformly in  $t > 0$ , for all  $N \in \mathbb{N}$ .

*Remark 2.7.* In the orthogonal case, (2.18) holds, for  $N$  sufficiently large, when the matrix  $\Sigma$  is away from the identity; see Theorem 1.2 in [34]. In this case the constants  $c$ ,  $C$  and the threshold for  $N$  in (2.18) depend on  $S_+$  and  $\mu_\sigma$ . Indeed, (1.17) and the assumption that the support of  $\mu_\sigma$  contains more than one point imply that  $\Sigma$  is separated away from the identity.

In Section 3, we will choose  $g$  in (2.15) to be the rescaled function  $f_{w_0}(\cdot) = N^{2\alpha} f(N^\alpha(\cdot - w_0))$ ; see (1.18). The local law in (2.17) together with (2.18) (with  $t/|w| \ll N^{-1}$ ) will allow us to choose  $\alpha \in (0, 1/2)$  as is asserted in Theorem 1.8. The details of the proof of Theorem 1.8, assuming Theorem 2.5, are carried out in Section 3. Our main task then is to prove Theorem 2.5. Actually, we will establish the local law in a more general setting; *c.f.*, Theorem 4.3. This will be accomplished in Sections 4-6 and we will separately outline the main ideas of this proof in Section 4.2. We begin with the proof of Theorem 2.2 in the next section.

### 3. PROOF OF THEOREM 1.8 AND COROLLARY 1.12

In this section, we prove Theorem 1.8 and Corollary 1.12, with the aid of Theorems 2.5 and 2.6. The use of Girko's hermitized matrices to derive local laws is a standard argument, see *e.g.*, [14, 36] for related models. Following [36], we use the identity (3.6) below to link the log-determinant of  $H^w$  with the Stieltjes transform  $m^w$ .



*Proof of Theorem 1.8.* For any  $\zeta \in \mathbb{C}$ , we denote

$$w \equiv w(\zeta) := w_0 + N^{-\alpha}\zeta. \quad (3.1)$$

Given  $f : \mathbb{C} \rightarrow \mathbb{R}$  satisfying the assumption of Theorem 1.8, we introduce the domain

$$\mathcal{D}_{w_0}(\alpha) \equiv \mathcal{D}_{w_0}(\alpha, f) := \left\{ \tilde{w} : N^\alpha(\tilde{w} - w_0) \in \text{supp}(f) \right\}. \quad (3.2)$$

According to (3.1),  $w \in \mathcal{D}_{w_0}(\alpha)$  is equivalent to  $\zeta \in \text{supp}(f)$ , in particular  $|\zeta| \leq C$  as  $f$  is compactly supported. Recall the notation  $f_{w_0}(\cdot)$  from Theorem 1.8. Using (2.15), we rewrite

$$\frac{1}{N} \sum_i f_{w_0}(\lambda_i(X)) = \frac{1}{2\pi} N^{2\alpha} \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \frac{1}{2N} \text{Tr} \log |H^w| \right) d^2\zeta. \quad (3.3)$$

Recalling the definitions in (1.8) and (1.13), we also have

$$\begin{aligned} \int_{\mathbb{C}} f_{w_0}(w) \rho_{\Sigma}(d^2w) &= \frac{1}{2\pi} \int_{\mathbb{C}} f_{w_0}(w) \Delta_w \left( \int_{\mathbb{R}} \log |u| \mu_{\Sigma, |w|}(du) \right) d^2w \\ &= \frac{1}{2\pi} N^{2\alpha} \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \int_{\mathbb{R}} \log |u| \mu_{\Sigma, |w|}(du) \right) d^2\zeta. \end{aligned} \quad (3.4)$$

Hence, we can write

$$\begin{aligned} \frac{1}{N} \sum_i f_{w_0}(\lambda_i(X)) - \int_{\mathbb{C}} f_{w_0}(w) \rho_{\Sigma}(d^2w) \\ = \frac{1}{2\pi} N^{2\alpha} \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \frac{1}{2N} \text{Tr} \log |H^w| - \int_{\mathbb{R}} \log |u| \mu_{\Sigma, |w|}(du) \right) d^2\zeta. \end{aligned} \quad (3.5)$$

We next use the following observation due to [36], Section 8. For any (large)  $K > 0$  and  $|w| > 0$ , we have

$$\frac{1}{2N} \text{Tr} \log |H^w| = \frac{1}{2N} \text{Tr} \log |(H^w - iK)| - \text{Im} \int_0^K m^w(i\eta) d\eta. \quad (3.6)$$

Analogously, we can also write, with the same  $K$ ,

$$\int_{\mathbb{R}} \log |u| \mu_{\Sigma, |w|}(du) = \int_{\mathbb{R}} \log |u - iK| \mu_{\Sigma, |w|}(du) - \text{Im} \int_0^K m_{\Sigma, |w|}(i\eta) d\eta. \quad (3.7)$$

Choosing  $K$  sufficiently large, say  $K = N^L$  for some large constant  $L$ , it is easy to see that

$$\left| \frac{1}{2N} \text{Tr} \log |(H^w - iK)| - \int_{\mathbb{R}} \log |u - iK| \mu_{\Sigma, |w|}(du) \right| \ll \frac{1}{N} \quad (3.8)$$

holds uniformly in  $w \in \mathcal{D}_{w_0}(\alpha)$ . Here we used the fact that  $\|H^w\| \leq C$  for some positive constant  $C$ , under *c.f.*, Assumption 1.1. The uniformity in  $w$  can be guaranteed by the fact that  $\mathcal{D}_{w_0}(\alpha)$  lies in a ball of finite (in fact  $CN^{-\alpha}$ ) radius since  $f$  is compactly supported. Hence, it suffices to show

$$\left| \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \text{Im} \int_0^{N^L} (m^w(i\eta) - m_{\Sigma, |w|}(i\eta)) d\eta \right) d^2\zeta \right| \prec \frac{\|\Delta f\|_{L^1(\mathbb{C})}}{N}. \quad (3.9)$$

To show (3.9), we decompose the integral with respect to  $\eta$  into two parts:

$$\int_0^{N^L} = \int_0^{N^{-L_1}} + \int_{N^{-L_1}}^{N^L}, \quad (3.10)$$

for sufficiently large constants  $L_1 > 1$  and  $L > 0$  to be chosen below. To control the first part, we use (2.18), while for the second part we use (2.17).

First, using the upper bound of  $m_{\Sigma, |w|}(i\eta)$  (*c.f.*, Corollary 2.4), we obtain

$$\left| \int_0^{N^{-L_1}} \text{Im} m_{\Sigma, |w|}(i\eta) d\eta \right| \leq \frac{1}{N}, \quad (3.11)$$

for  $L_1 > 1$ , uniformly in  $w \in \mathcal{D}_{w_0}(\alpha)$ . Hence, we have

$$\left| \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \int_0^{N^{-L_1}} \operatorname{Im} m_{\Sigma, |w|}(i\eta) d\eta \right) d^2\zeta \right| \leq C \frac{\|\Delta f\|_{L^1(\mathbb{C})}}{N}. \quad (3.12)$$

In addition, we observe that

$$\begin{aligned} & \mathbb{P} \left( \left| \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \int_0^{N^{-L_1}} \operatorname{Im} m^w(i\eta) d\eta \right) d^2\zeta \right| > \frac{\|\Delta f\|_{L^1(\mathbb{C})}}{N} \right) \\ & \leq \frac{N}{\|\Delta f\|_{L^1(\mathbb{C})}} \mathbb{E} \left| \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \int_0^{N^{-L_1}} \operatorname{Im} m^w(i\eta) d\eta \right) d^2\zeta \right| \\ & \leq \frac{N}{\|\Delta f\|_{L^1(\mathbb{C})}} \int_{\mathbb{C}} |(\Delta f)(\zeta)| \mathbb{E} \left( \int_0^{N^{-L_1}} \frac{\eta}{(\lambda_1^w)^2 + \eta^2} d\eta \right) d^2\zeta. \end{aligned} \quad (3.13)$$

Note that

$$\begin{aligned} & \mathbb{E} \left( \int_0^{N^{-L_1}} \frac{\eta}{(\lambda_1^w)^2 + \eta^2} d\eta \right) = \frac{1}{2} \mathbb{E} \log \left( 1 + (N^{L_1} \lambda_1^w)^{-2} \right) \\ & = \frac{1}{2} \int_0^\infty \mathbb{P} \left( \log \left( 1 + (N^{L_1} \lambda_1^w)^{-2} \right) \geq s \right) ds \\ & = \frac{1}{2} \int_0^\infty \mathbb{P} \left( \lambda_1^w \leq N^{-L_1} (e^s - 1)^{-\frac{1}{2}} \right) ds \\ & = \frac{1}{2} \left( \int_0^{N^{-L_1}} + \int_{N^{-L_1}}^1 + \int_1^\infty \right) \mathbb{P} \left( \lambda_1^w \leq N^{-L_1} (e^s - 1)^{-\frac{1}{2}} \right) ds. \end{aligned}$$

For the first integral, we use the trivial bound  $\mathbb{P}(\cdot) \leq 1$  to obtain

$$\int_0^{N^{-L_1}} \mathbb{P} \left( \lambda_1^w \leq N^{-L_1} (e^s - 1)^{-\frac{1}{2}} \right) ds \leq N^{-L_1}. \quad (3.14)$$

For the second part of the integral, using the crude bound  $(e^s - 1)^{-\frac{1}{2}} \leq s^{-\frac{1}{2}} \leq N^{\frac{L_1}{2}}$ ,  $s \in [N^{-L_1}, 1]$ , and (2.18), we estimate

$$\int_{N^{-L_1}}^1 \mathbb{P} \left( \lambda_1^w \leq N^{-L_1} (e^s - 1)^{-\frac{1}{2}} \right) ds \leq \int_{N^{-L_1}}^1 \mathbb{P} \left( \lambda_1^w \leq N^{-\frac{L_1}{2}} \right) ds \leq N^{-\frac{cL_1}{2} + C},$$

for some constants  $c > 0$  and  $C < \infty$ , for  $N$  sufficiently large. For the third part, using  $e^s - 1 > \frac{1}{2}e^s$ ,  $s > 1$ , and (2.18), we have

$$\begin{aligned} \int_1^\infty \mathbb{P} \left( \lambda_1^w \leq N^{-L_1} (e^s - 1)^{-\frac{1}{2}} \right) ds & \leq \int_1^\infty \mathbb{P} \left( \lambda_1^w \leq \sqrt{2} N^{-L_1} e^{-\frac{s}{2}} \right) ds \\ & \leq \frac{N^{-cL_1 + C}}{2} \int_1^\infty e^{-\frac{cs}{2}} ds \leq N^{-cL_1 + C}, \end{aligned} \quad (3.15)$$

for some constants  $c > 0$  and  $C < \infty$ . Combining (3.14)-(3.15), we obtain that there are positive constants  $c' > 0$  and  $C'$ , independent of  $L_1$  such that

$$\mathbb{E} \left( \int_0^{N^{-L_1}} \frac{\eta}{(\lambda_1^w)^2 + \eta^2} d\eta \right) \leq N^{-c'L_1 + C'}, \quad (3.16)$$

for  $N$  sufficiently large. In fact, the bound (3.16) is uniform in  $w \in \mathcal{D}_{w_0}(\alpha)$  since the constants  $c$  and  $C$  in Theorem 2.6 are uniform in  $t$  and  $w$ . Plugging (3.16) into (3.13), yields

$$\mathbb{P} \left( \left| \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \int_0^{N^{-L_1}} \operatorname{Im} m^w(i\eta) d\eta \right) d^2\zeta \right| \geq \frac{\|\Delta f\|_{L^1(\mathbb{C})}}{N} \right) \leq N^{-c'L_1 + C' + 1}, \quad (3.17)$$

for  $N$  sufficiently large (independent of  $L_1$ ). Choosing  $L_1$  large enough, the contribution of the first integral in (3.10) to (3.9) is within the claimed error.

To control the contributions from the second integral in (3.10), for any (large) constant  $L_1$ , we apply the local law for  $m^w$  in (2.17), uniform in  $w$ , to find

$$\begin{aligned} & \left| \int_{\mathbb{C}} (\Delta f)(\zeta) \left( \operatorname{Im} \int_{N^{-L_1}}^{N^L} (m^w(i\eta) - m_{\Sigma,|w|}(i\eta)) d\eta \right) d^2\zeta \right| \\ & \prec \int_{\mathbb{C}} |(\Delta f)(\zeta)| \left( \int_{N^{-L_1}}^{N^L} \frac{1}{N\eta} d\eta \right) d^2\zeta \prec \frac{\|\Delta f\|_{L^1(\mathbb{C})}}{N}. \end{aligned}$$

Combining (3.12) and (3.17), and choosing  $L_1$  sufficiently large, we get (3.9), which together with (3.5)-(3.8) concludes the proof of Theorem 1.8.  $\square$

*Proof of Corollary 1.12.* Let  $f : \mathbb{C} \rightarrow \mathbb{R}$  be smooth and supported on  $\mathcal{R}_\sigma^\tau$ ; see (1.16). It is straightforward following the proof of Theorem 1.8 to verify that (1.19) also holds with  $\alpha = 0$  and  $f_{w_0}$  replaced with  $f$  provided that  $\operatorname{supp} f \subset \mathcal{R}_\sigma^\tau$ ; *c.f.*, Remark 1.11. Thus under the assumptions of Corollary 1.12 it suffices to show that

$$\left| \int_{\mathcal{R}_\sigma^\tau} (\Delta f)(w) \int_{\mathbb{R}} \log |u| (\mu_{\Sigma,|w|}(du) - \mu_{\sigma,|w|}(du)) d^2w \right| \leq C \|\Delta f\|_{L^1(\mathbb{C})} d_L(\mu_\Sigma, \mu_\sigma),$$

for a constant  $C$  (depending on  $\tau$ ), to conclude its proof. From (3.7), it is sufficient to prove that

$$\left| \int_{\mathbb{R}} \log |u - i| (\mu_{\Sigma,|w|}(du) - \mu_{\sigma,|w|}(du)) \right| \leq C d_L(\mu_\Sigma, \mu_\sigma) \quad (3.18)$$

and

$$\left| \int_0^1 (m_{\Sigma,|w|}(i\eta) - m_{\sigma,|w|}(i\eta)) d\eta \right| \leq C d_L(\mu_\Sigma, \mu_\sigma), \quad (3.19)$$

uniformly for all  $w \in \mathcal{R}_\sigma^\tau$ , for  $N$  sufficiently large.

Inequality (3.18) follows from the continuity of the additive free convolution. More precisely, from Theorem 4.13 of [11], we know that  $d_L(\mu_{\Sigma,|w|}, \mu_{\sigma,|w|}) \leq d_L(\mu_\Sigma, \mu_\sigma)$ . Since  $\log |u - i|$  is a smooth function and  $\mu_{\Sigma,|w|}, \mu_{\sigma,|w|}$  are compactly supported, (3.18) follows.

To establish (3.19), we note that, for  $N$  sufficiently large,

$$\int_0^1 |m_{\Sigma,|w|}(i\eta) - m_{\sigma,|w|}(i\eta)| d\eta \leq \max_{\eta \in (0,1]} |m_{\Sigma,|w|}(i\eta) - m_{\sigma,|w|}(i\eta)| \leq C d_L(\mu_\Sigma, \mu_\sigma),$$

for all  $w$  with  $r_- + \tau \leq |w| \leq r_+ - \tau$ , with a constant depending on  $\tau$ . This follows directly from Theorem 2.7 of [2]. This shows (3.19).

So far we proved (1.20) for smooth functions  $f$ . Since  $\rho_\sigma$  is a Borel probability measure, see *e.g.*, Theorem 1.2, (1.20) extends to  $f \in C^2(\mathbb{C})$  supported in  $\mathcal{R}_\sigma^\tau$ . This completes the proof of Corollary 1.12.  $\square$

#### 4. LOCAL LAW FOR BLOCK ADDITIVE MODEL

In this section, we derive a local law for block additive random matrices in a slightly generalized setting; see Theorem 4.3 below. Theorem 2.5 is a direct consequence of this result.

First, note that the matrix  $H^w$  defined in (2.14) can be rewritten as

$$H^w = \begin{pmatrix} U & 0 \\ 0 & V \end{pmatrix} \begin{pmatrix} 0 & \Sigma \\ \Sigma & 0 \end{pmatrix} \begin{pmatrix} U^* & 0 \\ 0 & V^* \end{pmatrix} + \begin{pmatrix} 0 & -w \\ -w^* & 0 \end{pmatrix}, \quad (4.1)$$

where 0 is the  $N \times N$  matrix filled with zeros. In the following we consider a slightly more general problem by looking at random matrices  $H$  defined by

$$H := \begin{pmatrix} U & 0 \\ 0 & V \end{pmatrix} \begin{pmatrix} 0 & \Sigma \\ \Sigma^* & 0 \end{pmatrix} \begin{pmatrix} U^* & 0 \\ 0 & V^* \end{pmatrix} + \begin{pmatrix} 0 & \Xi \\ \Xi^* & 0 \end{pmatrix}, \quad (4.2)$$

where

$$\Sigma := \operatorname{diag}(\sigma_1, \dots, \sigma_N), \quad \Xi := \operatorname{diag}(\xi_1, \dots, \xi_N), \quad (4.3)$$

with  $\sigma_i, \xi_i \in \mathbb{C}$ ,  $i \in \llbracket 1, N \rrbracket$ . Here  $\Sigma$  and  $\Xi$  are deterministic diagonal matrices, while  $U$  and  $V$  are independent Haar unitary or Haar orthogonal matrices of degree  $N$  as before. Note that we allow in (4.3) for complex matrix elements in  $\Sigma$  and  $\Xi$ . In the sequel, we always assume that  $\Sigma$  and  $\Xi$  are bounded,

$$\|\Sigma\|, \|\Xi\| \leq C, \quad (4.4)$$

for some constant  $C$  independent of  $N$ . Denote the empirical density of their singular values by

$$\mu_\Sigma := \frac{1}{N} \sum_{i=1}^N \delta_{|\sigma_i|}, \quad \mu_\Xi := \frac{1}{N} \sum_{i=1}^N \delta_{|\xi_i|}. \quad (4.5)$$

Note that  $\mu_\Sigma$  and  $\mu_\Xi$  are probability measures on  $[0, \infty)$ . We assume that there are compactly supported probability measures  $\mu_\sigma$  and  $\mu_\xi$  such that

$$\sup_{N \geq N_0} (\mathrm{d}_L(\mu_\Sigma, \mu_\sigma) + \mathrm{d}_L(\mu_\Xi, \mu_\xi)) \leq 2b, \quad (4.6)$$

for a sufficiently small constant  $b > 0$  and sufficiently large  $N_0$ .

The following general regularity result is of interest.

**Lemma 4.1** (Theorem 4.1 in [7]). *Let  $\mu_1$  and  $\mu_2$  be Borel probability measures on  $\mathbb{R}$ , neither of them a point mass. Then the singular continuous part of  $\mu_1 \boxplus \mu_2$  vanishes. A point  $x \in \mathbb{R}$  is an atom of  $\mu_1 \boxplus \mu_2$  if and only if there are  $x_1, x_2 \in \mathbb{R}$  such that  $x = x_1 + x_2$  and  $\mu_1(\{x_1\}) + \mu_2(\{x_2\}) > 1$ . Moreover, the absolutely continuous part of  $\mu_1 \boxplus \mu_2$  is always nonzero, and its density is analytic wherever positive and finite.*

**Definition 4.2.** *For two Borel probability measures  $\mu_1$  on  $\mu_2$  on  $\mathbb{R}$  satisfying the assumptions of Lemma 4.1, we set*

$$\mathcal{B}_{\mu_1 \boxplus \mu_2} := \{x \in \mathbb{R} : 0 < f_{\mu_1 \boxplus \mu_2}(x) < \infty, \mu_1 \boxplus \mu_2(\{x\}) = 0\}, \quad (4.7)$$

where  $f_{\mu_1 \boxplus \mu_2}$  denotes the density function of  $\mu_1 \boxplus \mu_2$ . We call  $\mathcal{B}_\mu$  the bulk of  $\mu$ .

Let  $G \equiv G(z) := (H - z)^{-1}$  be the Green function of  $H$  at parameter  $z \in \mathbb{C}^+$ , and let

$$m_H(z) := \mathrm{tr} G(z) = \frac{1}{2N} \mathrm{Tr} G(z) \quad (4.8)$$

be the normalized trace of  $G(z)$ , which by the functional calculus agrees with the Stieltjes transform of the empirical eigenvalue distribution of  $H$ .

Given an interval  $\mathcal{I} \subset \mathbb{R}$  and  $0 \leq a \leq b$ , we introduce the domain

$$\mathcal{S}_{\mathcal{I}}(a, b) := \{z = E + i\eta \in \mathbb{C}^+ : E \in \mathcal{I}, a < \eta \leq b\}. \quad (4.9)$$

As before, we denote for a measure  $\mu$  on  $\mathbb{R}$  its symmetrization by  $\mu^{\mathrm{sym}}$ . The following is a key result of this paper.

**Theorem 4.3** (Strong law for  $H$ ). *Suppose that (4.4) holds. Let  $\mu_\sigma$  and  $\mu_\xi$  be two compactly supported probability measures on  $[0, \infty)$  such that neither  $\mu_\sigma^{\mathrm{sym}}$  nor  $\mu_\xi^{\mathrm{sym}}$  is a single point mass and at least of one of them is supported at more than two points. Fix some  $L > 0$  and let  $\mathcal{I}$  be any compact interval of the bulk  $\mathcal{B}_{\mu_\sigma^{\mathrm{sym}} \boxplus \mu_\xi^{\mathrm{sym}}}$ . Then there exists a (small) constant  $b_0 > 0$  and  $N_0 \in \mathbb{N}$ , depending only on  $\mu_\sigma$ ,  $\mu_\xi$ ,  $\mathcal{I}$  and the constant  $C$  in (4.4), such that whenever*

$$\sup_{N \geq N_0} (\mathrm{d}_L(\mu_\Sigma, \mu_\sigma) + \mathrm{d}_L(\mu_\Xi, \mu_\xi)) \leq 2b, \quad (4.10)$$

for some  $b \leq b_0$ , then

$$\left| m_H(z) - m_{\mu_\Sigma^{\mathrm{sym}} \boxplus \mu_\Xi^{\mathrm{sym}}}(z) \right| \prec \frac{1}{N\eta(1+\eta)} \quad (4.11)$$

holds uniformly on  $\mathcal{S}_{\mathcal{I}}(0, N^L)$ , for  $N$  sufficiently large depending only on  $\mu_\sigma$ ,  $\mu_\xi$ ,  $\mathcal{I}$ ,  $L$  and the constant  $C$  in (4.4). Moreover, there exists a constant  $\eta_M \geq 1$ , independent of  $N$ , such

that (4.11) holds uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ , for any compact interval  $\mathcal{I} \subset \mathbb{R}$ , for  $N$  sufficiently large depending only on  $\mu_\sigma, \mu_\xi, L$  and the constant  $C$  in (4.4).

Theorem 4.3 is proved in Sections 5-6 and Section 8. In fact in Section 8, we prove Theorem 4.3 for spectral parameters  $z \in \mathbb{C}^+$  with large imaginary parts,  $\eta$ . Here, large  $\eta$  means  $\eta \geq \eta_M$ , for some  $\eta_M \geq 1$  independent of  $N$  to be chosen below. The proof for large  $\eta$  relies on the Gromov–Milman concentration inequality for the full Haar measure in conjunction with identities for expectations of Green functions originating in the global  $U(N)$ -symmetry. These arguments are independent of the main line followed here and are hence postponed to Section 8. The results for large  $\eta$  serve as initial estimates in a bootstrap argument carried out in Sections 5-6 where we prove Theorem 4.3 in the complementary regime where  $\eta < \eta_M$ .

*Proof of Theorem 2.5.* Theorem 2.5 follows from Theorem 4.3 by choosing  $\mathcal{I} = \{0\}$ . The conditions of Theorem 4.3 require that the density of  $\mu_\sigma^{\text{sym}} \boxplus \mu_\xi^{\text{sym}}$  is uniformly bounded from below on the compact interval  $\mathcal{I}$ . For  $E = 0$ , this condition was verified in Theorem 2.2. This yields (2.17) uniformly for  $0 < \eta \leq N^L$ , with  $L > 1$  as in Theorem 4.3 for fixed  $w$  with  $|w| \in [r_- + \tau, r_+ - \tau]$ .

Next, we show that (2.17) can be strengthened to a uniform bound in  $w \in \mathcal{R}_\sigma^\tau := \{w \in \mathbb{C} : |w| \in [r_- + \tau, r_+ - \tau]\}$ . We introduce the lattice

$$\widehat{\mathcal{R}}_\sigma^\tau(L_1) := \mathcal{R}_\sigma^\tau \cap N^{-L_1} \{\mathbb{Z} \times i\mathbb{Z}\},$$

for some sufficiently large positive constant  $L_1$  such that  $L_1 \geq 2L$  (say). Using the definition of stochastic domination in Definition 1.6 and (2.17) for fixed  $w$ , we obtain

$$\max_{w \in \widehat{\mathcal{R}}_\sigma^\tau(L_1)} |m^w(z) - m_{\Sigma, |w|}(z)| \prec \frac{1}{N\eta},$$

uniformly in  $0 < \eta \leq N^L$ . To extend this bound to all of  $\mathcal{R}_\sigma^\tau$ , it suffices to show Lipschitz continuity of these quantities in  $w$ . We need that, for any  $w_1, w_2 \in \mathcal{R}_\sigma^\tau$  with  $|w_1 - w_2| \leq N^{-L_1}$  for sufficiently large  $L_1$ , one has

$$|m^{w_1}(z) - m^{w_2}(z)| \leq \frac{1}{N\eta}, \quad |m_{\Sigma, |w_1|}(z) - m_{\Sigma, |w_2|}(z)| \leq \frac{1}{N\eta}, \quad (4.12)$$

uniformly in  $0 < \eta \leq N^L$ . To show the first deterministic bound in (4.12), we use the bound

$$\begin{aligned} |m^{w_1}(z) - m^{w_2}(z)| &\leq \frac{|w_1 - w_2|}{2N} \text{Tr} |H^{w_1} - z|^{-1} |H^{w_2} - z|^{-1} \\ &\leq \frac{|w_1 - w_2|}{2\eta^2} \leq \frac{1}{2\eta} N^{-L_1+L} \leq \frac{1}{N\eta}, \end{aligned}$$

where  $|A| := \sqrt{A^*A}$ , for any square matrix  $A$ .

To show the second bound in (4.12), we use the stability of the Stieltjes transform of free additive convolution. Here it suffices to use the following bound (*c.f.*, (2.20) in [2] for instance)

$$|m_{\Sigma, |w_1|}(z) - m_{\Sigma, |w_2|}(z)| \leq \frac{C}{\eta} \left(1 + \frac{1}{\eta}\right) d_{\mathbb{L}} \left( \delta_{|w_1|}^{\text{sym}}, \delta_{|w_2|}^{\text{sym}} \right) \leq \frac{C}{\eta} \left(1 + \frac{1}{\eta}\right) |w_1 - w_2|,$$

for all  $z = E + i\eta \in \mathbb{C}^+$ , where  $C$  is a constant uniform in  $z$ . Using the assumptions  $|w_1 - w_2| \leq N^{-L_1}$  and  $L_1 \geq 2L$ , we get (4.12), which in turn establishes the desired uniformity of (2.17) in  $w \in \mathcal{R}_\sigma^\tau$ .

To complete the proof of (2.17), it remains to deal with the large  $\eta$  regime, *i.e.*, when  $\eta \geq N^L$ . For that we use the elementary (deterministic) estimates

$$m_H(i\eta) = -\frac{1}{i\eta} + O\left(\frac{1}{|\eta|^3}\right), \quad m_{\Sigma, |w|}(i\eta) = -\frac{1}{i\eta} + O\left(\frac{1}{|\eta|^3}\right), \quad (4.13)$$

as  $\eta \nearrow \infty$ , where we used a resolvent expansion of  $G$  together with  $\text{tr}H = 0$  and  $\|H\| \leq S_+$  (see (1.3)), and the large  $\eta$  expansion of the Stieltjes transform together with the fact

that  $\mu_{\Sigma,|w|}$  is symmetric and compactly supported. Thus for  $\eta \geq N^L$ , (2.17) follows from (4.13). Uniformity in  $w \in \mathcal{R}_\sigma^r$  is immediate.  $\square$

**4.1. Approximate subordination for block additive models.** In this subsection, we establish the matrix subordination for the Green function of  $H$ . To simplify notation, we introduce the block matrices

$$A := \begin{pmatrix} 0 & \Xi \\ \Xi^* & 0 \end{pmatrix}, \quad B := \begin{pmatrix} 0 & \Sigma \\ \Sigma^* & 0 \end{pmatrix}, \quad \mathcal{U} := \begin{pmatrix} U & 0 \\ 0 & V \end{pmatrix}. \quad (4.14)$$

Then we write (4.2) as

$$H = A + \tilde{B}, \quad \tilde{B} := \mathcal{U}B\mathcal{U}^*. \quad (4.15)$$

As before, we let  $G(z) := (H - z)^{-1}$  be the Green function of  $H$  at spectral parameter  $z \in \mathbb{C}^+$ . A simple consequence of the definition of  $G$  are the identities

$$\tilde{B}G(z) = I_{2N} - (A - z)G(z), \quad G(z)\tilde{B} = I_{2N} - G(z)(A - z). \quad (4.16)$$

Inspired by [33], see also [3, 10, 28], we introduce the approximate subordination functions

$$\omega_A^c(z) := z - \frac{\operatorname{tr} AG}{\operatorname{tr} G}, \quad \omega_B^c(z) := z - \frac{\operatorname{tr} \tilde{B}G}{\operatorname{tr} G}. \quad (4.17)$$

By these definitions and (4.16), we have

$$\omega_A^c(z) + \omega_B^c(z) - z = -\frac{1}{m_H(z)}. \quad (4.18)$$

Recall the measures  $\mu_\Sigma$  and  $\mu_\Xi$  of (4.5) as well as  $\mu_\sigma$  and  $\mu_\xi$  of (4.6). For their symmetrizations we introduce, hinting at (4.14), the shorthands

$$\mu_A \equiv \mu_\Xi^{\operatorname{sym}}, \quad \mu_B \equiv \mu_\Sigma^{\operatorname{sym}}, \quad \mu_\alpha \equiv \mu_\xi^{\operatorname{sym}}, \quad \mu_\beta \equiv \mu_\sigma^{\operatorname{sym}}. \quad (4.19)$$

Note that  $\mu_A$  and  $\mu_B$  are the empirical spectral distributions of  $A$  and  $B$ . We denote by  $\omega_A(z), \omega_B(z), \omega_\alpha(z), \omega_\beta(z)$  the subordination functions defined via (2.5) with the choices  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$  and  $(\mu_\alpha, \mu_\beta)$ , respectively.

The next result shows that the approximate subordination functions  $\omega_A^c$  and  $\omega_B^c$  are indeed good approximations to the subordination functions  $\omega_A$  and  $\omega_B$ . Moreover, it establishes the subordination for the diagonal Green function entries.

**Theorem 4.4.** *Under the conditions and with the notations of Theorem 4.3 the estimates*

$$|\omega_A^c(z) - \omega_A(z)| \prec \frac{1}{N\eta}, \quad |\omega_B^c(z) - \omega_B(z)| \prec \frac{1}{N\eta}, \quad (4.20)$$

hold uniformly on  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$ , for  $N$  sufficiently large depending only on  $\mu_\alpha, \mu_\beta, \mathcal{I}, L$  and the constant  $C$  in (4.4). Moreover, we have

$$\begin{aligned} \left| G_{ii}(z) - \frac{\omega_B(z)}{|\xi_i|^2 - (\omega_B(z))^2} \right| &\prec \frac{1}{\sqrt{N\eta}}, & \left| G_{i\hat{i}(z)} - \frac{\omega_B(z)}{|\xi_i|^2 - (\omega_B(z))^2} \right| &\prec \frac{1}{\sqrt{N\eta}}, \\ \left| G_{i\hat{i}(z)} - \frac{\xi_i}{|\xi_i|^2 - (\omega_B(z))^2} \right| &\prec \frac{1}{\sqrt{N\eta}}, & \left| G_{\hat{i}i}(z) - \frac{\bar{\xi}_i}{|\xi_i|^2 - (\omega_B(z))^2} \right| &\prec \frac{1}{\sqrt{N\eta}}, \end{aligned} \quad (4.21)$$

uniformly in  $i \in \llbracket 1, N \rrbracket$  and in  $z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)$ , where  $\hat{i} := i + N$ , for  $N$  sufficiently large depending only on  $\mu_\alpha, \mu_\beta, \mathcal{I}, L$  and the constant  $C$  in (4.4).

*Remark 4.5.* Some crucial properties of the subordination functions  $\omega_A$  and  $\omega_B$  are collected in Lemma A.2. Here, we mention that under the assumptions of Theorem 4.4, for  $N$  sufficiently large, the imaginary parts of the subordination functions,  $\operatorname{Im} \omega_A(z)$  and  $\operatorname{Im} \omega_B(z)$  are both bounded from below on  $z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)$ . This follows from Lemma A.2 and the assumption that  $\mathcal{I}$  is a compact interval in the bulk of  $\mu_\alpha \boxplus \mu_\beta$ . It then follows from (4.21) that  $|G_{ii}(z)| \prec 1$  and  $|G_{i\hat{i}(z)}| \prec 1$  uniformly on  $\mathcal{S}_{\mathcal{I}}(N^{-1+\gamma}, \eta_M)$ , for any  $\gamma > 0$ , and all  $i \in \llbracket 1, N \rrbracket$ . A direct

consequence of this result is that the eigenvectors associated with eigenvalues in the bulk are fully delocalized. More precisely, letting  $(\mathbf{u}_k)$  denote the  $\ell^2$ -normalized eigenvectors associated with the eigenvalues  $(\lambda_k)$ ,  $k \in \llbracket 1, 2N \rrbracket$ , we have

$$\max_{k: \lambda_k \in \mathcal{I}} \|\mathbf{u}_k\|_\infty \prec \frac{1}{\sqrt{N}}, \quad (4.22)$$

for any compact interval  $\mathcal{I}$  in the bulk of  $\mu_\alpha \boxplus \mu_\beta$ . For a proof of (4.22) from Theorem 4.4, we refer to the proof of Theorem 2.6 in [3].

**4.2. Outline of the strategy of proof.** The proof of the local law of Theorem 4.3 is carried out in three steps. In Step 1, we consider the large  $\eta$  regime, *i.e.*, we establish (4.11) on  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ , for some sufficiently large, but  $N$ -independent,  $\eta_M$ . In Step 2, we establish a weak local law for  $m^w$  in the small  $\eta$  regime, *i.e.*, we establish (4.11) with a weaker error bound on  $\mathcal{S}_{\mathcal{I}}(N^{-1+\gamma}, \eta_M)$ , for some small  $\gamma > 0$ ; see Theorem 5.1 below for the statement of the weak law. The extension to  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$  will follow directly from monotonicity of the Green function. This second step is based on a bootstrapping argument to reduce the spectral parameter  $\text{Im } z$ . Step 1 will provide the initial estimate to get the bootstrapping started. In Step 3, we use a fluctuation averaging argument together with the weak local law established in the second step to get (4.11) in its strong form.

Step 1 is carried out in Section 8. It builds on the celebrated Gromov–Milman concentration inequality whose application to random matrix theory is fairly standard [1]. For additive models of the form  $X + UYU^*$ , with deterministic  $X, Y \in M_N(\mathbb{C})$  and  $U$  Haar distributed on  $U(N)$  or on  $O(N)$  it was used in [33, 28, 2], and for the model block-additive model considered in this section in [10].

Step 2 is carried out in Section 5, where we prove Theorem 5.1. This proof has three major ingredients. First, we use a partial randomness decomposition of the Haar measure (see (5.2)) that enables us to take partial expectations of functions of the diagonal Green function entries  $G_{ii}$ ,  $G_{\tilde{i}\tilde{i}}$ . Exploiting concentration only for this partial randomness surpasses the more general but less flexible Gromov–Milman technique used in Step 1. Second, to compute the partial expectations of  $G_{ii}$ , we establish a system of self-consistent equations involving only two auxiliary quantities  $(S_{ij})$  and  $(T_{ij})$ ; see (5.16). In our previous work [3], we used a similar approach to derive the local law for  $X + UYU^*$ . For the model considered in this paper, we face with a new phenomenon causing several substantial difficulties. The main point is that for block additive models, we have less randomness originating in the Haar measure on  $U(N) \times U(N)$  than for the additive models with Haar measure on  $U(2N)$ . As a consequence, we have to control more quantities in the two blocks separately. Even more importantly, the coupling between the two blocks is provided solely by the diagonal matrix  $\Sigma$  without any randomness; see (4.2). Our proof shows that the randomness in the diagonal blocks and the deterministic off-diagonal blocks effectively make up for the lacking off-diagonal randomness.

To derive the aforementioned system of equations for  $(S_{ij})$ ,  $(T_{ij})$  and  $(G_{ii})$ , we use the partial decomposition of Haar measure in combination with recursive moment estimates; see *e.g.*, Lemma 5.3 for such a statement. Recursive moment estimates were used first in [29] to derive local laws for sparse Wigner matrices. They allow us to pass on cumbersome partial concentration estimates used in Section 5 of [3], and provide a conceptually clear approach to the weak local law for both models. Third, to connect the diagonal Green function entries with the subordination functions from Theorem 2.1, we rely on the optimal stability result for the subordination equations obtained in [2].

Step 3 is carried out in Section 6. In this section, we exploit the so-called fluctuation averaging mechanism to improve the estimates of Step 2. While the fluctuation averaging mechanism is, thanks to the independence of the matrix entries, well understood for Wigner type matrices (see *e.g.*, [19, 20]), dependencies among the entries of the Haar matrices mask this mechanism and its current understanding for matrix ensembles involving Haar matrices



is still rather poor. We gave a first result in [4] for additive models. In the present paper, we approach the fluctuation mechanics for block-additive models by first deriving a set of so-called ‘‘Ward identities’’ which will enable us to finish the proof of Theorem 4.3. Ward identities are relations among tracial quantities involving the Green function and the matrices  $A$  and  $B$ . In expectation, these relations can be derived using the invariance of Haar measure (see *e.g.*, (8.11) for a first example), yet we will require optimal estimates that hold with high probability; see *e.g.*, (5.21) and (6.3). These estimates are obtained using recursive moment estimates for carefully chosen quantities; see (5.19). Since we have less randomness coming from  $U(N) \times U(N)$  in the setup of block-additive models, more quantities need to be simultaneously controlled than in the additive models, resulting in a more sophisticated analysis.

**4.3. Notations.** We introduce some more notation used in the proof of Theorem 4.3.

*Notation for matrices:* In our analysis we also use the matrices

$$\mathcal{H} = B + \mathcal{U}^* A \mathcal{U} =: B + \tilde{A}, \quad \mathcal{G}(z) = (\mathcal{H} - z)^{-1}, \quad (4.23)$$

which are the analogues of  $H$  in (4.15) and of its Green function  $G(z)$ , obtained by switching the rôles of  $A$  and  $B$ , and also the rôles of  $\mathcal{U}$  and  $\mathcal{U}^*$ . Note that by cyclicity  $\text{Tr} G(z) = \text{Tr} \mathcal{G}(z)$ .

*Vector space notation:* For any index  $i \in \llbracket 1, N \rrbracket$ , we let  $\hat{i} \equiv i + N$ . We make the convention hereafter that the index  $i$  always runs from 1 to  $N$ , unless said otherwise. Thus the index  $\hat{i}$  runs from  $N + 1$  to  $2N$ . We denote by  $\sum_i^{(k)}$  the sum over  $i \in \llbracket 1, N \rrbracket \setminus \{k\}$ . We denote by  $\{e_i\}$  the canonical basis of  $\mathbb{C}^N$  while we denote by  $\{\hat{e}_i\}$  the canonical basis of  $\mathbb{C}^{2N}$ . We let  $\mathbf{0}$  denote the zero vector in either space. We use bold font for vectors and denote the components as  $\mathbf{v} = (v_i)$ .

The identity matrix in  $M_N(\mathbb{C})$ , respectively  $M_{2N}(\mathbb{C})$ , is denoted by

$$I \equiv I_N, \quad \hat{I} \equiv I_{2N}, \quad (4.24)$$

and we let

$$\hat{I}_1 := I \oplus 0, \quad \hat{I}_2 := 0 \oplus I \quad (4.25)$$

denote the block identities in  $M_N(\mathbb{C}) \oplus M_N(\mathbb{C})$ , where  $0$  represents the  $N \times N$  zero matrix.

For any matrix  $D \in M_n(\mathbb{C})$ ,  $n \geq 1$ , we let

$$\text{tr } D := \frac{1}{n} \text{Tr } D$$

denote the normalized trace of  $D$ . For  $D \in M_{2N}(\mathbb{C})$  we introduce the normalized partial traces

$$\tau_1(D) := \frac{1}{N} \sum_{i=1}^N D_{ii}, \quad \tau_2(D) := \frac{1}{N} \sum_{i=1}^N D_{\hat{i}\hat{i}}. \quad (4.26)$$

Using the block structure of  $H$ , it is easy to check that the Green function  $G(z)$  satisfies

$$\tau_1(G(z)) = \tau_2(G(z)), \quad z \in \mathbb{C}^+. \quad (4.27)$$

*$\Phi$ -system:* For our purposes it is convenient to recast (2.5) in a compact form: For generic probability measures  $\mu_1, \mu_2$  on  $\mathbb{R}$ , let the function  $\Phi_{\mu_1, \mu_2} : (\mathbb{C}^+)^3 \rightarrow \mathbb{C}^2$  be given by

$$\Phi_{\mu_1, \mu_2}(\omega_1, \omega_2, z) := \begin{pmatrix} F_{\mu_1}(\omega_2) - \omega_1 - \omega_2 + z \\ F_{\mu_2}(\omega_1) - \omega_1 - \omega_2 + z \end{pmatrix}. \quad (4.28)$$

Considering  $\mu_1, \mu_2$  as fixed, the equation

$$\Phi_{\mu_1, \mu_2}(\omega_1, \omega_2, z) = \mathbf{0}, \quad (4.29)$$

is equivalent to (2.5) and, by Theorem 2.1, there are unique analytic functions  $\omega_1, \omega_2 : \mathbb{C}^+ \rightarrow \mathbb{C}^+$ ,  $z \mapsto \omega_1(z), \omega_2(z)$  satisfying (2.4) that solve (4.29) in terms of  $z$ .

*Control parameters:* For  $z \in \mathbb{C}^+$ , we will use the following deterministic control parameter

$$\Psi \equiv \Psi(z) := \frac{1}{\sqrt{N\eta(1+\eta)}}, \quad \eta = \text{Im } z. \quad (4.30)$$

We further introduce, for  $z \in \mathbb{C}^+$  and  $i \in \llbracket 1, N \rrbracket$ , the random control parameters

$$\begin{aligned} \Lambda_{d;ii}(z) &:= \left| G_{ii} - \frac{\omega_B(z)}{|\xi_i|^2 - (\omega_B(z))^2} \right|, & \Lambda_{d;\hat{ii}}(z) &:= \left| G_{\hat{ii}} - \frac{\omega_B(z)}{|\xi_i|^2 - (\omega_B(z))^2} \right|, \\ \Lambda_{d;\hat{i}\hat{i}}(z) &:= \left| G_{\hat{i}\hat{i}} - \frac{\xi_i}{|\xi_i|^2 - (\omega_B(z))^2} \right|, & \Lambda_{d;i\hat{i}}(z) &:= \left| G_{i\hat{i}} - \frac{\bar{\xi}_i}{|\xi_i|^2 - (\omega_B(z))^2} \right|, \end{aligned}$$

$$\Lambda_d(z) := \max_{i \in \llbracket 1, N \rrbracket} \max_{k, l = i \text{ or } \hat{i}} \Lambda_{d;kl}(z). \quad (4.31)$$

We also define  $\Lambda_d^c(z)$  analogously by replacing  $\omega_B$  by  $\omega_B^c$  (c.f., (4.17)) in the definition of  $\Lambda_d(z)$ . We will often omit the variable  $z$  from the above notations when there is no confusion.

For notational simplicity, we do not follow the threshold  $N$  for which the estimates apply. Following the dependence of this threshold on the other parameters along the proofs one may easily verify the dependences stated in Theorem 4.3 and Theorem 4.4.

## 5. GREEN FUNCTION SUBORDINATION FOR SMALL $\eta$

Let  $\eta_M > 0$  be some sufficiently large constant, and for any given (small)  $\gamma > 0$ , we set

$$\eta_m \equiv \eta_m(\gamma) := N^{-1+\gamma}. \quad (5.1)$$

In this section, we prove a Green function subordination property in the regime  $\eta_m \leq \eta \leq \eta_M$ . The formal statement is given in Theorem 5.1 below. For definiteness, we work with the unitary setup in this section. The necessary modifications for the orthogonal case are stated in Appendix C. We start with the partial randomness decomposition of the Haar measure on  $U(N) \times U(N)$  announced in Subsection 4.2.

**5.1. Partial randomness decomposition of the Haar measure.** Let  $\mathbf{u}_i = (u_{i1}, \dots, u_{iN})'$  and  $\mathbf{v}_i = (v_{i1}, \dots, v_{iN})'$  be the  $i$ th columns of  $U$  and  $V$ , respectively. Let  $\theta_i^u$  and  $\theta_i^v$  be the arguments of  $u_{ii}$  and  $v_{ii}$ , respectively, and let  $\phi_i^a = e^{i\theta_i^a}$  for  $a = u, v$ . Our approach relies on the partial randomness decomposition of the Haar measure from [17, 32]:

$$U = -\phi_i^u R_i^u U^{(i)}, \quad V = -\phi_i^v R_i^v V^{(i)}. \quad (5.2)$$

Here  $U^{(i)}$  and  $V^{(i)}$  are unitary matrices with  $(i, i)$ -th entry equal 1, and their  $(i, i)$ -minors are independent, Haar distributed on  $\mathcal{U}(N-1)$ . In particular,  $U^{(i)} \mathbf{e}_i = V^{(i)} \mathbf{e}_i = \mathbf{e}_i$  and  $\mathbf{e}_i^* U^{(i)} = \mathbf{e}_i^* V^{(i)} = \mathbf{e}_i^*$ , where  $\mathbf{e}_i$  is the  $i$ -th coordinate vector. In addition,  $U^{(i)}$  is independent of  $\mathbf{u}_i$ , and  $V^{(i)}$  is independent of  $\mathbf{v}_i$ . Here  $R_i^u$  and  $R_i^v$  are reflections, defined as

$$R_i^a := I - \mathbf{r}_i^a (\mathbf{r}_i^a)^*, \quad a = u, v, \quad (5.3)$$

where

$$\mathbf{r}_i^u := \sqrt{2} \frac{\mathbf{e}_i + \bar{\phi}_i^u \mathbf{u}_i}{\|\mathbf{e}_i + \bar{\phi}_i^u \mathbf{u}_i\|_2}, \quad \mathbf{r}_i^v := \sqrt{2} \frac{\mathbf{e}_i + \bar{\phi}_i^v \mathbf{v}_i}{\|\mathbf{e}_i + \bar{\phi}_i^v \mathbf{v}_i\|_2}. \quad (5.4)$$

Note that  $R_i^u$  is independent of  $U^{(i)}$  and  $R_i^v$  is independent of  $V^{(i)}$ .

Set the  $(2N) \times (2N)$  matrices

$$\Phi_i := (\phi_i^u I) \oplus (\phi_i^v I), \quad \mathcal{R}_i := R_i^u \oplus R_i^v, \quad \mathcal{U}_i := U^{(i)} \oplus V^{(i)}. \quad (5.5)$$

With the above notations and the decompositions in (5.2), we have

$$\mathcal{U} = -\mathcal{R}_i \mathcal{U}_i \Phi_i. \quad (5.6)$$

Hence, for each  $i \in \llbracket 1, N \rrbracket$ , we can write

$$H = A + \widetilde{B} = A + \mathcal{R}_i \mathcal{U}_i \Phi_i B \Phi_i^* \mathcal{U}_i^* \mathcal{R}_i := A + \mathcal{R}_i \widetilde{B}^{(i)} \mathcal{R}_i, \quad (5.7)$$

where we introduced the notation

$$\widetilde{B}^{(i)} := \mathcal{U}_i \Phi_i B \Phi_i^* \mathcal{U}_i^*. \quad (5.8)$$

We further define the matrices

$$H^{(i)} := A + \widetilde{B}^{(i)}, \quad G^{(i)} := (H^{(i)} - z)^{-1}. \quad (5.9)$$

Since  $\mathbf{u}_i$  and  $\mathbf{v}_i$  are independent, uniformly distributed complex unit vectors, there exist independent normal vectors,  $\widetilde{\mathbf{g}}_i^u, \widetilde{\mathbf{g}}_i^v \sim \mathcal{N}_{\mathbb{C}}(0, \frac{1}{N} I_N)$  such that

$$\mathbf{u}_i = \frac{\widetilde{\mathbf{g}}_i^u}{\|\widetilde{\mathbf{g}}_i^u\|_2}, \quad \mathbf{v}_i = \frac{\widetilde{\mathbf{g}}_i^v}{\|\widetilde{\mathbf{g}}_i^v\|_2}.$$

We further define

$$\mathbf{g}_i^u := \overline{\phi}_i^u \widetilde{\mathbf{g}}_i^u, \quad \mathbf{h}_i^u := \frac{\mathbf{g}_i^u}{\|\mathbf{g}_i^u\|_2} = \overline{\phi}_i^u \mathbf{u}_i, \quad \ell_i^u := \frac{\sqrt{2}}{\|\mathbf{e}_i + \mathbf{h}_i^u\|_2}, \quad (5.10)$$

and define  $\mathbf{g}_i^v, \mathbf{h}_i^v$  and  $\ell_i^v$  analogously by replacing  $\mathbf{u}_i$  by  $\mathbf{v}_i$ . Note that for  $a = u$  or  $v$ ,  $g_{ik}^a$ 's for  $k \neq i$  are  $\mathcal{N}_{\mathbb{C}}(0, \frac{1}{N})$  variables and  $g_{ii}^a$  is  $\chi$ -distributed with  $\mathbb{E}[(g_{ii}^a)^2] = \frac{1}{N}$ . In addition, the components of  $\mathbf{g}_i^a$  are independent, and they are all independent of  $\phi_i^a$ . Hence,  $\mathbf{g}_i^a$  and  $\mathbf{h}_i^a$  are independent of  $\widetilde{B}^{(i)}$  (c.f., (5.8)), for  $a = u, v$ . With these notations, we can write

$$\mathbf{r}_i^a = \ell_i^a (\mathbf{e}_i + \mathbf{h}_i^a), \quad a = u, v, \quad (5.11)$$

where  $\mathbf{r}_i^a$  is defined in (5.4). Using Lemma A.1, it is elementary to check that, for  $a = u, v$ ,

$$\|\mathbf{g}_i^a\|_2 = 1 + \frac{1}{2} (\|\mathbf{g}_i^a\|_2^2 - 1) + O_{\prec}(\frac{1}{N}), \quad (\ell_i^a)^2 = \frac{1}{1 + \mathbf{e}_i^* \mathbf{h}_i^a} = 1 - g_{ii}^a + O_{\prec}(\frac{1}{N}), \quad (5.12)$$

where in the first estimate we used the fact  $|\|\mathbf{g}_i^a\|_2^2 - 1| \prec \frac{1}{\sqrt{N}}$ . In addition, by definition,  $R_i^a$  is a reflection sending  $\mathbf{e}_i$  to  $-\mathbf{h}_i^a$ , i.e.,

$$R_i^a \mathbf{e}_i = -\mathbf{h}_i^a, \quad R_i^a \mathbf{h}_i^a = -\mathbf{e}_i, \quad a = u, v. \quad (5.13)$$

We also denote by  $\hat{\mathbf{g}}_i^a$  the vector obtained from  $\mathbf{g}_i^a$  by replacing  $g_{ii}^a$  by 0, i.e.,

$$\hat{\mathbf{g}}_i^a := \mathbf{g}_i^a - g_{ii}^a \mathbf{e}_i, \quad a = u, v.$$

Correspondingly, we set

$$\hat{\mathbf{h}}_i^a := \frac{\hat{\mathbf{g}}_i^a}{\|\hat{\mathbf{g}}_i^a\|_2}, \quad a = u, v. \quad (5.14)$$

Recall the notation  $\mathbf{0}$  for the  $N \times 1$  null vector. Finally, for brevity, we set

$$\mathbf{k}_i^u := \begin{pmatrix} \hat{\mathbf{h}}_i^u \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{k}_i^v := \begin{pmatrix} \mathbf{0} \\ \hat{\mathbf{h}}_i^v \end{pmatrix}, \quad \hat{\mathbf{k}}_i^u := \begin{pmatrix} \hat{\mathbf{h}}_i^u \\ \mathbf{0} \end{pmatrix}, \quad \hat{\mathbf{k}}_i^v := \begin{pmatrix} \mathbf{0} \\ \hat{\mathbf{h}}_i^v \end{pmatrix}. \quad (5.15)$$

We move on to the formal statement of the Green function subordination.

**5.2. Green function subordination.** Recall the notation  $\{\hat{e}_i\}$  for the standard basis of  $\mathbb{C}^{2N}$ , and also the notation  $\hat{i} \equiv i + N$  for any  $i \in \llbracket 1, N \rrbracket$ . We introduce the following quantities for  $j = i, \hat{i}$ ,  $i \in \llbracket 1, N \rrbracket$ ,

$$\begin{aligned} S_{ij} &:= (\mathbf{k}_i^u)^* \tilde{B}^{(i)} G \hat{e}_j, & T_{ij} &:= (\mathbf{k}_i^u)^* G \hat{e}_j, \\ S_{\hat{i}j} &:= (\mathbf{k}_i^v)^* \tilde{B}^{(i)} G \hat{e}_j, & T_{\hat{i}j} &:= (\mathbf{k}_i^v)^* G \hat{e}_j, \end{aligned} \quad (5.16)$$

and

$$\mathring{S}_{ii} := (\mathring{\mathbf{k}}_i^u)^* \tilde{B}^{(i)} G \hat{e}_i = S_{ii} - \tilde{\sigma}_i h_{ii}^u G_{\hat{i}i}, \quad \mathring{T}_{ii} := (\mathring{\mathbf{k}}_i^u)^* G \hat{e}_i = T_{ii} - h_{ii}^u G_{ii}, \quad (5.17)$$

where  $\tilde{\sigma}_i = \phi_i^u \bar{\phi}_i^v \sigma_i$ , and  $\sigma_i$  is the  $i$ th diagonal entry of  $\Sigma$ , *c.f.*, (4.3). Here in (5.17) we used

$$\hat{e}_i^* \tilde{B}^{(i)} = \tilde{\sigma}_i \hat{e}_i^*, \quad \tilde{B}^{(i)} \hat{e}_i = \tilde{\sigma}_i \hat{e}_i, \quad i \in \llbracket 1, N \rrbracket, \quad (5.18)$$

which is checked from the definitions of  $\tilde{B}^{(i)}$  in (5.8),  $\mathcal{U}_i$  and  $\Phi_i$  in (5.5), and also  $B$  in (4.14).

Recall from (4.26) the notations for normalized partial traces  $\tau_1$  and  $\tau_2$  on  $M_{2N}(\mathbb{C})$ . Moreover, recall from (4.31) the definition of the control parameters  $\Lambda_{d;ii}(z)$ ,  $\Lambda_{d;\hat{i}\hat{i}}(z)$ ,  $\Lambda_{d;i\hat{i}}(z)$ ,  $\Lambda_{d;\hat{i}i}(z)$  and  $\Lambda_d(z)$ . We further introduce  $\Lambda_d^c(z)$  analogously by replacing  $\omega_B$  by  $\omega_B^c$  (*c.f.*, (4.17)) in the definition of  $\Lambda_d(z)$ . We will often omit the variable  $z$  from these notations.

In this section we will show that  $\Lambda_d(z)$ ,  $\Lambda_d^c(z)$  and  $\Lambda_T$  are of order  $\Psi$  with high probability; *i.e.*, matrix elements of the Green function can be expressed in terms of the subordination functions, up to a small random fluctuations of order  $\Psi$ . We will refer to these results as *Green function subordination*. The main tool is a high moment calculation and Gaussian integration by parts. However, we cannot directly estimate the high moments of  $T_{kl}$  and the formulas  $|G_{ij} - [\dots]|$  defining  $\Lambda_{d;ij}(z)$ . Instead, we introduce the following auxiliary quantities. For each  $i \in \llbracket 1, N \rrbracket$  and  $j = i$  or  $\hat{i}$ , let

$$\begin{aligned} \mathcal{P}_{ij} &\equiv \mathcal{P}_{ij}(z) := (\tilde{B}G)_{ij} \tau_1(G) - G_{ij} \tau_1(\tilde{B}G) + (G_{ij} + T_{ij}) \Upsilon_1, \\ \mathcal{P}_{\hat{i}j} &\equiv \mathcal{P}_{\hat{i}j}(z) := (\tilde{B}G)_{\hat{i}j} \tau_2(G) - G_{\hat{i}j} \tau_2(\tilde{B}G) + (G_{\hat{i}j} + T_{\hat{i}j}) \Upsilon_2, \\ \mathcal{K}_{ij} &\equiv \mathcal{K}_{ij}(z) := T_{ij} + \tau_1(G) (\tilde{\sigma}_i T_{\hat{i}j} + (\tilde{B}G)_{ij}) - \tau_1(G \tilde{B}) (G_{ij} + T_{ij}), \\ \mathcal{K}_{\hat{i}j} &\equiv \mathcal{K}_{\hat{i}j}(z) := T_{\hat{i}j} + \tau_2(G) (\tilde{\sigma}_i^* T_{ij} + (\tilde{B}G)_{\hat{i}j}) - \tau_2(G \tilde{B}) (G_{\hat{i}j} + T_{\hat{i}j}), \end{aligned} \quad (5.19)$$

where, with  $a = 1, 2$ ,

$$\Upsilon_a \equiv \Upsilon_a(z) := \tau_a(\tilde{B}G) + \tau_a(G) \tau_a(\tilde{B}G \tilde{B}) - \tau_a(G \tilde{B}) \tau_a(\tilde{B}G). \quad (5.20)$$

Using the invariance of the Haar measure, the following Ward identities

$$\mathbb{E} \Upsilon_a = 0, \quad a = 1, 2, \quad (5.21)$$

can be checked. However, we will also need to know that  $\Upsilon_a$  are small with high probability and not only in expectation in the following; see *e.g.*, (5.29) in Theorem 5.2 below.

We will compute their high moments of these auxiliary quantities  $\mathcal{P}$  and  $\mathcal{K}$  and from them we will conclude the estimates on the  $\Lambda$ 's. The careful choice of these auxiliary quantities  $\mathcal{P}$  and  $\mathcal{K}$  is essential for the proof. They have a built-in cancellation mechanism that makes the high moment calculation tractable, see (5.53)-(5.55) later.

Moreover, we recall the following matrices introduced in (4.23)

$$\mathcal{H} = B + \mathcal{U}^* A \mathcal{U} =: B + \tilde{A}, \quad \mathcal{G}(z) = (\mathcal{H} - z)^{-1}, \quad z \in \mathbb{C}^+,$$

which are the analogue of  $H$  in (4.15) and its Green function  $G(z)$ , obtained via swapping the rôles of  $A$  and  $B$ , and also the rôles of  $\mathcal{U}$  and  $\mathcal{U}^*$ . Note that the structure of  $\mathcal{H}$  is exactly the same as  $H$ , so we can define the  $\mathcal{H}$ -counterparts of all quantities we have introduced so far for  $H$ . We will not repeat the heavy notations of the partial randomness decomposition for  $\mathcal{H}$  as well, since we will not need all these details. We will only need to know that, accordingly, we can define  $\mathcal{G}_{ij}$ ,  $\mathcal{S}_{ij}$  and  $\mathcal{T}_{ij}$  by applying the same switching in the definitions of  $G_{ij}$ ,  $S_{ij}$  and  $T_{ij}$ .

Also note the following alternative definition of  $\omega_A^c$  and  $\omega_B^c$  in (4.17):

$$\omega_A^c(z) := z - \frac{\operatorname{tr} \tilde{A}\mathcal{G}}{\operatorname{tr} \mathcal{G}}, \quad \omega_B^c(z) := z - \frac{\operatorname{tr} B\mathcal{G}}{\operatorname{tr} \mathcal{G}}, \quad (5.22)$$

and the trivial fact  $\operatorname{tr} G = \operatorname{tr} \mathcal{G}$ .

In addition, replacing  $\xi_i, \omega_B, G_{ij}$  by  $\sigma_i, \omega_A, \mathcal{G}_{ij}$  respectively in (4.31), we define  $\tilde{\Lambda}_{d;ij}(z)$  and  $\tilde{\Lambda}_d(z)$  as the analogues of  $\Lambda_{d;ij}(z)$  and  $\Lambda_d(z)$ . For example

$$\tilde{\Lambda}_{d;ii}(z) = \left| \mathcal{G}_{ii} - \frac{\omega_A(z)}{|\sigma_i|^2 - (\omega_A(z))^2} \right|, \quad (5.23)$$

and

$$\tilde{\Lambda}_d(z) := \max_{i \in \llbracket 1, N \rrbracket} \max_{k, l = i \text{ or } \hat{i}} \tilde{\Lambda}_{d;kl}(z). \quad (5.24)$$

Similarly, we can also define  $\tilde{\Lambda}_d^c(z)$  and  $\tilde{\Lambda}_T(z)$  as the analogue of  $\Lambda_d^c(z)$  and  $\Lambda_T(z)$ , respectively. The analysis of the operator  $\mathcal{H}$  is very similar to that of  $H$ , but at some point it will be useful to work with them in tandem, so we will need to control both.

Our main aim in this section is to prove the following Green function subordination property. Recall the definition of the control parameter  $\Psi(z)$  from (4.30).

**Theorem 5.1.** *Suppose that the assumptions in Theorem 4.3 hold. Then*

$$\Lambda_d(z) \prec \Psi(z), \quad \tilde{\Lambda}_d(z) \prec \Psi(z), \quad \Lambda_T(z) \prec \Psi(z), \quad \tilde{\Lambda}_T(z) \prec \Psi(z) \quad (5.25)$$

*uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ , for any (large) constant  $\eta_M > 0$  and (small) constant  $\gamma > 0$ , in the definition of  $\eta_m$  (c.f., (5.1)). Moreover, the estimates*

$$\begin{aligned} |\omega_A^c(z) - \omega_A(z)| &\prec \Psi(z), & |\omega_B^c(z) - \omega_B(z)| &\prec \Psi(z), \\ |m_H(z) - m_{\mu_A \boxplus \mu_B}(z)| &\prec \Psi(z) \end{aligned} \quad (5.26)$$

*also hold uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ .*

The estimates on the tracial quantities and the subordination functions in (5.26) are weaker than the final result in Theorem 4.3 and Theorem 4.4. Later in Section 6, we will improve them. The estimates in (5.25) are, however, (believed to be) optimal.

In what follows, we will mainly work with  $\Lambda_d(z)$ . The discussion on  $\tilde{\Lambda}_d(z)$  is the same. First, we show the analogous estimate for  $\Lambda_d^c$  by assuming an a priori bound on  $\Lambda_d$  and  $\Lambda_T$ , for a fixed  $z \in \mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . This is the content of Theorem 5.2 below. A continuity argument in Subsection 5.5 then allows us to conclude Theorem 5.1 from Theorem 5.2.

**Theorem 5.2.** *Suppose that the assumptions in Theorem 4.3 hold. Let  $\eta_M > 0$  be a (large) constant and  $\gamma > 0$  be a (small) constant in (5.1). Fix a  $z \in \mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . Assume that*

$$\Lambda_d(z) \prec N^{-\frac{\gamma}{4}}, \quad \tilde{\Lambda}_d(z) \prec N^{-\frac{\gamma}{4}}, \quad \Lambda_T(z) \prec 1, \quad \tilde{\Lambda}_T(z) \prec 1. \quad (5.27)$$

*Then we have*

$$\begin{aligned} |\mathcal{P}_{ij}(z)| &\prec \Psi(z), & |\mathcal{P}_{\hat{i}j}(z)| &\prec \Psi(z), \\ |\mathcal{K}_{ij}(z)| &\prec \Psi(z), & |\mathcal{K}_{\hat{i}j}(z)| &\prec \Psi(z), \end{aligned} \quad (5.28)$$

*for all  $i \in \llbracket 1, N \rrbracket$  and  $j = i$  or  $\hat{i}$ . In addition, under (5.27) we also have*

$$|\Upsilon_1(z)| \prec \Psi(z), \quad |\Upsilon_2(z)| \prec \Psi(z), \quad (5.29)$$

*and*

$$\Lambda_d^c(z) \prec \Psi(z), \quad \Lambda_T(z) \prec \Psi(z). \quad (5.30)$$

*The same statements hold if we switch the rôles of  $A$  and  $B$ , and also the rôles of  $U$  and  $U^*$ , in all the conclusions from (5.28) to (5.30).*

Note that, since  $\eta_m \leq \eta \leq \eta_M$ , we have  $\Psi(z) \sim \frac{1}{\sqrt{N\eta}}$ .

The proof of Theorem 5.2 proceeds in two steps. In the first step, we establish in Subsection 5.3 recursive moment estimates for the quantities  $\mathcal{P}_{ii}$  and  $\mathcal{K}_{ii}$ . In the second step, carried out in Subsection 5.4, we use a local stability analysis to conclude Theorem 5.2 from the estimates established in Subsection 5.3.

**5.3. Recursive moment estimates for  $\mathcal{P}_{ii}$  and  $\mathcal{K}_{ii}$ .** In the proof of Theorem 5.2, assumption (5.27) is used to conclude that various  $G_{kl}$  and  $T_{kl}$  with  $k, l = i$  or  $\hat{i}$  are finite. More specifically, with the aid of assumption (5.27) and with the upper bound of  $|\omega_B|$  and the lower bound on  $\text{Im } \omega_B$  in (A.4) that together imply that  $\omega_B^2$  is away from the positive real axis so the denominators in the definition of  $\Lambda_{d,ij}$  do not vanish, we have

$$\max_{i \in [1, N]} \max_{k, l = i \text{ or } \hat{i}} |G_{kl}| < 1, \quad \max_{i \in [1, N]} \max_{k, l = i \text{ or } \hat{i}} |T_{kl}| < 1. \quad (5.31)$$

In addition, using the identities in (4.16), we can further get the bound

$$\max_{i \in [1, N]} \max_{k, l = i \text{ or } \hat{i}} |(XGY)_{kl}| < 1, \quad X, Y = \hat{I}, \text{ or } \tilde{B}. \quad (5.32)$$

Observe that

$$\frac{1}{N} \sum_i \frac{\omega_B(z)}{|\xi_i|^2 - (\omega_B(z))^2} = m_{\mu_A}(\omega_B(z)) = m_{\mu_A \boxplus \mu_B}(z), \quad (5.33)$$

where the first step follows from the definition of  $\mu_A$  in (4.19), and the second step follows from (2.5) with the choice  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$ . Then, (5.33) together with the first estimate in (5.27), (4.16), and the upper bound of  $|\omega_B|$  and the lower bound of  $\text{Im } \omega_B$  in (A.4) leads to the following estimates for tracial quantities

$$\begin{aligned} \tau_a(G) &= m_{\mu_A \boxplus \mu_B} + O_{<}(N^{-\frac{\gamma}{4}}), & a = 1, 2. \\ \tau_a(\tilde{B}G) &= (z - \omega_B)m_{\mu_A \boxplus \mu_B} + O_{<}(N^{-\frac{\gamma}{4}}), \\ \tau_a(G\tilde{B}) &= (z - \omega_B)m_{\mu_A \boxplus \mu_B} + O_{<}(N^{-\frac{\gamma}{4}}), \\ \tau_a(\tilde{B}G\tilde{B}) &= (\omega_B - z)(1 + (\omega_B - z)m_{\mu_A \boxplus \mu_B}) + O_{<}(N^{-\frac{\gamma}{4}}). \end{aligned} \quad (5.34)$$

Then, using the upper bound on  $|\omega_B|$  and the lower bound on  $\text{Im } \omega_B$  in (A.4), and the second identity in (5.33), we see that all these tracial quantities are stochastically dominated by 1, under assumption (5.27). Recalling  $\Upsilon_a$  from (5.20), we thus have under assumption (5.27) that

$$|\Upsilon_a(z)| < 1. \quad (5.35)$$

For (5.28), we only handle the estimate of  $\mathcal{P}_{ii}$  and  $\mathcal{K}_{ii}$  in detail. The others are similar. It suffices to show the high order moment estimate: for any fixed integer  $p \geq 1$ , we have

$$\mathbb{E}[|\mathcal{P}_{ii}|^{2p}] < \Psi^{2p}, \quad \mathbb{E}[|\mathcal{K}_{ii}|^{2p}] < \Psi^{2p}. \quad (5.36)$$

Let us introduce the notation

$$\mathbf{m}_i(k, l) := \mathcal{P}_{ii}^k \overline{\mathcal{P}_{ii}^l}, \quad \mathbf{n}_i(k, l) := \mathcal{K}_{ii}^k \overline{\mathcal{K}_{ii}^l}. \quad (5.37)$$

We will use the following notational conventions in the statement of the recursive moment estimates. The notation  $O_{<}(\Psi^k)$  for any given positive integer  $k$ , represents a generic (possibly)  $z$ -dependent random variable  $X \equiv X(z)$  that satisfies

$$X < \Psi^k, \quad \mathbb{E}[|X|^q] < \Psi^{qk}, \quad (5.38)$$

for any given positive integer  $q$ . In the sequel, we only check the first bound in (5.38) for various  $X$ 's, then the second bound is valid as well. Indeed, since the  $X$ 's we will encounter below are analogous to those in [4], we refer to the paragraph below (6.2) of [4] for a general reasoning why the second bound in (5.38) follows from the first one. Additionally, sometimes  $X$  will be of

the form  $1/|g|$  where  $g$  is an  $N$ -dimensional Gaussian random variable (see *e.g.*, (5.56)-(5.57)), whose  $q$ th moments are also integrable for any fixed  $q$  if  $N$  is large enough.

The main technical task in the proof of (5.36) is the following recursive moment estimate.

**Lemma 5.3** (Recursive moment estimate for  $\mathcal{P}_{ii}$  and  $\mathcal{K}_{ii}$ ). *Suppose the assumptions of Theorem 5.2 hold. For any fixed integer  $p \geq 1$ , and for any  $i \in \llbracket 1, N \rrbracket$ , we have*

$$\begin{aligned}\mathbb{E}[\mathbf{m}_i(p, p)] &= \mathbb{E}[O_{\prec}(\Psi)\mathbf{m}_i(p-1, p)] + \mathbb{E}[O_{\prec}(\Psi^2)\mathbf{m}_i(p-2, p)] \\ &\quad + \mathbb{E}[O_{\prec}(\Psi^2)\mathbf{m}_i(p-1, p-1)], \\ \mathbb{E}[\mathbf{n}_i(p, p)] &= \mathbb{E}[O_{\prec}(\Psi)\mathbf{n}_i(p-1, p)] + \mathbb{E}[O_{\prec}(\Psi^2)\mathbf{n}_i(p-2, p)] \\ &\quad + \mathbb{E}[O_{\prec}(\Psi^2)\mathbf{n}_i(p-1, p-1)],\end{aligned}\tag{5.39}$$

where we made the convention  $\mathbf{m}_i(0, 0) = \mathbf{n}_i(0, 0) = 1$  and  $\mathbf{m}_i(-1, 1) = \mathbf{n}_i(-1, 1) = 0$  if  $p = 1$ .

*Proof of Lemma 5.3.* According to the decomposition in (5.7), for  $i \in \llbracket 1, N \rrbracket$ , we have

$$(\tilde{B}G)_{ii} = \hat{\mathbf{e}}_i^* \mathcal{R}_i \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i = -((\mathbf{h}_i^u)^*, \mathbf{0}^*) \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i = -(\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i,\tag{5.40}$$

where in the second step we used (5.13), and in the last step we used the notation in (5.15).

Using (5.40), the definition in (5.3), and also the identity in (5.11), one can check

$$\begin{aligned}(\tilde{B}G)_{ii} &= -(\mathbf{k}_i^u)^* \tilde{B}^{(i)} (\hat{I} - \mathbf{r}_i^u (\mathbf{r}_i^u)^* \oplus \mathbf{r}_i^v (\mathbf{r}_i^v)^*) G \hat{\mathbf{e}}_i \\ &= -S_{ii} + (\mathbf{k}_i^u)^* \tilde{B}^{(i)} (\mathbf{r}_i^u (\mathbf{r}_i^u)^* \oplus \mathbf{r}_i^v (\mathbf{r}_i^v)^*) G \hat{\mathbf{e}}_i \\ &= -S_{ii} + (\mathbf{k}_i^u)^* \tilde{B}^{(i)} (0 \oplus \mathbf{r}_i^v (\mathbf{r}_i^v)^*) G \hat{\mathbf{e}}_i \\ &= -S_{ii} + (\ell_i^v)^2 (\mathbf{k}_i^u)^* \tilde{B}^{(i)} (\hat{\mathbf{e}}_i + \mathbf{k}_i^v) (\hat{\mathbf{e}}_i + \mathbf{k}_i^v)^* G \hat{\mathbf{e}}_i \\ &= -S_{ii} + (\ell_i^v)^2 (\tilde{\sigma}_i h_{ii}^u + (\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v) (G_{ii}^* + T_{ii}^*) \\ &=: -\hat{S}_{ii} + \varepsilon_{i1},\end{aligned}\tag{5.41}$$

where 0 in the third line is the  $N \times N$  zero matrix, and

$$\varepsilon_{i1} := \left( (\ell_i^v)^2 - 1 \right) \tilde{\sigma}_i h_{ii}^u + (\ell_i^v)^2 (\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v G_{ii}^* + (\ell_i^v)^2 (\tilde{\sigma}_i h_{ii}^u + (\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v) T_{ii}^*.\tag{5.42}$$

In the third step of (5.41) we used the fact  $(\mathbf{k}_i^u)^* \tilde{B}^{(i)} (\mathbf{r}_i^u (\mathbf{r}_i^u)^* \oplus 0) = 0$  which follows from the definition of  $\mathbf{k}_i^u$  and  $\tilde{B}^{(i)}$  in (5.15) and (5.8); in the fifth step we used the second identity in (5.18); and in the last step, we used (5.17). We note that

$$|\varepsilon_{i1}| \prec \Psi,\tag{5.43}$$

where we used (5.31) and the large deviation bound (A.1) to show that  $(\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v \prec N^{-1/2}$ .

Using integration by parts, we note that

$$\int_{\mathbb{C}} \bar{g} f(g, \bar{g}) e^{-\frac{|g|^2}{\sigma^2}} d^2 g = \sigma^2 \int_{\mathbb{C}} \partial_g f(g, \bar{g}) e^{-\frac{|g|^2}{\sigma^2}} d^2 g,\tag{5.44}$$

for differentiable functions  $f : \mathbb{C}^2 \rightarrow \mathbb{C}$  (recall that  $d^2 g$  is the Lebesgue measure on  $\mathbb{C}$ ).

According to the definitions in (5.3), (5.10), and the identity (5.11), one can check for  $k \neq i$ ,

$$\frac{\partial R_i^a}{\partial g_{ik}^a} = -\frac{(\ell_i^a)^2}{\|\mathbf{g}_i^a\|_2} \mathbf{e}_k (\mathbf{e}_i + \mathbf{h}_i^a)^* + \Delta_R^a(i, k), \quad a = u, v.\tag{5.45}$$

where

$$\begin{aligned}\Delta_R^a(i, k) &:= \frac{(\ell_i^a)^2}{2\|\mathbf{g}_i^a\|_2^2} \bar{g}_{ik}^a (\mathbf{e}_i (\mathbf{h}_i^a)^* + \mathbf{h}_i^a \mathbf{e}_i^* + 2\mathbf{h}_i^a (\mathbf{h}_i^a)^*) \\ &\quad - \frac{(\ell_i^a)^4}{2\|\mathbf{g}_i^a\|_2^3} g_{ii}^a \bar{g}_{ik}^a (\mathbf{e}_i + \mathbf{h}_i^a) (\mathbf{e}_i + \mathbf{h}_i^a)^*, \quad a = u, v.\end{aligned}\tag{5.46}$$



The  $\Delta_R^a(i, k)$ 's are irrelevant error terms. Their estimates will be presented separately in Appendix B. For convenience, we set for  $a = u, v$ ,

$$c_i^a := \frac{(\ell_i^a)^2}{\|\mathbf{g}_i^a\|_2} = \frac{1}{\|\mathbf{g}_i^a\|_2} - h_{ii}^a + O_{\prec}\left(\frac{1}{N}\right) = \|\mathbf{g}_i^a\|_2 - h_{ii}^a - (\|\mathbf{g}_i^a\|_2^2 - 1) + O_{\prec}\left(\frac{1}{N}\right), \quad (5.47)$$

where the last step follows from (5.12). Using (5.7), we have for  $k \neq i$

$$\frac{\partial G}{\partial g_{ik}^u} = -G \frac{\partial \tilde{B}}{\partial g_{ik}^u} G = -G \frac{\partial \mathcal{R}_i}{\partial g_{ik}^u} \tilde{B}^{(i)} \mathcal{R}_i G - G \mathcal{R}_i \tilde{B}^{(i)} \frac{\partial \mathcal{R}_i}{\partial g_{ik}^u} G. \quad (5.48)$$

According to (5.45) and the fact  $\mathcal{R}_i = R_i^u \oplus R_i^v$ , we have

$$\frac{\partial \mathcal{R}_i}{\partial g_{ik}^u} = -c_i^u \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* + \Delta_R^a(i, k) \oplus 0, \quad (5.49)$$

where 0 is the  $N \times N$  zero matrix. We also used that  $\partial R_i^v / \partial g_{ik}^u = 0$ . Plugging (5.49) into (5.48), for  $k \neq i$ , we can write

$$\frac{\partial G}{\partial g_{ik}^u} = c_i^u G \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i^* + (\mathbf{k}_i^u)^*) \tilde{B}^{(i)} \mathcal{R}_i G + c_i^u G \mathcal{R}_i \tilde{B}^{(i)} \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i^* + (\mathbf{k}_i^u)^*) G + \Delta_G^u(i, k), \quad (5.50)$$

where we set

$$\Delta_G^u(i, k) := -G(\Delta_R^u(i, k) \oplus 0) \tilde{B}^{(i)} \mathcal{R}_i G - G \mathcal{R}_i \tilde{B}^{(i)} (\Delta_R^u(i, k) \oplus 0) G. \quad (5.51)$$

With the above derivatives, we are ready to apply the integration by parts formula in (5.44). We start with the following

$$\begin{aligned} \mathbb{E}[\mathbf{m}_i(p, p)] &= \mathbb{E}[\mathcal{P}_{ii} \mathbf{m}_i(p-1, p)] = \mathbb{E}[(\tilde{B}G)_{ii} \tau_1(G) \mathbf{m}_i(p-1, p)] \\ &\quad + \mathbb{E}[( -G_{ii} \tau_1(\tilde{B}G) + (G_{ii} + T_{ii}) \Upsilon_1) \mathbf{m}_i(p-1, p)], \end{aligned} \quad (5.52)$$

$$\begin{aligned} \mathbb{E}[\mathbf{n}_i(p, p)] &= \mathbb{E}[\mathcal{K}_{ii} \mathbf{n}_i(p-1, p)] = \mathbb{E}[T_{ii} \mathbf{n}_i(p-1, p)] \\ &\quad + \mathbb{E}[(\tau_1(G) (\tilde{\sigma}_i T_{ij} + (\tilde{B}G)_{ij}) - \tau_1(G \tilde{B}) (G_{ij} + T_{ij})) \mathbf{n}_i(p-1, p)], \end{aligned} \quad (5.53)$$

which follow from the definitions in (5.19) and (5.37) directly. From (5.41) and (5.17), we have

$$\begin{aligned} \mathbb{E}[(\tilde{B}G)_{ii} \tau_1(G) \mathbf{m}_i(p-1, p)] &= -\mathbb{E}[\dot{S}_{ii} \tau_1(G) \mathbf{m}_i(p-1, p)] \\ &\quad + \mathbb{E}[\varepsilon_{i1} \tau_1(G) \mathbf{m}_i(p-1, p)], \end{aligned} \quad (5.54)$$

$$\mathbb{E}[T_{ii} \mathbf{n}_i(p-1, p)] = \mathbb{E}[\dot{T}_{ii} \mathbf{n}_i(p-1, p)] + \mathbb{E}[O_{\prec}(\Psi) \mathbf{n}_i(p-1, p)], \quad (5.55)$$

where we used the fact  $|h_{ii}| \prec N^{-\frac{1}{2}}$ , and also (5.31).

Now we will carefully compute the first terms in the right hand side of (5.54) and (5.55) with the integration by parts formula since both  $\dot{S}_{ii}$  and  $\dot{T}_{ii}$  explicitly contain a multiplicative Gaussian factor. We will then find that the leading term of the result of this calculation will exactly cancel the last quantities in the right side of equations in (5.52) and (5.53). This cancellation is the key point of the following tedious calculation and this is the main reason for defining the key quantities  $\mathcal{P}_{ii}$  and  $\mathcal{K}_{ii}$  in the form they are given in (5.19).

For the first term on the right side of (5.54), using the definition of  $\dot{S}_{ii}$  in (5.17) and the integration by parts formula in (5.44), we have

$$\begin{aligned} \mathbb{E}[\dot{S}_{ii} \tau_1(G) \mathbf{m}_i(p-1, p)] &= \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \mathbf{m}_i(p-1, p) \right] \\ &= \frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \frac{\partial (\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \tau_1(G) \mathbf{m}_i(p-1, p) \right] \\ &\quad + \frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \mathbf{m}_i(p-1, p) \right] \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \frac{\partial \tau_1(G)}{\partial g_{ik}^u} \mathbf{m}_i(p-1, p) \right] \\
& + \frac{p-1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \frac{\partial \mathcal{P}_{ii}}{\partial g_{ik}^u} \mathbf{m}_i(p-2, p) \right] \\
& + \frac{p}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \frac{\partial \overline{\mathcal{P}}_{ii}}{\partial g_{ik}^u} \mathbf{m}_i(p-1, p-1) \right]. \quad (5.56)
\end{aligned}$$

Analogously, we have

$$\begin{aligned}
\mathbb{E}[T_{ii}^\circ \mathbf{n}_i(p-1, p)] & = \frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \mathbf{n}_i(p-1, p) \right] \\
& + \frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \mathbf{n}_i(p-1, p) \right] \\
& + \frac{p-1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \frac{\partial \mathcal{K}_{ii}}{\partial g_{ik}^u} \mathbf{n}_i(p-2, p) \right] \\
& + \frac{p}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \frac{\partial \overline{\mathcal{K}}_{ii}}{\partial g_{ik}^u} \mathbf{n}_i(p-1, p-1) \right]. \quad (5.57)
\end{aligned}$$

We start from the first term on the right side of (5.56). Using (5.50), we have

$$\begin{aligned}
\frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} & = c_i^u \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i \\
& + c_i^u \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \mathcal{R}_i \tilde{B}^{(i)} \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* G \hat{\mathbf{e}}_i + \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} \Delta_G^u(i, k) \hat{\mathbf{e}}_i. \quad (5.58)
\end{aligned}$$

Let

$$\varepsilon_{i2} := \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} \Delta_G^u(i, k) \hat{\mathbf{e}}_i. \quad (5.59)$$

Note that

$$\frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_k = \tau_1(\tilde{B}^{(i)} G) - \frac{1}{N} (\tilde{B}^{(i)} G)_{ii} = \tau_1(\tilde{B} G) + O_{\prec}(\Psi^2), \quad (5.60)$$

where in the last step we used the second estimate in Corollary A.4 with the choice  $Q = \hat{I}_1$  (c.f., (4.25)),  $(\tilde{B}^{(i)} G)_{ii} = \tilde{\sigma}_i G_{ii}$  (c.f., (5.18)), and the bound in (5.31). Analogously, one shows

$$\frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \mathcal{R}_i \tilde{B}^{(i)} \hat{\mathbf{e}}_k = \tau_1(\tilde{B} G \tilde{B}) + O_{\prec}(\Psi^2). \quad (5.61)$$

Moreover, using (5.18), (5.13) and the fact  $\mathcal{R}_i^2 = \hat{I}$ , we also have the following observations

$$\begin{aligned}
\hat{\mathbf{e}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i & = \tilde{\sigma}_i \hat{\mathbf{e}}_i^* \mathcal{R}_i G \hat{\mathbf{e}}_i = -\tilde{\sigma}_i (\mathbf{k}_i^u)^* G \hat{\mathbf{e}}_i = -\tilde{\sigma}_i T_{ii}, \\
(\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i & = (\mathbf{k}_i^u)^* \mathcal{R}_i \tilde{B} G \hat{\mathbf{e}}_i = -\hat{\mathbf{e}}_i^* \tilde{B} G \hat{\mathbf{e}}_i = -(\tilde{B} G)_{ii}. \quad (5.62)
\end{aligned}$$

Plugging (5.60), (5.61) and (5.62) into (5.58), we obtain

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} &= -c_i^u \tau_1(\tilde{B}G) (\tilde{\sigma}_i T_{ii} + (\tilde{B}G)_{ii}) \\ &\quad + c_i^u \tau_1(\tilde{B}G\tilde{B})(G_{ii} + T_{ii}) + \varepsilon_{i2} + O_{\prec}(\Psi^2). \end{aligned} \quad (5.63)$$

Analogously to (5.63), we also have

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} &= -c_i^u \tau_1(G) (\tilde{\sigma}_i T_{ii} + (\tilde{B}G)_{ii}) \\ &\quad + c_i^u \tau_1(G\tilde{B})(G_{ii} + T_{ii}) + \varepsilon_{i3} + O_{\prec}(\Psi^2), \end{aligned} \quad (5.64)$$

where

$$\varepsilon_{i3} := \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \Delta_G^u(i, k) \hat{\mathbf{e}}_i.$$

The following estimates on  $\varepsilon_{i2}$  and  $\varepsilon_{i3}$  will be proved in Lemma B.1 in Appendix B.

$$|\varepsilon_{i2}| \prec \Psi^2, \quad |\varepsilon_{i3}| \prec \Psi^2. \quad (5.65)$$

Combining (5.63), (5.64) with an appropriate linear combination and using (5.65), we get

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \tau_1(G) - \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \tau_1(\tilde{B}G) \\ = -c_i^u (G_{ii} + T_{ii}) \left( \tau_1(\tilde{B}G) - \Upsilon_1 \right) + O_{\prec}(\Psi^2). \end{aligned} \quad (5.66)$$

Here we also used that the tracial quantities  $\tau_1(G)$ ,  $\tau_1(\tilde{B}G)$  and  $\Upsilon_1$  are stochastically dominated by 1, in light of (5.34). Applying (5.47), the fact  $\dot{T}_{ii} = T_{ii} - h_{ii}^u G_{ii}$  from (5.17), we can write

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \tau_1(G) &= -c_i^u (G_{ii} + T_{ii}) \left( \tau_1(\tilde{B}G) - \Upsilon_1 \right) \\ &\quad + \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \tau_1(\tilde{B}G) + O_{\prec}(\Psi^2) \\ &= -c_i^u (G_{ii} + T_{ii}) \left( \tau_1(\tilde{B}G) - \Upsilon_1 \right) + \dot{T}_{ii} \tau_1(\tilde{B}G) \\ &\quad + \left( \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} - \dot{T}_{ii} \right) \tau_1(\tilde{B}G) + O_{\prec}(\Psi^2) \\ &= -\|\mathbf{g}_i^u\|_2 \left( G_{ii} \tau_1(\tilde{B}G) - (G_{ii} + T_{ii}) \Upsilon_1 \right) + \left( \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} - \dot{T}_{ii} \right) \tau_1(\tilde{B}G) \\ &\quad + \varepsilon_{i4} + \varepsilon_{i5} + O_{\prec}(\Psi^2), \end{aligned} \quad (5.67)$$

where

$$\begin{aligned} \varepsilon_{i4} &:= \left( (1 - \|\mathbf{g}_i^u\|_2^2) \tau_1(\tilde{B}G) + (1 - \|\mathbf{g}_i^u\|_2^2 - h_{ii}) (\tau_1(\tilde{B}G) - \Upsilon_1) \right) T_{ii} \\ &\quad + (1 - \|\mathbf{g}_i^u\|_2^2 - h_{ii}) G_{ii} \Upsilon_1, \end{aligned} \quad (5.68)$$

$$\varepsilon_{i5} := (\|\mathbf{g}_i^u\|_2^2 - 1) G_{ii} \tau_1(\tilde{B}G). \quad (5.69)$$

Using  $\|\mathbf{g}_i^u\|_2 = 1 + O_{\prec}(\frac{1}{\sqrt{N}})$ , the estimates (5.31), (5.32) and (5.34), and Corollary A.4, we get

$$|\varepsilon_{i4}| \prec \frac{1}{\sqrt{N}}, \quad |\varepsilon_{i5}| \prec \frac{1}{\sqrt{N}}. \quad (5.70)$$

Notice that the first term in the right side of (5.67) will exactly cancel the explicit last term in the right side of (5.52). This cancellation is one of the main reasons behind the choice of the auxiliary quantity  $\mathcal{P}$ . Combining the first equation of (5.53), (5.54), (5.56) with (5.67), we get

$$\begin{aligned}
\mathbb{E}[\mathbf{m}_i(p, p)] &= \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \left(\hat{T}_{ii} - \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u}\right) \tau_1(\tilde{B}G) \mathbf{m}_i(p-1, p)\right] \\
&\quad - \frac{1}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{\partial\|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \mathbf{m}_i(p-1, p)\right] \\
&\quad - \frac{1}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \frac{\partial\tau_1(G)}{\partial g_{ik}^u} \mathbf{m}_i(p-1, p)\right] \\
&\quad - \frac{p-1}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \frac{\partial\mathcal{P}_{ii}}{\partial g_{ik}^u} \mathbf{m}_i(p-2, p)\right] \\
&\quad - \frac{p}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \frac{\partial\overline{\mathcal{P}}_{ii}}{\partial g_{ik}^u} \mathbf{m}_i(p-1, p-1)\right] \\
&\quad + \mathbb{E}\left[\left(\varepsilon_{i1} \tau_1(G) - \frac{\varepsilon_{i4} + \varepsilon_{i5}}{\|\mathbf{g}_i^u\|_2}\right) \mathbf{m}_i(p-1, p)\right] + \mathbb{E}[O_{\prec}(\Psi^2) \mathbf{m}_i(p-1, p)]. \quad (5.71)
\end{aligned}$$

Note that the sixth term on the right side can be estimated by  $\mathbb{E}[O_{\prec}(\Psi) \mathbf{m}_i(p-1, p)]$ , according to (5.43) and (5.70). This estimate is sufficient for the proof of Lemma 5.3. But here we keep the  $\varepsilon$ -terms explicit for further use.

In order to estimate the first term in the right side, similarly to (5.57), we can apply the integration by parts formula (5.44) to obtain

$$\begin{aligned}
&\mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \left(\hat{T}_{ii} - \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u}\right) \tau_1(\tilde{B}G) \mathbf{m}_i(p-1, p)\right] \\
&= \frac{1}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{\partial\|\mathbf{g}_i^u\|_2^{-2}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \tau_1(\tilde{B}G) \mathbf{m}_i(p-1, p)\right] \\
&\quad + \frac{p-1}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \frac{\partial\tau_1(\tilde{B}G)}{\partial g_{ik}^u} \mathbf{m}_i(p-1, p)\right] \\
&\quad + \frac{p-1}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \tau_1(\tilde{B}G) \frac{\partial\mathcal{P}_{ii}}{\partial g_{ik}^u} \mathbf{m}_i(p-2, p)\right] \\
&\quad + \frac{p}{N} \sum_k^{(i)} \mathbb{E}\left[\frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \tau_1(\tilde{B}G) \frac{\partial\overline{\mathcal{P}}_{ii}}{\partial g_{ik}^u} \mathbf{m}_i(p-1, p-1)\right]. \quad (5.72)
\end{aligned}$$

Notice the cancellation between the two terms in the bracket in the first line.

Next we consider the estimate of  $\mathbf{n}_i(p, p)$ ; especially we control the first term in the right side of (5.57). In addition, using (5.64), (5.65), and the facts  $\|\mathbf{g}_i^u\|_2 = 1 + O_{\prec}(\frac{1}{\sqrt{N}})$  and  $c_i^u = 1 + O_{\prec}(\frac{1}{\sqrt{N}})$ , we have

$$\frac{1}{N} \frac{1}{\|\mathbf{g}_i^u\|_2} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} = -\tau_1(G) (\tilde{\sigma}_i T_{ii} + (\tilde{B}G)_{ii}) + \tau_1(G \tilde{B}) (G_{ii} + T_{ii}) + O_{\prec}(\Psi). \quad (5.73)$$

Note that the result of this calculation exactly cancels the second term in the right side of (5.53). Hence, analogously to (5.71), combining (5.57), (5.65), (5.55), (5.53) and (5.73), we get

$$\begin{aligned}
\mathbb{E}[\mathbf{n}_i(p, p)] &= \frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \mathbf{n}_i(p-1, p) \right] \\
&\quad + \frac{p-1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \frac{\partial \mathcal{K}_{ii}}{\partial g_{ik}^u} \mathbf{n}_i(p-2, p) \right] \\
&\quad + \frac{p}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \frac{\partial \overline{\mathcal{K}}_{ii}}{\partial g_{ik}^u} \mathbf{n}_i(p-1, p-1) \right] \\
&\quad + \mathbb{E} \left[ O_{\prec}(\Psi) \mathbf{n}_i(p-1, p) \right]. \tag{5.74}
\end{aligned}$$

Hence, to prove the second equation of (5.39), it suffices to estimate the first three terms on the right side of (5.74). For the first equation of (5.39), with (5.43) and (5.70), it suffices to estimate the second to the fifth terms on the right side of (5.71), and the terms on the right side of (5.72). All these estimates can be derived from the following lemma.

**Lemma 5.4.** *Suppose that the assumptions in Theorem 5.2 hold. Set  $X_i = \hat{I}$  or  $\tilde{B}^{(i)}$ . Let  $Q$  be any deterministic diagonal matrix satisfying  $\|Q\| \leq C$  and  $X = \hat{I}$  or  $A$ . We have the following estimates*

$$\begin{aligned}
\frac{1}{N} \sum_k^{(i)} \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= O_{\prec} \left( \frac{1}{N} \right), & \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* X \frac{\partial G}{\partial g_{ik}^u} \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= O_{\prec}(\Psi^2), \\
\frac{1}{N} \sum_k^{(i)} \frac{\partial T_{ji}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= O_{\prec}(\Psi^2), & \frac{1}{N} \sum_k^{(i)} \frac{\partial \text{tr} Q X G}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= O_{\prec}(\Psi^4), \tag{5.75}
\end{aligned}$$

where  $j = i$  or  $\hat{i}$  in the third equation.

Assuming the validity of Lemma 5.4, we continue with the proof of Lemma 5.3. Recall that our task is to bound the terms on the right sides of (5.71), (5.72), (5.74). The second term in (5.71), the first term in (5.72) and the first term in (5.74) can all be estimated with the aid of first bound in (5.75). The estimates for the third term in (5.71) and the second term in (5.72) follow from the last bound in (5.75). Finally, the fourth term in (5.71), the third term in (5.72) and the second term in (5.74) together with their complex conjugate analogues can be estimated in a similar way, so we only present the details for the fourth term on the right side of (5.71) in the sequel.

Recall the definition of  $\mathcal{P}_{ii}$  from (5.19)

$$\mathcal{P}_{ii} = (\tilde{B}G)_{ii} \tau_1(G) - G_{ii} \tau_1(\tilde{B}G) + (G_{ii} + T_{ii}) \Upsilon_1.$$

Using (4.16), and recalling the definition of  $\Upsilon_1$  in (5.20), we can see that  $\mathcal{P}_{ii}$  is a combination of the terms of the following forms:  $T_{ii}$ ,  $(XG)_{ii}$  and  $\text{tr}(QXG)$ , for  $X = \hat{I}$  or  $A$ , and  $Q$  is certain deterministic diagonal matrix with  $\|Q\| \leq C$  for some positive constant  $C$ . For example:  $(\tilde{B}G)_{ii} = 1 + zG_{ii} - (AG)_{ii}$ , and

$$\begin{aligned}
\tau_1(G\tilde{B}) &= \tau_1(\hat{I} - G(A - z)) = 1 + z\tau_1(G) - \tau_1(GA) \\
&= 1 + 2z\text{tr}(\hat{I}_1 G) - 2\text{tr}(A\hat{I}_1 G) = 1 + 2z\text{tr}(\hat{I}_1 G) - 2\text{tr}(\hat{I}_2 AG).
\end{aligned}$$

Then, by the product rule for derivative, and the boundedness of all the partial traces (*c.f.*, (5.34)) and entries (*c.f.*, (5.31), (5.32)), we can apply the last three bounds in (5.75) to conclude that the fourth term on the right side of (5.71) is  $\mathbb{E}[O_{\prec}(\Psi^2) \mathbf{m}_i(p-2, p)]$ .

This completes the proof of Lemma 5.3, up to Lemma 5.4.  $\square$

*Proof of Lemma 5.4.* Since the sums in (5.75) are over  $k \neq i$ , it will be convenient to work in this proof with the following notations

$$I^{(i)} := I - \mathbf{e}_i \mathbf{e}_i^*, \quad \hat{I}_1^{(i)} := I^{(i)} \oplus 0, \quad (5.76)$$

where 0 is the  $N \times N$  zero matrix. We check the estimates in (5.75) one by one. For the first estimate, we have

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= -\frac{1}{2N} \frac{1}{\|\mathbf{g}_i^u\|_2^3} \sum_k^{(i)} \bar{g}_{ik}^u \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i \\ &= -\frac{1}{2N} \frac{1}{\|\mathbf{g}_i^u\|_2^2} (\hat{\mathbf{k}}_i^u)^* X_i G \hat{\mathbf{e}}_i = O_{\prec} \left( \frac{1}{N} \right), \end{aligned}$$

where in the last step we used that

$$(\hat{\mathbf{k}}_i^u)^* X_i G \hat{\mathbf{e}}_i \prec 1, \quad (5.77)$$

which would follow once we show  $|\hat{S}_{ii}| \prec 1$  and  $|\hat{T}_{ii}| = |T_{ii} - h_{ii}^u G_{ii}| \prec 1$  by (5.17). Since  $\hat{S}_{ii} = -(\tilde{B}G)_{ii} + O_{\prec}(\Psi)$  by (5.41), (5.43) and  $|(\tilde{B}G)_{ii}| \prec 1$  from (5.32), we get  $|\hat{S}_{ii}| \prec 1$ . The estimate  $|\hat{T}_{ii}| \prec 1$  follows from (5.31) and the fact  $|h_{ii}^u| \prec \frac{1}{\sqrt{N}}$ .

Next, we show the second estimate in (5.75). Using (5.50), we have

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_i^* X \frac{\partial G}{\partial g_{ik}^u} \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= c_i^u \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_i^* X G \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i \\ &\quad + c_i^u \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_i^* X G \mathcal{R}_i \tilde{B}^{(i)} \hat{\mathbf{e}}_k (\hat{\mathbf{e}}_i^* + (\mathbf{k}_i^u)^*) G \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i \\ &\quad + \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_i^* X \Delta_G^u(i, k) \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i \\ &= c_i^u \frac{1}{N} \hat{\mathbf{e}}_i^* X G \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i \\ &\quad + c_i^u \frac{1}{N} \hat{\mathbf{e}}_i^* X G \mathcal{R}_i \tilde{B}^{(i)} \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* G \hat{\mathbf{e}}_i \\ &\quad + \frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_i^* X \Delta_G^u(i, k) \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i, \end{aligned} \quad (5.78)$$

where we have used the notation introduced in (5.76).

From Lemma B.1 in Appendix B, we see that the last term on the right side of (5.78) is of order  $O_{\prec}(\Psi^2)$ . For the first two terms, we first claim that

$$|\hat{\mathbf{e}}_i^* X G \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i| \prec \frac{1}{\eta}, \quad |\hat{\mathbf{e}}_i^* X G \mathcal{R}_i \tilde{B}^{(i)} \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i| \prec \frac{1}{\eta}. \quad (5.79)$$

We prove the first estimate (5.79) as follows. Note that

$$\begin{aligned} \hat{\mathbf{e}}_i^* X G \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i &\leq \hat{\mathbf{e}}_i^* X |G|^2 X \hat{\mathbf{e}}_i + \hat{\mathbf{e}}_i^* G^* X_i \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i \\ &\leq \frac{1}{\eta} \text{Im}(XGX)_{ii} + \|X_i\|_2^2 \frac{1}{\eta} \text{Im} G_{ii}. \end{aligned} \quad (5.80)$$

Recall  $X = \hat{I}$  or  $A$ , and the fact  $(AGA)_{ii} = |\sigma_i|^2 G_{ii}^*$ . This together with (5.31) and the fact  $\|X_i\| \leq C$  since  $X_i = \hat{I}$  or  $\tilde{B}^{(i)}$  implies the first estimate in (5.79). The second estimate can be derived in a similar way.

Then, we recall from (5.62) that  $(\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i = -\tilde{\sigma}_i T_{ii} - (\tilde{B}G)_{ii}$ , and from the definition of  $T_{ij}$  in (5.16) that  $(\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* G \hat{\mathbf{e}}_i = G_{ii} + T_{ii}$ , which together with (5.31), (5.32)

and (5.79) imply that the first two terms on the right side of (5.78) are also of order  $O_{\prec}(\Psi^2)$ . This completes the second estimate in (5.75).

For the third estimate in (5.75), we present the details for  $j = i$  in the sequel. The case of  $j = \hat{i}$  is similar but simpler and we omit it. According to the definition of  $T_{ii}$  in (5.16), it suffices to show

$$\frac{1}{N} \sum_k^{(i)} \frac{\partial(\mathbf{k}_i^u)^*}{\partial g_{ik}^u} G \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i = O_{\prec}\left(\frac{1}{N}\right), \quad \frac{1}{N} \sum_k^{(i)} (\mathbf{k}_i^u)^* \frac{\partial G}{\partial g_{ik}^u} \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i = O_{\prec}(\Psi^2). \quad (5.81)$$

For the first estimate in (5.81), we have

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial(\mathbf{k}_i^u)^*}{\partial g_{ik}^u} G \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= -\frac{1}{2\|\mathbf{g}_i^u\|_2^2} \frac{1}{N} \sum_k^{(i)} \bar{h}_{ik}^u \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i (\mathbf{k}_i^u)^* G \hat{\mathbf{e}}_i \\ &= -\frac{1}{2\|\mathbf{g}_i^u\|_2^2} \frac{1}{N} (\mathbf{k}_i^u)^* X_i G \hat{\mathbf{e}}_i (\mathbf{k}_i^u)^* G \hat{\mathbf{e}}_i = O_{\prec}\left(\frac{1}{N}\right), \end{aligned}$$

where in the last step we used (5.31) and (5.77). The proof of the second estimate in (5.81) is similar to that for the second estimate in (5.75). It suffices to go through the discussion from (5.78) to (5.80) again, with the vector  $\hat{\mathbf{e}}_i^* X$  replaced by  $(\mathbf{k}_i^u)^*$ . The main differences are: instead of the last term of (5.78), we have

$$\frac{1}{N} \sum_k^{(i)} (\mathbf{k}_i^u)^* \Delta_G^u(i, k) \hat{\mathbf{e}}_i \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i, \quad (5.82)$$

and instead of the first term on the right side of (5.80), we have

$$\frac{1}{\eta} \text{Im}(\mathbf{k}_i^u)^* G \mathbf{k}_i^u. \quad (5.83)$$

The bound on (5.82) is stated in (B.3). For (5.83), we recall the identity (5.13) which implies  $\mathbf{k}_i^u = -\mathcal{R}_i \hat{\mathbf{e}}_i$ , the fact  $G = \mathcal{U}G\mathcal{U}^*$ , together with (5.6) and the fact  $\mathcal{R}_i^2 = \hat{I}$ . Then we have

$$(\mathbf{k}_i^u)^* G \mathbf{k}_i^u = \hat{\mathbf{e}}_i^* \mathcal{R}_i \mathcal{U}G\mathcal{U}^* \mathcal{R}_i \hat{\mathbf{e}}_i = \hat{\mathbf{e}}_i^* \mathcal{U}_i \Phi_i G \Phi_i^* \mathcal{U}_i^* \hat{\mathbf{e}}_i = \mathcal{G}_{ii}. \quad (5.84)$$

Similarly to (5.31), with the second bound in assumption (5.27), we can also show that

$$\max_{k,l} |\mathcal{G}_{kl}| \prec 1. \quad (5.85)$$

With these bounds for (5.82) and (5.83), we can show the second estimate of (5.81), which together with the first estimate in (5.81) implies the third bound in (5.75).

At the end, we show the last bound in (5.75). Applying (5.50), we have

$$\begin{aligned} \frac{\partial \text{tr} QXG}{\partial g_{ik}^u} &= \frac{1}{N} c_i^u (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G QXG \hat{\mathbf{e}}_k \\ &\quad + \frac{1}{N} c_i^u (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* G QXG \mathcal{R}_i \tilde{B}^{(i)} \hat{\mathbf{e}}_k + \text{tr} QX \Delta_G^u(i, k). \end{aligned} \quad (5.86)$$

Summing over  $k$  and using the notation in (5.76), we can write

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial \text{tr} QXG}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= \frac{c_i^u}{N^2} (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G QXG \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i \\ &\quad + \frac{c_i^u}{N^2} (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* G QXG \mathcal{R}_i \tilde{B}^{(i)} \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i + \frac{1}{N} \sum_k^{(i)} \text{tr} QX \Delta_G^u(i, k) \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i. \end{aligned} \quad (5.87)$$

The bound for the last term of the right side of (5.87) can be found in (B.4).



In the sequel, we bound the first two terms on the right side of (5.87). We only present the details for the first one; the second is estimated analogously. First, similarly to (5.62), we have

$$(\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i = -(\tilde{\sigma}_i(\mathbf{k}_i^v)^* + \hat{\mathbf{e}}_i^* \tilde{B}).$$

Then we can write

$$\begin{aligned} \frac{c_i^u}{N^2} (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathcal{R}_i G Q X G \hat{I}_1^{(i)} X_i G \hat{\mathbf{e}}_i &= -\frac{c_i^u}{N^2} (\tilde{\sigma}_i(\mathbf{k}_i^v)^* + \hat{\mathbf{e}}_i^* \tilde{B}) G Q X G \hat{I}_1 X_i G \hat{\mathbf{e}}_i \\ &\quad + \frac{c_i^u}{N^2} (\tilde{\sigma}_i(\mathbf{k}_i^v)^* + \hat{\mathbf{e}}_i^* \tilde{B}) G Q X G \hat{\mathbf{e}}_i \hat{\mathbf{e}}_i^* X_i G \hat{\mathbf{e}}_i. \end{aligned} \quad (5.88)$$

For the second term on the right side of (5.88), we use the bounds

$$|(\tilde{\sigma}_i(\mathbf{k}_i^v)^* + \hat{\mathbf{e}}_i^* \tilde{B}) G Q X G \hat{\mathbf{e}}_i| \prec \eta^{-2}, \quad |\hat{\mathbf{e}}_i^* X_i G \hat{\mathbf{e}}_i| \prec 1, \quad (5.89)$$

where in the first inequality we used the trivial bound  $\|G\| \leq \eta^{-1}$ , while in the second inequality we used the fact that  $X_i = \hat{I}$  or  $\tilde{B}^{(i)}$ , together with (5.18), and the first bound in (5.31). Using the bounds in (5.89), we see that the second term on the right side of (5.88) is of order  $O_{\prec}(\Psi^4)$ .

Now, we turn to the first term on the right side of (5.88). Note that

$$\begin{aligned} \left| \frac{1}{N^2} (\tilde{\sigma}_i(\mathbf{k}_i^v)^* + \hat{\mathbf{e}}_i^* \tilde{B}) G Q X G \hat{I}_1 X_i G \hat{\mathbf{e}}_i \right| &\leq \frac{C}{N^2 \eta} (\|(\mathbf{k}_i^v)^* G\|_2 + \|\hat{\mathbf{e}}_i^* \tilde{B} G\|_2) \|G \hat{\mathbf{e}}_i\|_2 \\ &\leq \frac{C}{N^2 \eta} (\|(\mathbf{k}_i^v)^* G\|_2^2 + \|\hat{\mathbf{e}}_i^* \tilde{B} G\|_2^2 + \|G \hat{\mathbf{e}}_i\|_2^2) \\ &\leq \frac{C}{N^2 \eta^2} (\text{Im}(\mathbf{k}_i^v)^* G \mathbf{k}_i^v + \text{Im} \hat{\mathbf{e}}_i^* \tilde{B} G \tilde{B} \hat{\mathbf{e}}_i + \text{Im} \hat{\mathbf{e}}_i^* G \hat{\mathbf{e}}_i). \end{aligned} \quad (5.90)$$

Similarly to (5.84), we have

$$(\mathbf{k}_i^v)^* G \mathbf{k}_i^v = \hat{\mathbf{e}}_i^* \mathcal{R}_i \mathcal{U} G \mathcal{U}^* \mathcal{R}_i \hat{\mathbf{e}}_i = \hat{\mathbf{e}}_i^* \mathcal{U}_i \Phi_i G \Phi_i^* \mathcal{U}_i^* \hat{\mathbf{e}}_i = \mathcal{G}_{ii}. \quad (5.91)$$

Combining (5.90) and (5.91), we obtain

$$\begin{aligned} \left| \frac{1}{N^2} (\tilde{\sigma}_i(\mathbf{k}_i^v)^* + \hat{\mathbf{e}}_i^* \tilde{B}) G Q X G \hat{I}_1 X_i G \hat{\mathbf{e}}_i \right| &\leq \frac{C}{N^2 \eta^2} (\text{Im} \mathcal{G}_{ii} + \text{Im}(\tilde{B} G \tilde{B})_{ii} + \text{Im} G_{ii}) \\ &= O_{\prec}(\Psi^4), \end{aligned}$$

where we also used (5.32) and (5.85). Hence the first term on the right side of (5.87) is  $O_{\prec}(\Psi^4)$ . The second term on the right side of (5.87) is bounded similarly. These bounds together with (B.4) yield the other estimates in (5.75). This completes the proof of Lemma 5.4.  $\square$

**5.4. Local stability analysis: proof of Theorem 5.2.** Having established Lemma 5.3, we move on to the local stability analysis in order to conclude the proof of Theorem 5.2.

*Proof of Theorem 5.2.* Applying Young's inequality, we obtain from (5.39) that for any given (small)  $\varepsilon > 0$ ,

$$\mathbb{E}[\mathbf{m}_i(p, p)] \leq 3 \frac{1}{2p} \mathbb{E}[N^{2p\varepsilon} \Psi^{2p}] + 3 \frac{2p-1}{2p} N^{-\frac{2p\varepsilon}{2p-1}} \mathbb{E}[\mathbf{m}_i(p, p)],$$

which implies  $\mathbb{E}[\mathbf{m}_i(p, p)] \prec \Psi^{2p}$ . Hence, we conclude the proof of the first estimate of (5.36).

The second estimate of (5.36) can be proved in the same way, with the aid of the second equation in (5.39). Then, applying Markov's inequality we get the first and the third estimates of (5.28) with  $j = i$ . The others in (5.28) are proved in an analogous way. We omit the details.

Next, we show that (5.28) together with the assumption (5.27) imply (5.29). To this end, we first show the following crude bound

$$\Lambda_T(z) \prec N^{-\frac{7}{4}} \quad (5.92)$$

under the assumption (5.27). We need the following equations for  $j = i, \hat{i}$ ,

$$\begin{aligned} T_{ij} &= -\tau_1(G)(\bar{\sigma}_i T_{\hat{i}j} + (\tilde{B}G)_{ij}) + \tau_1(G\tilde{B})(G_{ij} + T_{ij}) + O_{\prec}(\Psi), \\ T_{\hat{i}j} &= -\tau_2(G)(\bar{\sigma}_i^* T_{ij} + (\tilde{B}G)_{\hat{i}j}) + \tau_2(G\tilde{B})(G_{\hat{i}j} + T_{\hat{i}j}) + O_{\prec}(\Psi), \end{aligned} \quad (5.93)$$

which is just a rewriting of the second line of (5.28), according to the definition in (5.19).

Using the first identity in (4.16) and the definition of  $A$  in (4.14), we have

$$\begin{aligned} (\tilde{B}G)_{ii} &= 1 + zG_{ii} - \xi_i G_{\hat{i}i}, & (\tilde{B}G)_{\hat{i}\hat{i}} &= -\xi_i G_{\hat{i}\hat{i}} + zG_{\hat{i}\hat{i}}, \\ (\tilde{B}G)_{\hat{i}i} &= -\bar{\xi}_i G_{ii} + zG_{\hat{i}i}, & (\tilde{B}G)_{i\hat{i}} &= 1 + zG_{i\hat{i}} - \bar{\xi}_i G_{i\hat{i}}. \end{aligned} \quad (5.94)$$

Applying the assumption on  $\Lambda_d$  in (5.27), and also the lower bound of  $\text{Im } \omega_B$  and the upper bound on  $|\omega_B|$  in (A.4), we can get from (5.94) that

$$\begin{aligned} (\tilde{B}G)_{ii} &= \frac{(z - \omega_B)\omega_B}{|\xi_i|^2 - \omega_B^2} + O_{\prec}(N^{-\frac{\gamma}{4}}), & (\tilde{B}G)_{\hat{i}\hat{i}} &= \frac{(z - \omega_B)\xi_i}{|\xi_i|^2 - \omega_B^2} + O_{\prec}(N^{-\frac{\gamma}{4}}), \\ (\tilde{B}G)_{\hat{i}i} &= \frac{(z - \omega_B)\bar{\xi}_i}{|\xi_i|^2 - \omega_B^2} + O_{\prec}(N^{-\frac{\gamma}{4}}), & (\tilde{B}G)_{i\hat{i}} &= \frac{(z - \omega_B)\omega_B}{|\xi_i|^2 - \omega_B^2} + O_{\prec}(N^{-\frac{\gamma}{4}}). \end{aligned} \quad (5.95)$$

This together with (5.34), leads to the following estimates for  $j = i, \hat{i}$ ,

$$\begin{aligned} -\tau_1(G)(\tilde{B}G)_{ij} + \tau_1(G\tilde{B})G_{ij} &= O_{\prec}(N^{-\frac{\gamma}{4}}), \\ -\tau_2(G)(\tilde{B}G)_{\hat{i}j} + \tau_2(G\tilde{B})G_{\hat{i}j} &= O_{\prec}(N^{-\frac{\gamma}{4}}), \end{aligned}$$

which together with (5.93) implies

$$\begin{aligned} (1 - \tau_1(G\tilde{B}))T_{ij} + \tau_1(G)\bar{\sigma}_i T_{\hat{i}j} &= O_{\prec}(N^{-\frac{\gamma}{4}}), \\ (1 - \tau_2(G\tilde{B}))T_{\hat{i}j} + \tau_2(G)\bar{\sigma}_i^* T_{ij} &= O_{\prec}(N^{-\frac{\gamma}{4}}), \quad j = i, \hat{i}. \end{aligned} \quad (5.96)$$

Solving  $T_{ij}$  from the equations in (5.96), we get

$$((1 - \tau_1(G\tilde{B}))(1 - \tau_2(G\tilde{B})) - |\sigma_i|^2 \tau_1(G)\tau_2(G))T_{ij} = O_{\prec}(N^{-\frac{\gamma}{4}}). \quad (5.97)$$

Using the assumption on  $\Lambda_T$  in (5.27), and also (5.34), we obtain from (5.97) that

$$((1 + (\omega_B - z)m_{\mu_A \boxplus \mu_B})^2 - |\sigma_i|^2 m_{\mu_A \boxplus \mu_B}^2)T_{ij} = O_{\prec}(N^{-\frac{\gamma}{4}}). \quad (5.98)$$

Further, observe that

$$(1 + (\omega_B - z)m_{\mu_A \boxplus \mu_B})^2 - |\sigma_i|^2 m_{\mu_A \boxplus \mu_B}^2 = m_{\mu_A \boxplus \mu_B}^2 (\omega_A - |\sigma_i|)(\omega_A + |\sigma_i|), \quad (5.99)$$

which follows from the second equation in (2.5) with  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$ . Then by (A.4) and the fact  $m_{\mu_A \boxplus \mu_B} = m_{\mu_A}(\omega_B)$ , we see that  $|T_{ij}| \prec N^{-\frac{\gamma}{4}}$  for  $j = i, \hat{i}$ . Analogously, one can show  $|T_{\hat{i}j}| \prec N^{-\frac{\gamma}{4}}$ . This completes the proof of the crude bound (5.92).

With (5.92), we can now proceed to the proof of (5.29). We consider the average of  $\mathcal{P}_{ii}$  over  $i \in \llbracket 1, N \rrbracket$ , and use (5.28) to obtain

$$\Upsilon_1 \cdot \frac{1}{N} \sum_{i=1}^N (G_{ii} + T_{ii}) = \frac{1}{N} \sum_{i=1}^N \mathcal{P}_{ii} = O_{\prec}(\Psi). \quad (5.100)$$

By the first estimate in (5.34), the fact  $m_{\mu_A \boxplus \mu_B} = m_{\mu_A}(\omega_B)$ , the lower bound on  $\text{Im } \omega_B$  in (A.4), and also the crude bound (5.92), we can see that

$$\left| \frac{1}{\frac{1}{N} \sum_{i=1}^N (G_{ii} + T_{ii})} \right| = \left| \frac{1}{m_{\mu_A}(\omega_B) + O_{\prec}(N^{-\frac{\gamma}{4}})} \right| \prec 1. \quad (5.101)$$

Then the first estimate in (5.29) follows from (5.100) and (5.101) immediately. The second one can be verified similarly.

Finally, using (5.28) and (5.29), we can prove (5.30) as follows. Recall the definition in (5.19). Applying (5.27)-(5.29), we obtain, for  $j = i, \hat{i}$ ,

$$(\tilde{B}G)_{ij} = G_{ij} \frac{\tau_1(\tilde{B}G)}{\tau_1(G)} + O_{\prec}(\Psi), \quad (\tilde{B}G)_{\hat{i}j} = G_{\hat{i}j} \frac{\tau_2(\tilde{B}G)}{\tau_2(G)} + O_{\prec}(\Psi). \quad (5.102)$$

Using (5.94) and (5.102) we get the following system of equations,

$$\begin{aligned} 1 - \xi_i G_{ii} + \omega_{B,1}^c G_{ii} &= O_{\prec}(\Psi), & -\xi_i G_{\hat{i}\hat{i}} + \omega_{B,1}^c G_{\hat{i}\hat{i}} &= O_{\prec}(\Psi), \\ -\bar{\xi}_i G_{ii} + \omega_{B,2}^c G_{ii} &= O_{\prec}(\Psi), & 1 - \bar{\xi}_i G_{\hat{i}\hat{i}} + \omega_{B,2}^c G_{\hat{i}\hat{i}} &= O_{\prec}(\Psi), \end{aligned} \quad (5.103)$$

where we used the notation introduced in (8.20). Solving (5.103) we find

$$\begin{aligned} G_{ii} &= \frac{\omega_{B,2}^c}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} + O_{\prec}(\Psi), & G_{\hat{i}\hat{i}} &= \frac{\xi_i}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} + O_{\prec}(\Psi), \\ G_{\hat{i}\hat{i}} &= \frac{\bar{\xi}_i}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} + O_{\prec}(\Psi), & G_{ii} &= \frac{\omega_{B,1}^c}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} + O_{\prec}(\Psi). \end{aligned} \quad (5.104)$$

From (5.34), we see that

$$\omega_{B,a}^c = \omega_B + O_{\prec}(N^{-\frac{3}{4}}), \quad a = 1, 2. \quad (5.105)$$

The first estimate of (5.30) could be verified from (5.104), if we could show

$$\omega_{B,a}^c = \omega_B^c + O_{\prec}(\Psi), \quad a = 1, 2. \quad (5.106)$$

To this end, we use  $\tau_1(G(z)) = \tau_2(G(z))$ ; *c.f.*, (4.27). From (8.20) and (4.27), we also have

$$\omega_{B,1}^c + \omega_{B,2}^c = 2\omega_B^c. \quad (5.107)$$

Then, averaging the first and the fourth equations of (5.104) over  $i \in \llbracket 1, N \rrbracket$ , we get

$$\omega_{B,2}^c \frac{1}{N} \sum_{i=1}^N \frac{1}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} = \omega_{B,1}^c \frac{1}{N} \sum_{i=1}^N \frac{1}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} + O_{\prec}(\Psi), \quad (5.108)$$

where we also used (4.27). We further claim that

$$\left( \frac{1}{N} \sum_{i=1}^N \frac{1}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \right)^{-1} \prec 1, \quad (5.109)$$

which together with (5.108) implies that

$$\omega_{B,2}^c = \omega_{B,1}^c + O_{\prec}(\Psi). \quad (5.110)$$

Combining (5.110) with (5.107), we get (5.106). Hence, it suffices to show (5.109). To this end, we use (5.105). Then we have

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \frac{1}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} &= \frac{1}{N} \sum_{i=1}^N \frac{1}{|\xi_i|^2 - \omega_B^2 + O_{\prec}(N^{-\frac{3}{4}})} \\ &= \frac{1}{N} \sum_{i=1}^N \frac{1}{|\xi_i|^2 - \omega_B^2} + O_{\prec}(N^{-\frac{3}{4}}) = \omega_B^{-1} m_{\mu_A}(\omega_B) + O_{\prec}(N^{-\frac{3}{4}}), \end{aligned}$$

where in the first step above, we used the upper bound of  $|\omega_B|$  in (A.4); in the second step, we used again the fact that  $|\xi_i|^2 - \omega_B^2$  is away from 0 due to the lower bound of  $\text{Im } \omega_B$  in (A.4); and the last step follows from (5.33). Then the fact  $\|A\| \leq C$  (*c.f.*, (4.4)), the lower bound of  $\text{Im } \omega_B$  and the upper bound on  $|\omega_B|$  in (A.4), we can get (5.109). Hence, we conclude the proof of the first estimate of (5.30).

For the second estimate in (5.30), we need to go through the proof of (5.92) again, but this time with the a priori input (5.27) replaced by the first estimate of (5.30). Therefore, with (5.30), we can get

$$\left( (1 + (\omega_B^c - z)m_A(\omega_B^c))^2 - |\sigma_i|^2(m_A(\omega_B^c))^2 \right) T_{ij} = O_{\prec}(\Psi), \quad (5.111)$$

which is the analogue of (5.98). Then, by the estimates in (5.34) and the definition in (5.22), it is not difficult to check that the coefficient of  $T_{ij}$  above can be approximated by (5.99), up to an error  $O_{\prec}(N^{-\frac{7}{4}})$ . Hence, we can improve the estimate to  $|T_{ij}| \prec \Psi$  for  $j = i, \hat{i}$ . Similarly, we can prove the same bound for  $T_{\hat{i}j}$ . This completes the second estimate of (5.30). Hence, we conclude the proof of Theorem 5.2.  $\square$

**5.5. Continuity argument: Proof of Theorem 5.1.** Having derived Theorem 5.2, we prove Theorem 5.1 using a continuity argument similar to [22].

*Proof of Theorem 5.1.* First, we show that  $\Lambda_d^c(z)$  in (5.30) can be replaced by  $\Lambda_d(z)$ . This means, we have to control the difference between  $(\omega_A, \omega_B)$  and  $(\omega_A^c, \omega_B^c)$  as described in (5.26); this estimate will follow from the stability of the system  $\Phi_{\mu_A, \mu_B}(\omega_A, \omega_B, z) = 0$ , (c.f., (4.29) with  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$ ). We will use the dual pair of subordination equations, i.e., when we analyze  $\mathcal{H}$  instead of  $H$ . Recall the notations introduced in (4.23), and also  $\tilde{\Lambda}_d$  and  $\tilde{\Lambda}_T$  as the analogue of  $\Lambda_d$  and  $\Lambda_T$ , respectively, see the explanation around (5.23). For any  $\delta \in [0, 1]$  and  $z \in \mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ , we introduce the following event

$$\Theta(z, \delta) := \{ \Lambda_d(z) \leq \delta, \quad \tilde{\Lambda}_d(z) \leq \delta, \quad \Lambda_T(z) \leq 1, \quad \tilde{\Lambda}_T(z) \leq 1 \}. \quad (5.112)$$

With the above notation, we have the following lemma.

**Lemma 5.5.** *Suppose that the assumptions in Theorem 4.3 hold. Let  $\eta_M > 0$  be a sufficiently large constant and  $\gamma > 0$  be a small constant in the definition (5.1). For any  $\varepsilon$  with  $0 < \varepsilon \leq \frac{\gamma}{8}$  and for any  $D > 0$ , there exists a positive integer  $N_2(D, \varepsilon)$  such that the following holds: For any fixed  $z \in \mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$  there exists an event  $\Omega(z) \equiv \Omega(z, D, \varepsilon)$  with*

$$\mathbb{P}(\Omega(z)) \geq 1 - N^{-D}, \quad \forall N \geq N_2(D, \varepsilon), \quad (5.113)$$

such that if the estimate

$$\mathbb{P}(\Theta(z, N^{-\frac{7}{4}})) \geq 1 - N^{-D}(1 + N^5(\eta_M - \eta)), \quad \eta = \text{Im } z, \quad (5.114)$$

holds for all  $D > 0$  and  $N \geq N_1(D, \gamma, \varepsilon)$ , for some threshold  $N_1(D, \gamma, \varepsilon)$ , then we also have

$$\Theta(z, N^{-\frac{7}{4}}) \cap \Omega(z) \subset \Theta(z, \frac{N^\varepsilon}{\sqrt{N}\eta}), \quad (5.115)$$

for all  $N \geq N_3(D, \gamma, \varepsilon) := \max \{ N_1(D, \gamma, \varepsilon), N_2(D, \varepsilon) \}$ .

*Proof of Lemma 5.5.* In this proof, we fix  $z \in \mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . According to the definition of  $\prec$  in Definition 1.6, we see from the assumption (5.114) that

$$\Lambda_d(z) \prec N^{-\frac{7}{4}}, \quad \tilde{\Lambda}_d(z) \prec N^{-\frac{7}{4}}, \quad \Lambda_T(z) \prec 1, \quad \tilde{\Lambda}_T(z) \prec 1. \quad (5.116)$$

We apply Theorem 5.2; by the estimates on  $\Lambda_d^c$  and on  $\Lambda_T$  in (5.30) and their analogues for  $\tilde{\Lambda}_d^c$  and  $\tilde{\Lambda}_T$ , we have

$$\Lambda_d^c(z) \prec \Psi, \quad \tilde{\Lambda}_d^c(z) \prec \Psi, \quad \Lambda_T(z) \prec \Psi, \quad \tilde{\Lambda}_T(z) \prec \Psi. \quad (5.117)$$

Now, we state the conclusions in (5.117) in a more explicit quantitative form, with the quantitative assumption (5.114). To this end, we need a more quantitative version of Lemma 5.3. Let  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  be a smooth cutoff function s.t

$$\varphi(x) = 1 \text{ if } |x| \leq K, \quad \varphi(x) = 0 \text{ if } |x| \geq 2K, \quad \sup_{x \in \mathbb{R}} |\varphi'(x)| \leq CK^{-1} \quad (5.118)$$

for some sufficiently large constant  $K > 0$ . Let

$$\begin{aligned} \Gamma_i \equiv \Gamma_i(z) := & \sum_{a,b=i,\hat{i}} (|G_{ab}|^2 + |\mathcal{G}_{ab}|^2 + |T_{ab}|^2 + |\mathcal{T}_{ab}|^2) \\ & + \sum_{a=1,2} (|\tau_a(G)|^2 + |\tau_a(\tilde{B}G)|^2 + |\tau_a(G\tilde{B})|^2 + |\tau_a(\tilde{B}G\tilde{B})|^2). \end{aligned} \quad (5.119)$$

Note that for a given  $i$ , all the a priori bounds we needed in the proof of Lemma 5.3 are the  $O_{\prec}(1)$  bound for  $G_{ab}$ ,  $\mathcal{G}_{ab}$ ,  $T_{ab}$ ,  $\mathcal{T}_{ab}$  with  $a, b = i, \hat{i}$  and the tracial quantities in (5.119). The  $O_{\prec}(1)$  bound for  $(XGY)_{ab}$  with  $X, Y = \hat{I}$  or  $\tilde{B}$  were also used (see  $(\tilde{B}X\tilde{B})_{ii}$  in (5.90) for instance), but they can be derived from the bound of  $G_{ab}$ 's by using (4.16). Recall the definitions of  $\mathbf{m}_i$  and  $\mathbf{n}_i$  in (5.37). We now introduce modifications of  $\mathbf{m}_i$  and  $\mathbf{n}_i$  by setting

$$\tilde{\mathbf{m}}_i(p, q) := \mathbf{m}_i(p, q)(\varphi(\Gamma_i))^{p+q}, \quad \tilde{\mathbf{n}}_i(p, q) := \mathbf{n}_i(p, q)(\varphi(\Gamma_i))^{p+q}.$$

In addition, for any  $\varepsilon' > 0$ , let  $\hat{\Omega}(z) = \hat{\Omega}(z, \varepsilon')$  be the event that all the concentration estimates of the components or quadratic forms of  $\mathbf{h}_i^u$  and  $\mathbf{h}_i^v$  in the proof of Lemma 5.3 hold with precision  $N^{\varepsilon'}$ . For instance, we used the large deviation bound (A.1) to bound  $(\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v$  in (5.41) by  $O_{\prec}(N^{-\frac{1}{2}})$ , in the proof of Lemma 5.3. Now we can bound it more quantitatively by  $\frac{N^{\varepsilon'}}{\sqrt{N}}$  on  $\hat{\Omega}(z)$ . Now we claim that

$$\mathbb{E}[\tilde{\mathbf{m}}_i(p, p)] = \mathbb{E}[\mathbf{c}_1 \tilde{\mathbf{m}}_i(p-1, p)] + \mathbb{E}[\mathbf{c}_2 \tilde{\mathbf{m}}_i(p-2, p)] + \mathbb{E}[\mathbf{c}_3 \tilde{\mathbf{m}}_i(p-1, p-1)] \quad (5.120)$$

with some random variables  $\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3$ , satisfying

$$|\mathbf{c}_1| \leq C \frac{N^{\varepsilon'}}{\sqrt{N\eta}}, \quad |\mathbf{c}_2| \leq C \frac{N^{2\varepsilon'}}{N\eta}, \quad |\mathbf{c}_3| \leq C \frac{N^{2\varepsilon'}}{N\eta}, \quad \text{on } \hat{\Omega}(z), \quad (5.121)$$

for some positive constant  $C$  which may depend on  $K$  in (5.118). In addition, the  $\mathbf{c}_i$ 's also admit trivial deterministic bounds of order  $\eta^{-k}$ , for some constant  $k > 0$ . Moreover, for any  $D' > 0$ , there exists  $N(D', \varepsilon')$ , such that if  $N \geq N(D', \varepsilon')$

$$\mathbb{P}(\hat{\Omega}(z)) \geq 1 - N^{-D'}.$$

Observe that (5.120) is just a more explicit version of (5.39), considering that  $\hat{\Omega}(z)$  holds with high probability. The proof of the more quantitative estimate (5.120) with (5.121) is basically the same as the proof of the non-quantitative one in (5.39).

The price for introducing  $\varphi(\Gamma_i)$  into  $\tilde{\mathbf{m}}_i$  is that it creates additional terms in the integration by parts. However, they are absorbed into the first term in the right side of (5.120). For instance, in the analogue of the step (5.56), except for replacing  $\mathbf{m}_i$  by  $\tilde{\mathbf{m}}_i$ , we will have an additional term

$$\frac{1}{N} \sum_k^{(i)} \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} (\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i) \frac{\partial \varphi(\Gamma_i)}{\partial g_{ik}^u} \tau_1(G) \tilde{\mathbf{m}}_i(p-1, p) \right].$$

For example, one term of  $\frac{\partial \varphi(\Gamma_i)}{\partial g_{ik}^u}$  is

$$\varphi'(\Gamma_i) \frac{\partial |G_{ii}|^2}{\partial g_{ik}^u} = \varphi'(\Gamma_i) \frac{\partial G_{ii}}{\partial g_{ik}^u} \overline{G_{ii}} + \varphi'(\Gamma_i) \frac{\partial \overline{G_{ii}}}{\partial g_{ik}^u} G_{ii}.$$

Using the second estimate in (5.75),

$$\frac{1}{N} \sum_k^{(i)} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \frac{\partial |G_{ii}|^2}{\partial g_{ik}^u} = O\left(\frac{N^{\varepsilon'}}{\sqrt{N\eta}}\right), \quad \text{on } \{\varphi'(\Gamma_i) \neq 0\} \cap \hat{\Omega}(z).$$

It is also easy to check that the other terms in  $\frac{\partial \varphi(\Gamma_i)}{\partial g_{ik}^u}$  give the same bound. Therefore, we have (5.120).

Using Young's inequality to (5.120), we can get

$$\begin{aligned}\mathbb{E}[\tilde{\mathbf{m}}_i(p, p)] &\leq C_p N^{2p\varepsilon'} \left( \mathbb{E}[|\mathbf{c}_1|^{2p}] + \mathbb{E}[|\mathbf{c}_2|^p] + \mathbb{E}[|\mathbf{c}_3|^p] \right) \\ &\leq C_p N^{2p\varepsilon'} \left( \left( \frac{N^{\varepsilon'}}{\sqrt{N\eta}} \right)^{2p} + N^{-D'} \eta^{-2kp} \right),\end{aligned}$$

which implies by Markov's inequality that

$$\mathbb{P}\left( |\mathcal{P}_{ii}\varphi(\Gamma_i)| \geq \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}} \right) \leq C_p \left( \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}} \right)^{-2p} N^{2p\varepsilon'} \left( \left( \frac{N^{\varepsilon'}}{\sqrt{N\eta}} \right)^{2p} + N^{-D'} \eta^{-2kp} \right). \quad (5.122)$$

For the given  $\varepsilon > 0$  in Lemma 5.5, by first choosing  $\varepsilon' = \varepsilon'(\varepsilon)$  to be smaller than  $\frac{\varepsilon}{8}$ , and then choosing  $p = p(\varepsilon, D)$  to be sufficiently large, we get

$$C_p \left( \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}} \right)^{-2p} N^{2p\varepsilon'} \left( \frac{N^{\varepsilon'}}{\sqrt{N\eta}} \right)^{2p} \leq \frac{1}{2} N^{-D}. \quad (5.123)$$

Then, by further choosing  $D' = D'(\varepsilon, D)$  sufficiently large, we can guarantee

$$C_p \left( \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}} \right)^{-2p} N^{2p\varepsilon'} N^{-D'} \eta^{-2kp} \leq \frac{1}{2} N^{-D}. \quad (5.124)$$

With these choices of  $\varepsilon'$  and  $D'$ , we now set  $N_2(D, \varepsilon) := N(D', \varepsilon')$ .

Further, by (5.122)-(5.124), there exists an event  $\Omega(z)$ , such that

$$\mathbb{P}(\Omega(z)) \geq 1 - N^{-D}, \quad N \geq N_2(D, \varepsilon)$$

and

$$|\mathcal{P}_{ii}\varphi(\Gamma_i)| \leq \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}}, \quad \text{on } \Omega(z).$$

This now implies that  $|\mathcal{P}_{ii}| \leq \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}}$  on  $\Theta(z, N^{-\frac{\gamma}{4}}) \cap \Omega(z)$ . Similarly, by working on  $\tilde{\mathbf{n}}_i$ , we can get  $|\mathcal{K}_{ii}| \leq \frac{N^{\frac{\varepsilon}{4}}}{\sqrt{N\eta}}$  on  $\Theta(z, N^{-\frac{\gamma}{4}}) \cap \Omega(z)$ .

The same bound can be obtained for  $\mathcal{P}_{ij}, \mathcal{P}_{\hat{i}j}, \mathcal{K}_{ij}$  and  $\mathcal{K}_{\hat{i}j}$  for  $j = i, \hat{i}$ . The remaining argument is the same as the proof of (5.30) in Theorem 5.2. The only change is, instead of the notation  $\prec$ , we use the deterministic  $\leq$ , but restricting onto the event  $\Theta(z, N^{-\frac{\gamma}{4}}) \cap \Omega(z)$ .

More specifically, the quantitative proof of (5.117) yields that

$$\Lambda_d^c(z) \leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}}, \quad \tilde{\Lambda}_d^c(z) \leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}}, \quad \Lambda_T(z) \leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}}, \quad \tilde{\Lambda}_T(z) \leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}} \quad (5.125)$$

hold on the event  $\Theta(z, N^{-\frac{\gamma}{4}}) \cap \Omega(z)$ , for all  $N \geq N_3(D, \gamma, \varepsilon)$ .

Therefore, by the definitions of  $\Lambda_d^c$  and  $\tilde{\Lambda}_d^c$ , we have

$$\begin{aligned}\left| G_{ii} - \frac{\omega_B^c}{|\xi_i|^2 - (\omega_B^c)^2} \right| &\leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}}, & \left| \mathcal{G}_{ii} - \frac{\omega_A^c}{|\sigma_i|^2 - (\omega_A^c)^2} \right| &\leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}}, \\ \left| G_{\hat{i}\hat{i}} - \frac{\omega_B^c}{|\xi_i|^2 - (\omega_B^c)^2} \right| &\leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}}, & \left| \mathcal{G}_{\hat{i}\hat{i}} - \frac{\omega_A^c}{|\sigma_i|^2 - (\omega_A^c)^2} \right| &\leq \frac{N^{\frac{\varepsilon}{2}}}{\sqrt{N\eta}},\end{aligned} \quad (5.126)$$

for all  $i \in \llbracket 1, N \rrbracket$ , on the event  $\Theta(z, N^{-\frac{\gamma}{4}}) \cap \Omega(z)$  for all  $N \geq N_3(D, \gamma, \varepsilon)$ . Averaging the above estimates over  $i$ , we obtain the system of equations

$$\begin{aligned}m_H(z) &= m_A(\omega_B^c(z)) + r_A(z), \\ m_H(z) &= m_B(\omega_A^c(z)) + r_B(z), \\ \omega_A^c(z) + \omega_B^c(z) &= z - \frac{1}{m_H(z)},\end{aligned} \quad (5.127)$$

where the error terms  $r_A(z)$  and  $r_B(z)$  satisfy  $|r_A(z)|, |r_B(z)| \leq \frac{CN^{\frac{5}{2}}}{\sqrt{N\eta}}$  on the event  $\Theta(z, N^{-\frac{7}{4}}) \cap \Omega(z)$  for all  $N \geq N_3(D, \gamma, \varepsilon)$ . Here the last equation in (5.127) follows from the definition (4.17) or (5.22). From the definition of  $\Theta(z, \delta)$  in (5.112), (4.17) or (5.22), and the equations in (2.5) with  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$ , it is not difficult to check that

$$|\omega_A^c - \omega_A| \leq CN^{-\frac{7}{4}}, \quad |\omega_B^c - \omega_B| \leq CN^{-\frac{7}{4}}$$

hold on  $\Theta(z, N^{-\frac{7}{4}})$ . In particular, with the help of (A.4), this guarantees that the imaginary parts of  $\omega_A^c$  and  $\omega_B^c$  are separated away from zero, hence so are  $m_A(\omega_B^c)$  and  $m_B(\omega_A^c)$ . This allows us to rewrite (5.127) as

$$\|\Phi_{\mu_A, \mu_B}(\omega_A^c, \omega_B^c, z)\| = \tilde{r}(z), \quad (5.128)$$

where  $\tilde{r}(z) = (\tilde{r}_A(z), \tilde{r}_B(z))'$  satisfy  $|\tilde{r}_A(z)|, |\tilde{r}_B(z)| \leq \frac{CN^{\frac{5}{2}}}{\sqrt{N\eta}}$  on the event  $\Theta(z, N^{-\frac{7}{4}}) \cap \Omega(z)$  for all  $N \geq N_3(D, \gamma, \varepsilon)$ . Applying the stability of the system  $\Phi_{\mu_A, \mu_B}(\omega_A, \omega_B, z) = 0$  (see Theorem 4.1 of [2]), we obtain

$$|\omega_A^c - \omega_A| \leq \frac{CN^{\frac{5}{2}}}{\sqrt{N\eta}}, \quad |\omega_B^c - \omega_B| \leq \frac{CN^{\frac{5}{2}}}{\sqrt{N\eta}}, \quad (5.129)$$

on the event  $\Theta(z, N^{-\frac{7}{4}}) \cap \Omega(z)$  for all  $N \geq N_3(D, \gamma, \varepsilon)$ . Substituting (5.129) into the definition of  $\Lambda_d^c$  and  $\tilde{\Lambda}_d^c$ , we see that the first two inequalities in (5.125) imply similar bounds for  $\Lambda_d$  and  $\tilde{\Lambda}_d$ . This completes the proof of Lemma 5.5.  $\square$

With Lemma 5.5, the remaining proof of Theorem 5.1 closely follows that for Theorem 2.5 in [3], so we will only sketch the argument. We start with the result with large  $\eta = \eta_M$  for some large but fixed positive constant  $\eta_M$ . More specifically, from Lemma 8.1, we see that

$$\Lambda_d(E + i\eta_M) \prec \frac{1}{\sqrt{N\eta_M^4}}, \quad \tilde{\Lambda}_d(E + i\eta_M) \prec \frac{1}{\sqrt{N\eta_M^4}}, \quad (5.130)$$

for any fixed  $E \in \mathbb{R}$ . The second estimate in (5.130) can be obtained from Lemma 8.1 since one can apply this lemma to  $\mathcal{H}$  as well. In addition, using the trivial bound  $\|G\| \leq \frac{1}{\eta}$  and inequality  $|\mathbf{x}^* G \mathbf{y}| \leq \|G\| \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$ , we also have

$$\Lambda_T(E + i\eta_M) \leq \frac{1}{\eta_M}, \quad \tilde{\Lambda}_T(E + i\eta_M) \leq \frac{1}{\eta_M}, \quad (5.131)$$

for any fixed  $E \in \mathcal{B}_{\mu_\alpha \boxplus \mu_\beta}$ . According to the definition of  $\Theta(z, \delta)$  in (5.112), (5.130) and (5.131), we see that for any fixed  $E \in \mathcal{B}_{\mu_\alpha \boxplus \mu_\beta}$  and  $D > 0$ ,

$$\mathbb{P}(\Theta(E + i\eta_M, N^{-\frac{3\gamma}{8}})) \geq 1 - N^{-D}, \quad (5.132)$$

holds for all  $N \geq N_0(D, \gamma)$  for some positive integer  $N_0(D, \gamma)$ .

Starting with (5.132), we conduct a standard continuity argument, whose setup is best suited to our problem in the form presented in [3]. Specifically, we do a bootstrap by reducing  $\eta$  in very small steps,  $N^{-5}$  (say), starting from  $\eta_M$  and successively control the probability of the ‘‘good’’ events  $\Theta$ . Recall the event  $\Omega(z)$  in Lemma 5.5. The main task is to show for any fixed  $E \in \mathcal{I}$  and any  $\eta \in [\eta_m, \eta_M]$ ,

$$\Theta(E + i\eta, N^{-\frac{3\gamma}{8}}) \cap \Omega(E + i(\eta - N^{-5})) \subset \Theta(E + i(\eta - N^{-5}), N^{-\frac{3\gamma}{8}}), \quad (5.133)$$

which is the analogue of (7.20) of [3]. To see this inclusion, one first uses the Lipschitz continuity of the Green function,  $\|G(z) - G(z')\| \leq N^2|z - z'|$ , and of the subordination functions, *c.f.*, (A.4), to obtain

$$\Theta(E + i\eta, N^{-\frac{3\gamma}{8}}) \subset \Theta(E + i(\eta - N^{-5}), N^{-\frac{7}{4}}). \quad (5.134)$$



Then (5.134) together with (5.115) implies (5.133). Using (5.133) recursively, one goes from  $\eta_M$  down to  $\eta_m$ , step by step. The remaining proof of (5.25), based on (5.133) and Lemma 5.5, is the same as the counterpart in [3] (c.f., (7.20)-(7.25) therein). We omit the details.

With (5.25), we can prove (5.26) in the sequel. The first two inequalities in (5.26) have already been proved in (5.129) with a fixed  $\eta$ , under (5.116). The uniformity then follows from (5.25) which holds uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . Then the last inequality in (5.26) follows from the first two, together with the last equation in (5.127) and the second equation in (2.5) with  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$ . This completes the proof of Theorem 5.1.  $\square$

## 6. STRONG LAW FOR SMALL $\eta$

In this section, we prove the strong law, *i.e.*, Theorem 4.3, for  $z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)$ . It suffices to work on the regime  $z \in \mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$  at first. The extension to  $z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)$  will be easy. Our main task is to establish the fluctuation averaging for the quantities  $\mathcal{P}_{ij}$  defined in (5.19).

**Lemma 6.1** (Fluctuation averaging). *Suppose that the assumptions in Theorem 4.3 hold. Let  $\eta_M > 0$  be any (large) constant and  $\gamma > 0$  be any (small) constant in the definition of  $\eta_m$  (c.f., (5.1)). For any fixed integer  $p \geq 1$ , and deterministic numbers  $d_1, \dots, d_N \in \mathbb{C}$  satisfying  $\max_{i \in \llbracket 1, N \rrbracket} |d_i| \leq 1$ , we have*

$$\begin{aligned} \left| \frac{1}{N} \sum_{i=1}^N d_i \mathcal{P}_{ii} \right| &\prec \Psi^2, & \left| \frac{1}{N} \sum_{i=1}^N d_i \mathcal{P}_{\hat{i}\hat{i}} \right| &\prec \Psi^2, \\ \left| \frac{1}{N} \sum_{i=1}^N d_i \mathcal{P}_{\hat{i}\hat{i}} \right| &\prec \Psi^2, & \left| \frac{1}{N} \sum_{i=1}^N d_i \mathcal{P}_{\hat{i}\hat{i}} \right|^2 &\prec \Psi^2 \end{aligned} \quad (6.1)$$

uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ .

We will often use the following improvement of (5.34),

$$\begin{aligned} \tau_a(G) &= m_{\mu_A \boxplus \mu_B} + O_{\prec}(\Psi), \\ \tau_a(\tilde{B}G) &= (z - \omega_B) m_{\mu_A \boxplus \mu_B} + O_{\prec}(\Psi), \\ \tau_a(G\tilde{B}) &= (z - \omega_B) m_{\mu_A \boxplus \mu_B} + O_{\prec}(\Psi), \\ \tau_a(\tilde{B}G\tilde{B}) &= (\omega_B - z)(1 + (\omega_B - z) m_{\mu_A \boxplus \mu_B}) + O_{\prec}(\Psi), \quad a = 1, 2, \end{aligned} \quad (6.2)$$

which can be proved in the same way as (5.34), but with the first inequality in (5.27) replaced by the first inequality in (5.25), as the input of the proof.

In the next Section 6.1 we will show how to prove Theorem 4.3 on  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$  with the aid of Lemma 6.1. Then, in Section 6.2 we will prove Lemma 6.1.

**6.1. Proof of Theorem 4.3 on  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$ .** To prove the strong law from Lemma 6.1, first of all, we need to derive that the estimates

$$|\Upsilon_1| \prec \Psi^2, \quad |\Upsilon_2| \prec \Psi^2 \quad (6.3)$$

hold uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . These are the strongest high probability bounds related to the Ward identities in (5.29). To see (6.3), we choose  $d_i = 1$  for all  $i \in \llbracket 1, N \rrbracket$  in (6.1). From the definition of  $\mathcal{P}_{ii}$  in (5.19), we get

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \mathcal{P}_{ii} &= \frac{\Upsilon_1}{N} \sum_{i=1}^N (G_{ii} + T_{ii}) = \Upsilon_1 \left( \tau_1(G) + \frac{1}{N} \sum_{i=1}^N T_{ii} \right) \\ &= \Upsilon_1 \left( m_{\mu_A \boxplus \mu_B} + O_{\prec}(\Psi) \right), \end{aligned} \quad (6.4)$$

where in the last step we used (6.2) and the third inequality in (5.25). Then, using the lower bound of  $\text{Im } m_{\mu_A \boxplus \mu_B} = \text{Im } m_{\mu_A}(\omega_B)$  inherited from the lower bound of  $\text{Im } \omega_B$  in (A.4), and also the first bound in (6.1), we can easily see  $|\Upsilon_1| \prec \Psi^2$  from (6.4). Similarly, we can also

show  $|\Upsilon_2| \prec \Psi^2$ . Notice that *a posteriori* we could have defined  $\mathcal{P}_{ij}$  in (5.19) without the last term involving  $\Upsilon_a$  with  $a = 1, 2$ , since we are interested only up to  $O_{\prec}(\Psi^2)$  precision. We do not, however, know how to prove directly that  $\Upsilon_a = O_{\prec}(\Psi^2)$  without first proving a fluctuation averaging result (6.1) involving the quantity  $\mathcal{P}_{ij}$  with  $\Upsilon_a$ . The correct choice of  $\mathcal{P}_{ij}$  is the essential idea of the entire proof.

Plugging (6.3) back to the definition of  $\mathcal{P}_{ii}$ ,  $\mathcal{P}_{i\hat{i}}$ ,  $\mathcal{P}_{\hat{i}i}$  and  $\mathcal{P}_{\hat{i}\hat{i}}$  in (5.19), we obtain from (6.1),

$$\begin{aligned} & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ij} \tau_1(\tilde{B}G) - (\tilde{B}G)_{ij} \tau_1(G) \right) \right| \prec \Psi^2, \\ & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{i\hat{j}} \tau_2(\tilde{B}G) - (\tilde{B}G)_{i\hat{j}} \tau_2(G) \right) \right| \prec \Psi^2, \quad j = i, \hat{i}, \end{aligned} \quad (6.5)$$

for any deterministic numbers  $d_1, \dots, d_N \in \mathbb{C}$  satisfying  $|d_i| \lesssim 1$ , which is a shorthand notation for  $|d_i| \leq C$  with some constant  $C$ . While Lemma 6.1 was formulated for  $|d_i| \leq 1$ , it clearly holds as long as  $|d_i| \lesssim 1$ . Recall the notation introduced in (8.20). We claim that the following estimates can be derived from (5.94) and (6.5):

$$\begin{aligned} & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} - \frac{\omega_{B,2}^c}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \right) \right| \prec \Psi^2, \\ & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{i\hat{i}} - \frac{\bar{\xi}_i}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \right) \right| \prec \Psi^2, \\ & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{\hat{i}i} - \frac{\omega_{B,1}^c}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \right) \right| \prec \Psi^2, \\ & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{\hat{i}\hat{i}} - \frac{\xi_i}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \right) \right| \prec \Psi^2. \end{aligned} \quad (6.6)$$

We derive the first estimate in (6.6), the others are proven similarly. We write

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} - \frac{\omega_{B,2}^c}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \right) \\ &= \frac{1}{N} \sum_{i=1}^N \frac{d_i}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} \left( G_{ii} (|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c) - \omega_{B,2}^c \right). \end{aligned}$$

Applying Theorem 5.1, and (5.106) along its proof, it is easy to check that

$$\omega_{B,a}^c = \omega_B + O_{\prec}(\Psi), \quad a = 1, 2, \quad (6.7)$$

hence

$$\omega_{B,1}^c \omega_{B,2}^c = \omega_B^2 + O_{\prec}(\Psi), \quad G_{ii} (|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c) - \omega_{B,2}^c = O_{\prec}(\Psi). \quad (6.8)$$

Moreover, from the lower bound on  $\text{Im} \omega_B$  from (A.4) and the first estimate of (6.8), we have

$$\frac{1}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c} = \frac{1}{|\xi_i|^2 - \omega_B^2} + O_{\prec}(\Psi). \quad (6.9)$$

Then, in light of (A.4), (6.8) and (6.9), it suffices to check

$$\frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} (|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c) - \omega_{B,2}^c \right) = O_{\prec}(\Psi^2), \quad (6.10)$$

for any deterministic numbers  $d_1, \dots, d_N \in \mathbb{C}$  satisfying  $|d_i| \lesssim 1$  (here we redefined  $d_i$  to  $d_i/(|\xi_i|^2 - \omega_B^2)$ ). Using (5.94), we can write

$$\begin{aligned} G_{ii}(|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c) - \omega_{B,2}^c &= -\frac{\omega_{B,2}^c}{\tau_1(G)} \left( (\tilde{B}G)_{ii} \tau_1(G) - G_{ii} \tau_1(\tilde{B}G) \right) \\ &\quad - \frac{\xi_i}{\tau_2(G)} \left( (\tilde{B}G)_{ii} \tau_2(G) - G_{ii} \tau_2(\tilde{B}G) \right). \end{aligned} \quad (6.11)$$

Then, from (6.2) and (6.7), we see that

$$\begin{aligned} \frac{\omega_{B,2}^c}{\tau_1(G)} &= \frac{\omega_B}{m_{\mu_A \boxplus \mu_B}} + O_{\prec}(\Psi), & \frac{\xi_i}{\tau_2(G)} &= \frac{\xi_i}{m_{\mu_A \boxplus \mu_B}} + O_{\prec}(\Psi), \\ (\tilde{B}G)_{ii} \tau_1(G) - G_{ii} \tau_1(\tilde{B}G) &= O_{\prec}(\Psi), & (\tilde{B}G)_{ii} \tau_2(G) - G_{ii} \tau_2(\tilde{B}G) &= O_{\prec}(\Psi), \end{aligned} \quad (6.12)$$

where the second line follows from (5.102). Thus combining (6.11), (6.12) and (6.5) yields (6.10), which implies (6.6) according to the discussion above.

Notice that in this argument it was essential that  $G_{ii}$  was approximated in (6.6) not by  $\omega_B/(|\xi_i|^2 - \omega_B^2)$  or by  $\omega_{B,1}^c/(|\xi_i|^2 - (\omega_B^c)^2)$  but by

$$G_{ii} \approx \frac{\omega_{B,2}^c}{|\xi_i|^2 - \omega_{B,1}^c \omega_{B,2}^c},$$

since this latter approximation is precise up to  $O_{\prec}(\Psi^2)$  after averaging, while the previous ones are *a priori* correct only with an error  $O_{\prec}(\Psi)$ .

Next, we show that (6.6) nevertheless holds if we approximate  $G_{ii}$  by  $\omega_{B,1}^c/(|\xi_i|^2 - (\omega_B^c)^2)$ . Choosing all  $d_i = 1$  in the first and third inequalities in (6.6) and applying (4.27), we note that

$$\omega_{B,a}^c = \omega_B^c + O_{\prec}(\Psi^2), \quad a = 1, 2,$$

so the first approximation in (6.7) is actually one order better. Thus we get from (6.6) that

$$\begin{aligned} \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} - \frac{\omega_B^c}{|\xi_i|^2 - (\omega_B^c)^2} \right) \right| &\prec \Psi^2, & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} - \frac{\bar{\xi}_i}{|\xi_i|^2 - (\omega_B^c)^2} \right) \right| &\prec \Psi^2, \\ \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} - \frac{\omega_B^c}{|\xi_i|^2 - (\omega_B^c)^2} \right) \right| &\prec \Psi^2, & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( G_{ii} - \frac{\xi_i}{|\xi_i|^2 - (\omega_B^c)^2} \right) \right| &\prec \Psi^2. \end{aligned} \quad (6.13)$$

Further, recalling the definitions of  $\mathcal{H}$  and  $\mathcal{G}$  in (4.23). Switching the rôles of  $A$  and  $B$ , and also the rôles of  $U$  and  $U^*$  in the above discussions, we have

$$\begin{aligned} \left| \frac{1}{N} \sum_{i=1}^N d_i \left( \mathcal{G}_{ii} - \frac{\omega_A^c}{|\sigma_i|^2 - (\omega_A^c)^2} \right) \right| &\prec \Psi^2, & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( \mathcal{G}_{ii} - \frac{\bar{\sigma}_i}{|\sigma_i|^2 - (\omega_A^c)^2} \right) \right| &\prec \Psi^2, \\ \left| \frac{1}{N} \sum_{i=1}^N d_i \left( \mathcal{G}_{ii} - \frac{\omega_A^c}{|\sigma_i|^2 - (\omega_A^c)^2} \right) \right| &\prec \Psi^2, & \left| \frac{1}{N} \sum_{i=1}^N d_i \left( \mathcal{G}_{ii} - \frac{\sigma_i}{|\sigma_i|^2 - (\omega_A^c)^2} \right) \right| &\prec \Psi^2. \end{aligned} \quad (6.14)$$

Applying (6.13) and (6.14) to average over the diagonal entries of the Green functions  $G$  and  $\mathcal{G}$ , and also using the fact  $\text{tr } G(z) = \text{tr } \mathcal{G}(z) = m_H(z)$ , we see that

$$m_H(z) = \int_{\mathbb{R}} \frac{\omega_B^c}{x^2 - (\omega_B^c)^2} d\mu_{\Xi}(x) + O_{\prec}(\Psi^2) = \int_{\mathbb{R}} \frac{\omega_A^c}{x^2 - (\omega_A^c)^2} d\mu_{\Sigma}(x) + O_{\prec}(\Psi^2).$$

From this, using

$$\frac{\omega_B^c}{x^2 - (\omega_B^c)^2} = \frac{1}{2} \left[ \frac{1}{x - \omega_B^c} + \frac{1}{-x - \omega_B^c} \right],$$

we can get

$$m_H(z) = m_A(\omega_B^c(z)) + O_{\prec}(\Psi^2) = m_B(\omega_A^c(z)) + O_{\prec}(\Psi^2), \quad (6.15)$$

where we used the fact  $\mu_A \equiv \mu_{\Xi}^{\text{sym}}$  and  $\mu_B \equiv \mu_{\Sigma}^{\text{sym}}$ , in light of (4.14). In addition, we also have (4.18). Summarizing these estimates, we have  $\Phi_{\mu_A, \mu_B}(\omega_A^c, \omega_B^c, z) = O_{\prec}(\Psi^2)$ , *i.e.*, compared with (5.128), we improved the error in the approximate subordination equations.

Similarly to the proof of Lemma 5.5, we use the stability of the system  $\Phi_{\mu_A, \mu_B}(\omega_A, \omega_B, z) = 0$  again, but with the improved error  $\Psi^2$ . We also note that the estimates from Theorem 5.1 and Lemma 6.1 used in the above discussion hold uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . Hence, we can conclude the proof of Theorem 4.3 on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ .

At the end, we extend (4.11) from  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$  to  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$ . The extension relies on a standard use of the monotonicity of the Green function: For all  $i \in \llbracket 1, N \rrbracket$  and  $j = i$  or  $\hat{i}$ , we have

$$|G'_{jj}(z)| = \left| \sum_{k=1}^{2N} G_{jk}(z)G_{kj}(z) \right| \leq \sum_{k=1}^{2N} |G_{jk}(z)|^2 = \frac{\text{Im } G_{jj}(z)}{\eta},$$

where the last step follows from the spectral decomposition. In addition, note that the function  $s \mapsto s \text{Im } G_{jj}(E + is)$  is monotonically increasing. This implies that for any  $\eta \in (0, \eta_m)$ ,

$$\begin{aligned} |G_{jj}(E + i\eta) - G_{jj}(E + i\eta_m)| &\leq \int_{\eta}^{\eta_m} \frac{s \text{Im } G_{jj}(E + is)}{s^2} ds \\ &\leq 2 \frac{\eta_m}{\eta} \text{Im } G_{jj}(E + i\eta_m) \leq C \frac{N^\gamma}{N\eta} \leq CN^\gamma \Psi^2, \end{aligned} \quad (6.16)$$

with high probability, for any  $E \in \mathcal{I}$ . Here we used  $|G_{jj}(E + i\eta_m)| \prec 1$  which follows from the first bound in (5.25). On the other hand, for any  $i \in \llbracket 1, N \rrbracket$ , we also have

$$\left| \frac{\omega_B(E + i\eta)}{|\xi_i|^2 - \omega_B^2(E + i\eta)} - \frac{\omega_B(E + i\eta_m)}{|\xi_i|^2 - \omega_B^2(E + i\eta_m)} \right| \leq C(\eta_m - \eta) \leq \Psi^2, \quad (6.17)$$

$\eta \in (0, \eta_m]$ ,  $E \in \mathcal{I}$ , for sufficiently small  $\gamma$ , which follows from the upper bound of  $\omega'_B(z)$ , the lower bound of  $|\xi_i|^2 - \omega_B^2(z)$  which follows from the lower bound of  $\text{Im } \omega_B$ , and also the upper bound of  $\omega_B$ , in Lemma A.2. Combining (6.16) and (6.17), and using (5.33), we conclude that (4.11) holds uniformly on  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$ . This completes the proof of Theorem 4.3 on  $\mathcal{S}_{\mathcal{I}}(0, \eta_M)$ .

Hence, what remains is to prove Lemma 6.1.

**6.2. Proof of Lemma 6.1.** Since the proofs for the four estimates in (6.1) are nearly the same, we only present the details for the first one. First of all, from (5.25) and (5.29) we have

$$|T_{ii}\Upsilon_1| \prec \Psi^2. \quad (6.18)$$

Hence, it suffices to bound the weighted average of the following slight modifications of  $\mathcal{P}_{ii}$ 's:

$$\mathcal{Q}_{ii} \equiv \mathcal{Q}_{ii}(z) := (\tilde{B}G)_{ii}\tau_1(G) - G_{ii}\tau_1(\tilde{B}G) + G_{ii}\Upsilon_1, \quad i \in \llbracket 1, N \rrbracket. \quad (6.19)$$

Then we introduce the notation

$$\mathbf{m}(k, l) := \left( \frac{1}{N} \sum_{i=1}^N d_i \mathcal{Q}_{ii} \right)^k \left( \frac{1}{N} \sum_{i=1}^N \overline{d_i \mathcal{Q}_{ii}} \right)^l.$$

Similarly to Lemma 5.3, the main technical task is the following recursive moment estimate.

**Theorem 6.2** (Recursive moment estimate). *Suppose that the assumptions in Theorem 4.3 hold. Let  $\eta_M > 0$  be any (large) constant and  $\gamma > 0$  in (5.1) be any (small) constant. For any fixed integer  $p \geq 1$ , we have*

$$\begin{aligned} \mathbb{E}[\mathbf{m}(p, p)] &= \mathbb{E}[O_{\prec}(\Psi^2)\mathbf{m}(p-1, p)] + \mathbb{E}[O_{\prec}(\Psi^4)\mathbf{m}(p-2, p)] \\ &\quad + \mathbb{E}[O_{\prec}(\Psi^4)\mathbf{m}(p-1, p-1)], \end{aligned} \quad (6.20)$$

uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ , where we made the convention  $\mathbf{m}(0, 0) = 1$  and  $\mathbf{m}(-1, 1) = 0$  if  $p = 1$ .

The reason why we prefer to work with  $\mathcal{Q}_{ii}$  instead of  $\mathcal{P}_{ii} = \mathcal{Q}_{ii} + T_{ii}\Upsilon_i$  is as follows. To prove Theorem 6.2, we will follow a similar strategy as the proof of Lemma 5.3. In Lemma 5.3 and its proof, we worked on  $\mathcal{P}_{ii}$  directly. The derivative  $\frac{\partial T_{ii}}{\partial g_{ik}^u}$  was necessary for the proof of Lemma 5.3, *c.f.*, (5.75). However, in the proof of Theorem 6.2, we would need to consider the derivative  $\frac{\partial T_{ii}}{\partial g_{jk}^u}$  for all  $j \neq i$  if we carry the term  $T_{ii}\Upsilon_1$  from  $\mathcal{P}_{ii}$  in the discussion. Unfortunately, the dependence of the factor  $(\mathbf{k}_i^u)^*$  in  $T_{ii}$  (*c.f.*, (5.16)) on  $g_{jk}^u$  for  $j \neq i$  is difficult to capture. On the other hand, at this stage of the proof we already have the bound (6.18) available and this allows us to drop the term  $T_{ii}\Upsilon_1$  from the beginning.

With the aid of Theorem 6.2, one can prove Lemma 6.1.

*Proof of Lemma 6.1.* Similarly to the proof of (5.28) for  $\mathcal{P}_{ii}$  from Lemma 5.3, one can apply Young's inequality to (6.20) and get  $|\frac{1}{N} \sum_{i=1}^N \mathcal{Q}_{ii}| \prec \Psi^2$ , which together with (6.18) implies the first bound in (6.1). The other three in (6.1) can be verified analogously. Hence, we completed the proof of Lemma 6.1.  $\square$

*Proof of Theorem 6.2.* Hence, we start with the averaged analogue of (5.71), but with  $\mathcal{P}_{ii}$ 's replaced by  $\mathcal{Q}_{ii}$ 's. In particular, the term  $T_{ii}$  is missing. Following the proof of (5.71) with these modifications, we obtain

$$\begin{aligned}
\mathbb{E}[\mathbf{m}(p, p)] &= \frac{1}{N} \sum_i d_i \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \left( \dot{T}_{ii} - \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \right) \tau_1(\tilde{B}G) \mathbf{m}(p-1, p) \right] \\
&\quad - \frac{1}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \mathbf{m}(p-1, p) \right] \\
&\quad - \frac{1}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{\partial \tau_1(G)}{\partial g_{ik}^u} \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \mathbf{m}(p-1, p) \right] \\
&\quad - \frac{p-1}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \left( \frac{1}{N} \sum_j d_j \frac{\partial \mathcal{Q}_{jj}}{\partial g_{ik}^u} \right) \mathbf{m}(p-2, p) \right] \\
&\quad - \frac{p}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i \tau_1(G) \left( \frac{1}{N} \sum_j \bar{d}_j \frac{\partial \overline{\mathcal{Q}}_{jj}}{\partial g_{ik}^u} \right) \mathbf{m}(p-1, p-1) \right] \\
&\quad + \frac{1}{N} \sum_i d_i \mathbb{E} \left[ \left( \varepsilon_{i1} \tau_1(G) - \frac{\varepsilon_{i4} + \varepsilon_{i5}}{\|\mathbf{g}_i^u\|_2} \right) \mathbf{m}(p-1, p) \right] + \mathbb{E}[O_{\prec}(\Psi^2) \mathbf{m}(p-1, p)]. \tag{6.21}
\end{aligned}$$

In addition, we also have the averaged analogue of (5.72):

$$\begin{aligned}
&\frac{1}{N} \sum_i d_i \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2} \left( \dot{T}_{ii} - \frac{1}{N} \sum_k^{(i)} \frac{\partial(\hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} \right) \tau_1(\tilde{B}G) \mathbf{m}(p-1, p) \right] \\
&= \frac{1}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{\partial \|\mathbf{g}_i^u\|_2^{-2}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \tau_1(\tilde{B}G) \mathbf{m}(p-1, p) \right] \\
&\quad + \frac{1}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{\partial \tau_1(\tilde{B}G)}{\partial g_{ik}^u} \frac{1}{\|\mathbf{g}_i^u\|_2^2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \mathbf{m}(p-1, p) \right] \\
&\quad + \frac{p-1}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2^2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \tau_1(\tilde{B}G) \left( \frac{1}{N} \sum_j d_j \frac{\partial \mathcal{Q}_{jj}}{\partial g_{ik}^u} \right) \mathbf{m}(p-2, p) \right]
\end{aligned}$$

$$+ \frac{p}{N^2} \sum_i \sum_k^{(i)} d_i \mathbb{E} \left[ \frac{1}{\|\mathbf{g}_i^u\|_2^2} \hat{\mathbf{e}}_k^* G \hat{\mathbf{e}}_i \tau_1(\tilde{B}G) \left( \frac{1}{N} \sum_j \tilde{d}_j \frac{\partial \overline{Q_{jj}}}{\partial g_{ik}^u} \right) \mathbf{m}(p-1, p-1) \right]. \quad (6.22)$$

Hence, to show (6.20), it suffices to estimate the second to the fifth terms on the right side of (6.21), and the terms on the right side of (6.22). First, we notice that

$$\varepsilon_{i4} = O_{\prec}(\Psi^2), \quad (6.23)$$

which can be seen from (5.25), (5.29), and the facts  $|\|\mathbf{g}_i^u\|_2^2 - 1| \prec \frac{1}{\sqrt{N}}$  and  $|h_{ii}^u| \prec \frac{1}{\sqrt{N}}$ . All the other desired estimates can be derived from the following lemma.

**Lemma 6.3.** *Suppose that the assumptions in Theorem 4.3 hold. Let  $\eta_M > 0$  be any (large) constant and  $\gamma > 0$  in (5.1) be any (small) constant. Let  $\hat{d}_1, \dots, \hat{d}_N \in \mathbb{C}$  be deterministic numbers with the bound  $\max_i |\hat{d}_i| \lesssim 1$  and let  $\tilde{d}_1, \dots, \tilde{d}_N \in \mathbb{C}$  be (possibly random) numbers with the bound  $\max_i |\tilde{d}_i| \prec 1$  for all  $i \in \llbracket 1, N \rrbracket$ . Let  $Q$  be any deterministic diagonal matrix satisfying  $\|Q\| \leq C$  and  $X = \hat{I}$  or  $A$ , set  $X_i = \hat{I}$  or  $\tilde{B}^{(i)}$ , and let*

$$\mathbf{x}_i, \mathbf{y}_i = \begin{pmatrix} \hat{\mathbf{g}}_i^u \\ \mathbf{0} \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} \mathbf{0} \\ \hat{\mathbf{g}}_i^v \end{pmatrix}.$$

We have the estimates

$$\begin{aligned} \frac{1}{N^2} \sum_{i=1}^N \sum_k^{(i)} \tilde{d}_i \frac{\partial \|\mathbf{g}_i^u\|_2^{-1}}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= O_{\prec} \left( \frac{1}{N} \right), \\ \frac{1}{N^2} \sum_{i=1}^N \sum_k^{(i)} \tilde{d}_i \frac{\partial \text{tr} Q X G}{\partial g_{ik}^u} \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i &= O_{\prec}(\Psi^4), \end{aligned} \quad (6.24)$$

uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ . In addition, we also have

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \hat{d}_i \mathbb{E} \left[ \left( \mathbf{x}_i^* X_i \mathbf{y}_i - \mathbb{E}_i \left[ \mathbf{x}_i^* X_i \mathbf{y}_i \right] \right) \mathbf{m}(p-1, p) \right] \\ = \mathbb{E} \left[ O_{\prec}(\Psi^2) \mathbf{m}(p-1, p) \right] + \mathbb{E} \left[ O_{\prec}(\Psi^4) \mathbf{m}(p-2, p) \right] \\ + \mathbb{E} \left[ O_{\prec}(\Psi^4) \mathbf{m}(p-1, p-1) \right], \end{aligned} \quad (6.25)$$

uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_m, \eta_M)$ , where  $\mathbb{E}_i$  denotes the expectation with respect to  $\hat{\mathbf{g}}_i^u$  and  $\hat{\mathbf{g}}_i^v$ .

With Lemma 6.3, we can proceed to the proof of Theorem 6.2 as follows. First of all, for any diagonal matrix  $Q = \text{diag}(q_1, \dots, q_{2N})$ , using the first estimate in (5.25), we have

$$\begin{aligned} \text{tr} QG &= \frac{1}{2N} \sum_{i=1}^N (q_i + q_i) \frac{\omega_B(z)}{|\xi_i|^2 - (\omega_B(z))^2} + O_{\prec}(\Psi), \\ \text{tr} QAG &= \frac{1}{2N} \sum_{i=1}^N (q_i + q_i) \frac{|\xi_i|^2}{|\xi_i|^2 - (\omega_B(z))^2} + O_{\prec}(\Psi). \end{aligned}$$

Using the upper bound of  $\omega_B$  and the lower bound of  $\text{Im} \omega_B$  in (A.4), we can see that

$$|\text{tr} QXG| \prec 1, \quad (6.26)$$

for diagonal  $Q$  with  $\|Q\| \leq C$  and  $X = \hat{I}$  or  $A$ . Note that all partial traces such as  $\tau_1(G)$ ,  $\tau_1(\tilde{B}G)$  can be written as a linear combination of terms of the form  $\text{tr} QXG$  with the aid of the identities in (4.16), and thus for these partial traces we have

$$\tau_1(G) = O_{\prec}(1), \quad \tau_1(\tilde{B}G) = O_{\prec}(1).$$

These bounds together with the first estimate in (6.24), imply the desired estimates for the second term on the right side of (6.21) and the first term on the right side of (6.22).

Next, notice that

$$\frac{1}{N} \sum_{j=1}^N d_j \mathcal{Q}_{jj} = \text{tr}(D\tilde{B}G)\tau_1(G) - \text{tr}(DG)\tau_1(\tilde{B}G) + \text{tr}(DG)\Upsilon_1,$$

where we denoted the deterministic diagonal matrix  $D := \text{diag}(d_1, \dots, d_N) \oplus 0$ , with 0 the  $N \times N$  zero matrix. In addition, using (4.16), we can see that  $\frac{1}{N} \sum_j d_j \mathcal{Q}_{jj}$  is a polynomial of  $\frac{1}{N} \sum_j d_j T_{jj}$  and the terms of the form  $\text{tr} QXG$  for some diagonal  $Q$  with  $\|Q\| \leq C$  and  $X = \hat{I}$  or  $A$ . Here we also used the fact that  $\tau_a(\mathcal{D}) = \text{tr}(\hat{I}_a \mathcal{D})$  for any  $\mathcal{D} \in M_{2N}(\mathbb{C})$  and  $a = 1, 2$ , where  $\hat{I}_a$  is defined in (4.25). Then the last two estimates in (6.24), (6.26), together with the chain rule, imply that

$$\frac{1}{N^3} \sum_{i=1}^N \sum_k^{(i)} \tilde{d}_i \hat{e}_k^* X_i G \hat{e}_i \sum_{j=1}^N d_j \frac{\partial \mathcal{Q}_{jj}}{\partial g_{ik}^u} = O_{\prec}(\Psi^4). \quad (6.27)$$

Similarly, we can prove the same bound if we replace  $\mathcal{Q}_{jj}$ 's by  $\bar{\mathcal{Q}}_{jj}$ 's. Hence, the desired estimates for the third to the fifth terms on the right side of (6.21), and the last three terms on the right side of (6.22) can be obtained from the second estimate in (6.24).

Hence, what remains is to estimate the sixth term in (6.21). First, according to (6.23), we can neglect  $\varepsilon_{i4}$ . Then we recall the definition of  $\varepsilon_{i1}$  from (5.43). Using the estimates of  $G_{ii}$  and  $T_{ii}$  from the first and the third inequalities in (5.25), and the estimates

$$\ell_i^v = 1 + O_{\prec}\left(\frac{1}{\sqrt{N}}\right), \quad |h_{ii}^u| \prec \frac{1}{\sqrt{N}}, \quad |(\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v| \prec \frac{1}{\sqrt{N}},$$

we see that

$$\varepsilon_{i1} = \frac{\bar{\xi}_i}{|\bar{\xi}_i|^2 - \omega_B^2} (\mathbf{k}_i^u)^* \tilde{B}^{(i)} \mathbf{k}_i^v + O_{\prec}(\Psi^2) = \frac{\bar{\xi}_i}{|\bar{\xi}_i|^2 - \omega_B^2} (\boldsymbol{\ell}_i^u)^* \tilde{B}^{(i)} \boldsymbol{\ell}_i^v + O_{\prec}(\Psi^2), \quad (6.28)$$

where we introduced the notations

$$\boldsymbol{\ell}_i^u := \begin{pmatrix} \hat{g}_i^u \\ \mathbf{0} \end{pmatrix}, \quad \boldsymbol{\ell}_i^v := \begin{pmatrix} \mathbf{0} \\ \hat{g}_i^v \end{pmatrix}.$$

Then, recall the definition of  $\varepsilon_{i5}$  from (5.69). Applying the estimate of  $G_{ii}$  from the first inequality in (5.25), and the second formula in (6.2), and the fact  $\|\mathbf{g}_i^u\|_2^2 = \|\boldsymbol{\ell}_i^u\|_2^2 + O_{\prec}(\frac{1}{N}) = 1 + O_{\prec}(\frac{1}{\sqrt{N}})$ , we also have

$$\varepsilon_{i5} = \frac{(z - \omega_B) m_{\mu_A} \boxplus \mu_B \omega_B}{|\bar{\xi}_i|^2 - \omega_B^2} (\|\boldsymbol{\ell}_i^u\|_2^2 - 1) + O_{\prec}(\Psi^2). \quad (6.29)$$

Note that both of the first terms on the right side of (6.28) and (6.29) are of the form  $\hat{d}_i(\mathbf{x}_i^* X_i \mathbf{y}_i - \mathbb{E}_i[\mathbf{x}_i^* X_i \mathbf{y}_i])$  for some deterministic  $\hat{d}_i$  with  $|\hat{d}_i| \lesssim 1$ . Hence, using (6.25), we get the desired bound for the sixth term of (6.21). This completes the proof of Theorem 6.2 up to the proof of Lemma 6.3.  $\square$

*Proof of Lemma 6.3.* The first estimate in (6.24) follows directly from the first estimate in (5.75). The second estimate of (6.24) is a weighted average of the last estimate in (5.75).

Hence, what remains is to prove (6.25). We only show the details for the case  $\mathbf{x}_i = \boldsymbol{\ell}_i^u$  and  $\mathbf{y}_i = \boldsymbol{\ell}_i^v$ . The others are similar. Notice that in this case,  $\mathbb{E}_i[\mathbf{x}_i^* X_i \mathbf{y}_i] = 0$ . Using the integration by parts formula (5.44) again, we have

$$\frac{1}{N} \sum_i \hat{d}_i \mathbb{E} \left[ (\boldsymbol{\ell}_i^u)^* X_i \boldsymbol{\ell}_i^v \mathfrak{m}(p-1, p) \right] = \frac{1}{N} \sum_i \sum_k^{(i)} \hat{d}_i \mathbb{E} \left[ \bar{g}_{ik}^u \hat{e}_k^* X_i \boldsymbol{\ell}_i^v \mathfrak{m}(p-1, p) \right]$$

$$\begin{aligned}
&= \frac{p-1}{N^2} \sum_i \sum_k^{(i)} \hat{d}_i \mathbb{E} \left[ \hat{e}_k^* X_i \ell_i^v \frac{1}{N} \sum_j d_j \frac{\partial \mathcal{Q}_{jj}}{\partial g_{ik}^u} \mathbf{m}(p-2, p) \right] \\
&\quad + \frac{p}{N^2} \sum_i \sum_k^{(i)} \hat{d}_i \mathbb{E} \left[ \hat{e}_k^* X_i \ell_i^v \frac{1}{N} \sum_j \bar{d}_j \frac{\partial \bar{\mathcal{Q}}_{jj}}{\partial g_{ik}^u} \mathbf{m}(p-1, p-1) \right].
\end{aligned}$$

Hence, it suffices to show

$$\left| \frac{1}{N^3} \sum_i \sum_k^{(i)} \hat{d}_i \hat{e}_k^* X_i \ell_i^v \sum_j d_j \frac{\partial \mathcal{Q}_{jj}}{\partial g_{ik}^u} \right| \prec \Psi^4 \quad (6.30)$$

and its complex conjugate analogue. The proof of (6.30) is nearly the same as that of (6.27). Hence, we omit it. Therefore, we completed the proof of Lemma 6.3.  $\square$

## 7. PROOF OF THEOREM 2.2

Theorem 2.2 will directly follow from a more detailed result, Theorem 7.1, below. Recall the definitions of  $s_+ < \infty$  from (1.4) and of  $r_\pm$  from (1.6). Recall further that we assumed  $\text{supp } \mu_1 \subset [-s_+, s_+]$ .

Given  $r \in (r_-, r_+)$ , we define

$$\sigma_- \equiv \sigma_-(r) := \left( \frac{r^2 - r_-^2}{r_+^2 - r_-^2} \right)^{\frac{1}{2}}, \quad \sigma_+ \equiv \sigma_+(r) := \left( \frac{r_+^2}{r_+^2 - r^2} \right)^{\frac{1}{2}}, \quad (7.1)$$

and we note that  $\sigma_-(r) \in (0, 1)$  and  $\sigma_+(r) \in (1, \infty)$ .

The main result of this section is the following bound on  $\omega_2$  along the imaginary axis.

**Theorem 7.1.** (Bounds on  $\omega_2$ ) *Let  $r \in (r_-, r_+)$ . Then there are strictly positive constants  $b_- \equiv b_-(\mu_1, r) > 0$  and  $s_- \equiv s_-(\mu_1, r) > 0$  such that*

$$C_2^{-1} \sigma_- s_- b_- \min \left\{ 1, \frac{\sigma_- s_-}{\eta} \right\} \leq |\omega_2(i\eta) - i\eta| \leq C_2 \min \left\{ \sigma_+ s_+, \frac{r^2}{\eta} \right\}, \quad (7.2)$$

for all  $\eta \geq 0$ , for a numerical constant  $C_2 > 1$  (independent of  $\mu_1$  and  $r$ ).

*Remark 7.2.* It follows from the proof of Theorem 7.1, that  $s_-(\mu_1, r)$  is a monotonic increasing function in the variable  $r \in (r_-, r_+)$ , with  $0 < s_-(\mu_1, r) \leq s_+$ . The  $r$  dependence of  $b_-(\mu_1, r)$  (and  $a_-(\mu_1, r)$  in (7.19)) is then explicit in terms of  $s_-(\mu_1, r)$  and  $r$ . This allows us to obtain uniform bounds for  $r \in [r_- + \tau, r_+ - \tau]$ , with a fixed  $\tau > 0$ , in Theorem 2.2.

With this proposition we can easily establish all necessary bounds on the free convolution measure and both associated subordination functions stated in Theorem 2.2.

*Proof of Theorem 2.2.* Using (2.10) and the facts  $\omega_1(i\eta) = -\overline{\omega_1(i\eta)}$  and  $\omega_2(i\eta) = -\overline{\omega_2(i\eta)}$ , Theorem 2.2 follows readily. Indeed, the upper bound on  $|\omega_2(i\eta)|$  follows from the upper bound in (7.2), the lower bound on  $\text{Im } \omega_2(i\eta)$  follows from the lower bound in (7.2) for small  $\eta$  and from  $\text{Im } \omega_2(i\eta) \geq \eta$ , for large  $\eta$ . The upper and lower bounds on  $\text{Im } \omega_2(i\eta) - i\eta$  then imply a lower and upper bound on  $|\omega_1(i\eta)|$  by (2.10). Finally, (2.9) controls  $F_{\mu_1}(\omega_2(i\eta)) = -1/m_{\mu_1 \boxplus \mu_2}(i\eta)$  from above and  $\text{Im } F_{\mu_1}(\omega_2(i\eta)) \geq \text{Im } (\omega_2(i\eta) - i\eta)$  controls it from below by (7.2), which yields (2.13).  $\square$

**7.1. Proof of Theorem 7.1.** For the sake of simplicity of presentation, the proof of Theorem 7.1 is accomplished in a sequence of lemmas.

**Lemma 7.3.** *Let  $\mu_1$  be as in Theorem 7.1. Then there exists a symmetric (non-negative) Borel measure  $\tilde{\mu}_1$  such that*

$$F_{\mu_1}(\omega) - \omega = \frac{-r^2}{\omega} + \int_{\mathbb{R}} \frac{d\tilde{\mu}_1(x)}{x - \omega}, \quad \omega \in \mathbb{C}^+, \quad (7.3)$$



with  $\tilde{\mu}_1(\mathbb{R}) = r_+^2 - r_-^2$ ,  $\text{supp } \tilde{\mu}_1 \subset [-s_+, s_+]$  and  $\tilde{\mu}_1(\{0\}) = 0$ .

*Proof.* Since  $\mu_1$  is a symmetric probability measure,  $F_{\mu_1} : \mathbb{C}^+ \rightarrow \mathbb{C}^+$  satisfies (2.3) and there exists a symmetric Borel measure,  $\hat{\mu}_1$ , such that  $F_{\mu_1}$  admits the Nevanlinna representation

$$F_{\mu_1}(\omega) = \omega + \int_{\mathbb{R}} \frac{d\hat{\mu}_1(x)}{x - \omega}, \quad \omega \in \mathbb{C}^+; \quad (7.4)$$

see *e.g.*, Proposition 2.2 in [31]. We observe that  $\hat{\mu}_1(\mathbb{R}) = r_+^2$ . Indeed, expanding the Stieltjes transform  $m_{\mu_1}$  around complex infinity we find

$$\begin{aligned} m_{\mu_1}(\omega) &= \int_{\mathbb{R}} \frac{d\mu_1(x)}{x - \omega} = \int_{\mathbb{R}} \left( \frac{-1}{\omega} - \frac{x^2}{\omega^3} + O(|\omega|^{-5}) \right) d\mu_1(x) \\ &= -\frac{1}{\omega} - \frac{r_+^2}{\omega^3} + O(|\omega|^{-5}), \end{aligned}$$

as  $|\omega| \rightarrow \infty$  in  $\mathbb{C}^+$ , where we used that  $\mu_1$  is symmetric. Thus  $F_{\mu_1}(\omega) - \omega = -r_+^2/\omega + O(|\omega|^{-3})$  in the same limit and we conclude by comparing with (7.4) that  $\hat{\mu}_1(\mathbb{R}) = r_+^2$ .

Since  $\mu_1$  is symmetric, its Stieltjes transform satisfies  $m_{\mu_1}(i\eta) = i \text{Im } m_{\mu_1}(i\eta)$ ,  $\eta > 0$ . We then obtain

$$\begin{aligned} \lim_{\eta \searrow 0} i\eta(F_{\mu_1}(i\eta) - i\eta) &= -\lim_{\eta \searrow 0} \frac{\eta}{\text{Im } m_{\mu_1}(i\eta)} \\ &= -\lim_{\eta \searrow 0} \left( \int_{\mathbb{R}} \frac{d\mu_1(x)}{x^2 + \eta^2} \right)^{-1} = -r_-^2, \end{aligned} \quad (7.5)$$

by the definition of  $r_-$  in (1.6). Comparing with (7.4), we conclude that  $\hat{\mu}_1(\{0\}) = r_-^2$ , since for any Borel measure  $\nu$  we have  $\lim_{\eta \searrow 0} \eta \text{Im } m_{\nu}(E + i\eta) = \nu(\{E\})$ , for all  $E \in \mathbb{R}$ . Setting  $\tilde{\mu}_1 := \hat{\mu}_1 - r_-^2 \delta_0$  we get (7.3). Clearly,  $\tilde{\mu}_1$  is a symmetric (non-negative) Borel measure with  $\tilde{\mu}_1(\mathbb{R}) = r_+^2 - r_-^2$  satisfying  $\text{supp } \tilde{\mu}_1 \subset [-s_+, s_+]$ . This concludes the proof of the lemma.  $\square$

We now introduce  $s_- \in \mathbb{R}^+$  as

$$s_- := \sup \left\{ x \in \mathbb{R}^+ : \int_0^x d\tilde{\mu}_1(x) \leq \frac{r^2 - r_-^2}{8} \right\}. \quad (7.6)$$

Note that, for  $r > r_-$ ,  $s_-$  is strictly positive since  $\mu_1$  is symmetric and we assume that  $\mu_1$  is supported at least at three points. (We assume that  $\mu_{\sigma}$  is supported at least at two points. Thus  $\mu_1 = \mu_{\sigma}^{\text{sym}}$  is supported at least at three points). Note that since  $\text{supp } \tilde{\mu}_1 \subset [-s_+, s_+]$ , thus  $\tilde{\mu}_1([0, s_+]) = \tilde{\mu}_1(\mathbb{R}^+) = \frac{1}{2}(r_+^2 - r_-^2) > \frac{1}{8}(r^2 - r_-^2)$ , we have  $s_- \leq s_+$ .

Equation (2.9) for  $\omega_2(z)$ , when combined with (7.3), reads

$$F_{\mu_1}(\omega_2(z)) - \omega_2(z) = \frac{-r_-^2}{\omega_2(z)} + \int_{\mathbb{R}} \frac{d\tilde{\mu}_1(x)}{x - \omega_2(z)} = -z - \frac{r^2}{\omega_2(z) - z}, \quad (7.7)$$

$z \in \mathbb{C}^+$ . We then rewrite this last equation as

$$\begin{aligned} (r^2 - r_-^2)\omega_2(z) + r_-^2 z \\ = -z(\omega_2(z) - z)\omega_2(z) - (\omega_2(z) - z)\omega_2(z) \int_{\mathbb{R}} \frac{d\tilde{\mu}_1(x)}{x - \omega_2(z)}, \end{aligned} \quad (7.8)$$

$z \in \mathbb{C}^+$ . Our first goal is to show that  $\text{Im } \omega_2(0) \equiv \lim_{\eta \searrow 0} \text{Im } \omega_2(i\eta)$  is strictly positive.

**Lemma 7.4.** *Let  $\mu_1$  and  $r \in (r_-, r_+)$  be as in Theorem 7.1. Then,*

$$\text{Im } \omega_2(0) > \frac{\sqrt{3}}{2} \sigma_{-s_-}. \quad (7.9)$$

*Proof.* By Theorem 2.3 of [6],  $\omega_2(z)$  extends continuously to the real line. Choosing  $z = i\eta$  in (7.8) we can assume that  $\lim_{\eta \searrow 0} \operatorname{Im} \omega_2(i\eta) < \infty$ . By symmetry,  $\omega_2(i\eta) = -\overline{\omega_2(i\eta)}$ , we know that  $\omega_2(0)$  is purely imaginary. Assume first that  $\operatorname{Im} \omega_2(0) > 0$ . Taking the limit  $\eta \searrow 0$  in (7.8) and dividing through  $\omega_2(0)$  we get

$$(r^2 - r_-^2) = -\omega_2(0) \int_{\mathbb{R}} \frac{d\tilde{\mu}_1(x)}{x - \omega_2(0)} = \int_{\mathbb{R}} \frac{|\omega_2(0)|^2 d\tilde{\mu}_1(x)}{x^2 + |\omega_2(0)|^2}, \quad (7.10)$$

where we used that  $\omega_2(0)$  is purely imaginary. Recalling  $s_-$  in (7.6), we further get

$$\begin{aligned} \int_{\mathbb{R}} \frac{|\omega_2(0)|^2 d\tilde{\mu}_1(x)}{x^2 + |\omega_2(0)|^2} &\leq 2 \int_0^{s_-} \frac{|\omega_2(0)|^2 d\tilde{\mu}_1(x)}{x^2 + |\omega_2(0)|^2} + 2 \int_{s_-}^{s_+} \frac{|\omega_2(0)|^2 d\tilde{\mu}_1(x)}{x^2 + |\omega_2(0)|^2} \\ &\leq 2 \int_0^{s_-} d\tilde{\mu}_1(x) + 2|\omega_2(0)|^2 \int_{s_-}^{s_+} \frac{d\tilde{\mu}_1(x)}{x^2} \\ &\leq \frac{r^2 - r_-^2}{4} + |\omega_2(0)|^2 \frac{r_+^2 - r_-^2}{s_-^2}, \end{aligned} \quad (7.11)$$

where we also used that  $\tilde{\mu}_1(\mathbb{R}) = r_+^2 - r_-^2$ . Hence from (7.11) and (7.10), we conclude that

$$3 \frac{r^2 - r_-^2}{4} \leq |\omega_2(0)|^2 \frac{r_+^2 - r_-^2}{s_-^2}. \quad (7.12)$$

Thus, we get

$$\frac{\sqrt{3}}{2} \sigma_- s_- \leq \operatorname{Im} \omega_2(0), \quad (7.13)$$

provided that  $\operatorname{Im} \omega_2(0) > 0$ , where we used that  $|\omega_2(0)| = \operatorname{Im} \omega_2(0)$ .

To conclude the proof, we need to show that  $\lim_{\eta \searrow 0} \operatorname{Im} \omega_2(i\eta) > 0$ . Arguing by contradiction, we assume that  $\lim_{\eta \searrow 0} \operatorname{Im} \omega_2(i\eta) = 0$ . Choose an arbitrary  $\epsilon > 0$ . Letting  $\eta > 0$  be sufficiently small, we can assure that

$$\left| \omega_2(i\eta) \int_{\mathbb{R}} \frac{d\tilde{\mu}_1(x)}{x - \omega_2(i\eta)} \right| = \left| \int_{\mathbb{R}} \frac{|\omega_2(i\eta)|^2 d\tilde{\mu}_1(x)}{x^2 + |\omega_2(i\eta)|^2} \right| \leq \epsilon, \quad (7.14)$$

where we first used that  $\omega_2(i\eta)$  is purely imaginary and then used that 0 is not an atom of the measure  $\tilde{\mu}_1$ . We thus obtain from (7.8) and (7.14) that

$$|(r^2 - r_-^2)\omega_2(i\eta)| \leq |(r^2 - r_-^2)\omega_2(i\eta) + r_-^2 i\eta| \leq \eta |\omega_2(i\eta)|^2 + |\omega_2(i\eta)|\epsilon,$$

for  $\eta > 0$  sufficiently small, where we used  $r > r_-$ ,  $|\omega_2(i\eta)| = \operatorname{Im} \omega_2(i\eta)$ ,  $\operatorname{Im} \omega_2(i\eta) \geq \eta$  (c.f., Theorem 2.1), so  $|\omega_2(i\eta) - i\eta| \leq |\omega_2(i\eta)|$ . Choosing  $\epsilon = (r^2 - r_-^2)/2$ , we get

$$|(r^2 - r_-^2)\omega_2(i\eta)| \leq 2\eta |\omega_2(i\eta)|^2, \quad (7.15)$$

for  $\eta > 0$  sufficiently small, i.e., we have  $\operatorname{Im} \omega_2(i\eta) \geq (r^2 - r_-^2)/(2\eta)$ . Since  $r^2 - r_-^2 > 0$ , we get a contraction with the assumption that  $\lim_{\eta \searrow 0} \omega_2(i\eta) = 0$ . We thus conclude that  $\lim_{\eta \searrow 0} \omega_2(i\eta) > 0$ . This completes the proof of the lemma.  $\square$

We are now prepared to prove the lower bound in (7.2). Recall  $s_- > 0$  from (7.6).

**Lemma 7.5.** *Let  $\mu_1$  and  $r \in (r_-, r_+)$  be as in Theorem 7.1. Then, there is a strictly positive constant  $b_- \equiv b_-(\mu_1, r) > 0$  such that*

$$|\omega_2(i\eta) - i\eta| \geq C^{-1} \sigma_- s_- b_- \min \left\{ 1, \frac{\sigma_- s_-}{\eta} \right\}, \quad (7.16)$$

for all  $\eta \geq 0$ , where  $C > 1$  is a numerical constant (independent of  $\mu_1$  and  $r$ ).

*Proof.* Using the definition of  $s_-$  in (7.6), we write (7.8), the defining equation for  $\omega_2(z)$ , as an equation with a free variable  $\omega$ :

$$\frac{r^2 - r_-^2}{\omega} + \int_{|x| \leq s_-} \frac{d\tilde{\mu}_1(x)}{(x - \omega)} + \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{(x - \omega)} = -z - \frac{r^2 z}{(\omega - z)\omega}, \quad (7.17)$$

$z \in \mathbb{C}^+$ , whose unique solution on the upper half plane gives  $\omega = \omega_2(z)$ . Note that the third term on the left side has the expansion

$$\int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{(x - \omega)} = \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{x^2} \omega + \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{x^3(x - \omega)} \omega^3, \quad (7.18)$$

for  $|\omega| < s_-$ , where we used that  $\tilde{\mu}_1$  is symmetric to get the second line. Let

$$a_- \equiv a_-(\mu_1, r) := \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{x^2}. \quad (7.19)$$

Note that by the definition of  $s_-$  in (7.6) we have the bound

$$0 < \frac{3}{4} \frac{r^2 - r_-^2}{s_+^2} \leq a_- \leq \frac{r_+^2 - r_-^2}{s_-^2}. \quad (7.20)$$

Let moreover

$$\hat{\omega} := i \left( \frac{r^2 - r_-^2}{a_-} \right)^{\frac{1}{2}}. \quad (7.21)$$

Note that from (7.20), we have

$$\sigma_- s_- \leq |\hat{\omega}|. \quad (7.22)$$

Using the definitions of  $\hat{\omega}$  in (7.21) and of  $a_-$  in (7.19), we rewrite (7.17) as

$$(-\hat{\omega}^2 + \omega^2)(\omega - z) = -\frac{r^2 z}{a_-} + \psi(\omega, z), \quad z \in \mathbb{C}^+, \quad (7.23)$$

where we further introduced the shorthand notation

$$\begin{aligned} \psi(\omega, z) := & -\frac{\omega(\omega - z)z}{a_-} - \frac{\omega^4(\omega - z)}{a_-} \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{x^3(x - \omega)} \\ & - \frac{\omega(\omega - z)}{a_-} \int_{|x| \leq s_-} \frac{d\tilde{\mu}_1(x)}{x - \omega}. \end{aligned} \quad (7.24)$$

Next, we abbreviate  $t_- := s_- \sigma_- < s_-$  and define

$$b_- \equiv b_-(\mu_1, r) := \min \left\{ 1, a_-, a_- \frac{t_-^2}{r^2} \right\} > 0. \quad (7.25)$$

Then we introduce the sets

$$\mathcal{F}_- := \left\{ \omega = i|\omega| \in \mathbb{C}^+ : |\omega|^2 \leq \frac{7}{8} t_-^2 \right\} \quad (7.26)$$

$$\mathcal{E}_-^{(1)} := \left\{ z = i\eta \in \mathbb{C}^+ : \eta \leq \frac{t_- b_-}{64} \right\}. \quad (7.27)$$

For  $\omega \in \mathcal{F}_-$  we bound the last term in the definition of  $\psi(\omega, z)$  in (7.24) as

$$\begin{aligned} \frac{|\omega||\omega - z|}{a_-} \left| \int_{|x| \leq s_-} \frac{d\tilde{\mu}_1(x)}{x - \omega} \right| &= \frac{|\omega - z|}{a_-} \left| \int_{|x| \leq s_-} \frac{|\omega|^2 d\tilde{\mu}_1(x)}{x^2 + |\omega|^2} \right| \\ &\leq \frac{|\omega - z|}{4a_-} (r^2 - r_-^2) = \frac{|\hat{\omega}|^2}{4} |\omega - z|, \quad z \in \mathbb{C}^+, \end{aligned} \quad (7.28)$$

where we used that  $\tilde{\mu}_1$  is symmetric and the definitions of  $s_-$  in (7.6) and of  $\hat{\omega}$  in (7.21).

For  $\omega \in \mathcal{F}_-$  we bound the second but last term in the definition of  $\psi(\omega, z)$  as

$$\frac{|\omega|^4 |\omega - z|}{a_-} \left| \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{x^3(x - \omega)} \right| \leq \frac{|\omega|^4 |\omega - z|}{a_-} \int_{|x| > s_-} \frac{d\tilde{\mu}_1(x)}{x^2} \frac{1}{|s_-|^2} \leq \frac{7}{8} |\omega|^2 |\omega - z|, \quad (7.29)$$

where we used that  $|x - \omega| \geq |x| \geq s_-$ , as  $\omega$  lies on the imaginary axis, the definition of  $a_-$  in (7.19) and  $t_- = \sigma_- s_- \leq s_-$ .

For the first term on the right side of  $\psi(\omega, z)$ , we get for  $z \in \mathcal{E}_-^{(1)}$  the bound

$$\frac{|\omega(\omega - z)z|}{a_-} \leq \frac{|\omega - z| |\omega| t_-}{64}, \quad (7.30)$$

where we used that  $64|z| \leq t_- b_- \leq t_- a_-$  on  $\mathcal{E}_-^{(1)}$ .

Combining these estimates, we get that, for  $\omega \in \mathcal{F}_-$  and  $z \in \mathcal{E}_-^{(1)}$ ,

$$|\psi(\omega, z)| \leq \frac{|\omega| t_-}{64} |\omega - z| + \frac{7}{8} |\omega|^2 |\omega - z| + \frac{|\widehat{\omega}|^2}{4} |\omega - z|. \quad (7.31)$$

For the first term on the right side of (7.23) we note for  $z \in \mathcal{E}_-^{(1)}$  the bound

$$\frac{r^2 |z|}{a_-} \leq \frac{t_- r^2}{64 a_-} b_- \leq \frac{t_-^3}{64}, \quad (7.32)$$

where we used that  $b_- r^2 / a_- \leq t_-^2$  on  $\mathcal{E}_-^{(1)}$  as follows from (7.25).

Thus, for  $\omega$  a solution to (7.23) in  $\mathcal{F}_-$  with  $z \in \mathcal{E}_-^{(1)}$ , at least one of the following holds

$$|-\widehat{\omega}^2 + \omega^2| |\omega - z| \leq \frac{t_-^3}{32} \quad (7.33)$$

or

$$|-\widehat{\omega}^2 + \omega^2| |\omega - z| \leq \frac{|\omega| t_-}{32} |\omega - z| + \frac{7}{4} |\omega|^2 |\omega - z| + \frac{|\widehat{\omega}|^2}{2} |\omega - z|. \quad (7.34)$$

First, assume that (7.34) holds. Then we either have  $\omega - z = 0$ , or

$$|-\widehat{\omega}^2 + \omega^2| \leq \frac{t_-^2}{64} + \frac{|\omega|^2}{64} + \frac{7}{4} |\omega|^2 + \frac{|\widehat{\omega}|^2}{2}. \quad (7.35)$$

We then absorb the last term on the right side into the left side to get

$$\frac{1}{2} |\widehat{\omega}^2| \leq \frac{t_-^2}{64} + |\omega|^2 + \frac{|\omega|^2}{64} + \frac{7}{4} |\omega|^2. \quad (7.36)$$

We thus find

$$|\widehat{\omega}^2| < \frac{t_-^2}{32} + 6|\omega|^2. \quad (7.37)$$

Since  $|\widehat{\omega}| \geq t_-$  by (7.22), we thus obtain that in this case that either  $\omega - z = 0$  or

$$|\omega| > \frac{3}{8} t_-. \quad (7.38)$$

Second, assume that (7.33) holds. Then we can estimate, using that  $\omega \in \mathcal{F}_-$ ,

$$|\omega - z| \leq \frac{t_-^3}{32} \frac{1}{|-\widehat{\omega}^2 + \omega^2|} \leq \frac{2}{8} t_-, \quad (7.39)$$

where we used that  $|\widehat{\omega}| \geq t_-$  and  $|\omega|^2 \leq 7t_-^2/8$  on  $\mathcal{F}_-$ . Since  $|z| \leq t_-/64$  for  $z \in \mathcal{E}_-^{(1)}$ , we find  $|\omega| < 5t_-/16$  in this case.

We conclude that for any  $z \in \mathcal{E}_-^{(1)}$  a solution  $\omega(z)$  to (7.23) in  $\mathcal{F}_-$  satisfies either

$$|\omega(z)| < \frac{5}{16} t_- \quad \text{or} \quad |\omega(z)| > \frac{6}{16} t_-. \quad (7.40)$$

Also note that if a solution  $\omega(z)$  of (7.23) satisfies  $\omega(z) \notin \mathcal{F}_-$  for some  $z \in \mathcal{E}_-^{(1)}$ , then the second alternative in (7.40) holds trivially.

Now, since the subordination function  $\eta \mapsto \omega_2(i\eta)$  (extends to) a continuous function on  $[0, \infty)$  by Theorem 2.3 of [6], we can conclude from (7.40) that

$$|\omega_2(z)| > \frac{3}{8}t_-, \quad z \in \mathcal{E}_-^{(1)}, \quad (7.41)$$

since we already showed in (7.9) that  $|\omega_2(0)| \geq \sqrt{3}t_-/2$ . This proves the lower bound in (7.16) for the small  $\eta$  regime.

Next, we introduce the domain which will handle the regime complementary to  $\mathcal{E}_-^{(1)}$ ,

$$\mathcal{E}_-^{(2)} := \left\{ z = i\eta \in \mathbb{C}^+ : \eta \geq \frac{t_- b_-}{64} \right\}. \quad (7.42)$$

We claim that  $\eta \mapsto \eta \cdot (\operatorname{Im} \omega_2(i\eta) - \eta)$  is a monotone increasing function for  $\eta \in \mathbb{R}^+$ . Indeed since the analytic function  $\omega_2 : \mathbb{C}^+ \rightarrow \mathbb{C}^+$  satisfies (2.4) and  $\omega_2(i\eta) = -\overline{\omega_2(i\eta)}$ , it has the Nevanlinna representation

$$\omega_2(z) = z + \int_{\mathbb{R}} \frac{d\nu_2(x)}{x - z}, \quad z \in \mathbb{C}^+, \quad (7.43)$$

where  $\nu_2$  is a finite symmetric Borel measure. The claim follows directly by considering the imaginary part of (7.43) for  $z$  along the positive imaginary axis. Hence, for any  $\eta \geq \eta_0 > 0$ ,

$$\operatorname{Im} \omega_2(i\eta) - \eta \geq \frac{\eta_0}{\eta} (\operatorname{Im} \omega_2(i\eta_0) - \eta_0). \quad (7.44)$$

Choosing  $\eta_0 = t_- b_- / 64$  on the boundary of  $\mathcal{E}_-^{(1)}$ , we can apply (7.41) for  $z_0 = i\eta_0$ , and we obtain the estimate

$$\operatorname{Im} \omega_2(i\eta) - \eta \geq \frac{\eta_0}{\eta} \left( \frac{3t_-}{8} - \frac{t_- b_-}{64} \right) \geq \frac{2t_- \eta_0}{8\eta} \geq \frac{2t_-^2 b_-}{256\eta}, \quad i\eta \in \mathcal{E}_-^{(2)}, \quad (7.45)$$

where we used the definition of  $b_-$  in (7.25) to get the second inequality. Combining (7.41) and (7.45) we get the bound (7.16).  $\square$

We move on to the upper bound in (7.2).

**Lemma 7.6.** *Let  $\mu_1$  and  $r \in (r_-, r_+)$  be as in Theorem 7.1. Then,*

$$|\omega_2(i\eta) - i\eta| \leq C \min \left\{ \sigma_+ s_+, \frac{r^2}{\eta} \right\}, \quad (7.46)$$

for all  $\eta \geq 0$ , for a numerical constant  $C < \infty$  (independent of  $\mu_1$  and  $r$ ).

*Proof.* Using (7.4) we write

$$F_{\mu_1}(\omega) - \omega = -\frac{r_+^2}{\omega} - \frac{1}{\omega^2} \int_{\mathbb{R}} \frac{x^2 d\widehat{\mu}_1(x)}{x - \omega}, \quad \omega \in \mathbb{C}^+. \quad (7.47)$$

For  $z \in \mathbb{C}^+$ , we write (7.8) with  $\omega_2(z)$  replaced by the free variable  $\omega$  as

$$-\frac{r_+^2}{\omega} - \frac{\chi(\omega)}{\omega^2} = -z - \frac{r^2}{\omega - z}, \quad (7.48)$$

where we introduced the short hand notation

$$\chi(\omega) := \int_{\mathbb{R}} \frac{x^2 d\widehat{\mu}_1(x)}{x - \omega}.$$

From (7.48), we find that

$$\omega_2(z) = \frac{r_+^2 - r^2 + z^2}{2z} \left( 1 - \left( 1 - \frac{4z^2 r_+^2 - 4z\chi(\omega_2(z)) \frac{\omega_2(z) - z}{\omega_2(z)}}{(r_+^2 - r^2 + z^2)^2} \right)^{\frac{1}{2}} \right), \quad (7.49)$$

where we choose the square root such that  $\operatorname{Im} \omega_2(z) \geq \operatorname{Im} z$ .

Abbreviate  $t_+ := \sigma_+ s_+$  and partition the positive imaginary axis by introducing

$$\begin{aligned}\mathcal{E}_+^{(1)} &:= \left\{ z = i\eta \in \mathbb{C}^+ : 0 < \eta \leq \frac{1}{4} \frac{r_+^2 - r^2}{t_+} \right\}, \\ \mathcal{E}_+^{(2)} &:= \left\{ z = i\eta \in \mathbb{C}^+ : \frac{1}{4} \frac{r_+^2 - r^2}{t_+} < \eta \leq (r_+^2 - r^2)^{1/2} \right\}, \\ \mathcal{E}_+^{(3)} &:= \left\{ z = i\eta \in \mathbb{C}^+ : (r_+^2 - r^2)^{1/2} < \eta \right\}.\end{aligned}\tag{7.50}$$

We will prove the bound in (7.46) separately for these three regimes.

Choose  $z \in \mathcal{E}_+^{(1)}$  first. We will argue by contradiction that  $\text{Im } \omega_2(z) \leq 2t_+$  for this domain. Assuming that  $\omega \in \mathbb{C}^+$  with  $\text{Im } \omega > 2t_+$ , we have the simple bound

$$\left| \frac{\chi(\omega)}{r_+^2 - r^2} \right| = \left| \frac{1}{r_+^2 - r^2} \int_{\mathbb{R}} \frac{x^2 d\hat{\mu}_1(x)}{x - \omega} \right| \leq \frac{s_+^2 r_+^2}{r_+^2 - r^2} \frac{1}{2t_+} = \frac{t_+}{2},\tag{7.51}$$

where we used  $\hat{\mu}_1(\mathbb{R}) = r_+^2$ , and  $|x| \leq s_+$  on the support of  $\hat{\mu}_1$  in the first inequality and  $t_+ \geq s_+$  in the second. Now, for  $z \in \mathcal{E}_+^{(1)}$ , we have

$$\frac{4|z|^2 r_+^2}{(r_+^2 - r^2 - |z|^2)^2} \leq \frac{r_+^2}{16t_+^2} \frac{(r_+^2 - r^2)^2}{(r_+^2 - r^2 - |z|^2)^2} \leq \frac{16}{15^2 \sigma_+^2} < \frac{1}{10},\tag{7.52}$$

where we use that  $t_+ = \sigma_+ s_+$ ,  $\sigma_+ > 1$ ,  $s_+ \geq r_+$  and

$$\frac{|z|^2}{r_+^2 - r^2} \leq \frac{r_+^2 - r^2}{16t_+^2} < \frac{r_+^2}{16t_+^2} < \frac{1}{16}, \quad z \in \mathcal{E}_+^{(1)},\tag{7.53}$$

since  $t_+ > r_+$ . For  $z \in \mathcal{E}_+^{(1)}$ , we further have

$$\frac{4|z\chi(\omega_2(z))|}{(r_+^2 - r^2 - |z|^2)^2} \left| \frac{\omega_2(z) - z}{\omega_2(z)} \right| \leq \frac{1}{2t_+} \frac{t_+}{2} \frac{(r_+^2 - r^2)^2}{(r_+^2 - r^2 - |z|^2)^2} < \frac{3}{10},\tag{7.54}$$

where we used (7.51), and  $|\omega_2(z) - z| \leq |\omega_2(z)|$  and  $z \in \mathcal{E}_+^{(1)}$  to get the first inequality.

We then obtain from (7.49), upon expanding the square root using (7.52) and (7.54) that

$$|\omega_2(z)| \leq 2 \frac{|z| r_+^2}{r_+^2 - r^2 - |z|^2} + 2 \frac{|\chi(\omega_2)|}{r_+^2 - r^2 - |z|^2} \leq 2 \frac{16|z| r_+^2}{15(r_+^2 - r^2)} + \frac{16}{15} t_+ < 2t_+,\tag{7.55}$$

where we used (7.52), (7.53) and (7.54) to get the second inequality, and that  $z \in \mathcal{E}_+^{(1)}$  and  $r_+ \leq t_+$  to get the third. However, (7.55) yields a contradiction with the assumption that  $\text{Im } \omega_2(z) > 2t_+$ . We can therefore conclude that

$$\text{Im } \omega_2(z) \leq 2t_+, \quad z \in \mathcal{E}_+^{(1)}.\tag{7.56}$$

Choose now  $z \in \mathcal{E}_+^{(2)}$ . Starting from (7.49), we estimate

$$\begin{aligned}|\omega_2(z)| &\leq \frac{1}{2|z|} \left( 2(r_+^2 - r^2) + 2|z|r_+ + 2 \left( \frac{|z| s_+^2 r_+^2}{|\omega_2(z)|} \right)^{\frac{1}{2}} \right) \\ &\leq \frac{4t_+}{(r_+^2 - r^2)} (r_+^2 - r^2) + 2 \frac{r_+}{2} + 2 \left( \frac{s_+^2 r_+^2}{4|z\omega_2(z)|} \right)^{\frac{1}{2}} \\ &\leq 5t_+ + 2 \left( \frac{t_+ s_+^2 r_+^2}{(r_+^2 - r^2) |\omega_2(z)|} \right)^{\frac{1}{2}}, \quad z \in \mathcal{E}_+^{(2)},\end{aligned}$$

where we used  $z^2 < 0$ ,  $r_+ \leq t_+$ ,  $(r_+^2 - r^2)/(4t_+) \leq |z|$ ,  $|\omega_2(z) - z| \leq |\omega_2(z)|$  and

$$|\chi(\omega_2(z))| \leq \frac{s_+^2 r_+^2}{\text{Im } \omega_2(z)},$$

with  $\operatorname{Im} \omega_2(z) = |\omega_2(z)|$ , for  $z \in \mathcal{E}_+^{(2)}$ . Thus at least one of the following bounds holds

$$|\omega_2(z)| \leq 10t_+ \quad \text{or} \quad |\omega_2(z)| \leq 2 \left( \frac{4t_+ s_+^2 r_+^2}{(r_+^2 - r^2) |\omega_2(z)|} \right)^{\frac{1}{2}}.$$

In the latter case we find that

$$|\omega_2(z)| \leq \left( \frac{16t_+ s_+^2 r_+^2}{r_+^2 - r^2} \right)^{\frac{1}{3}} = (16\sigma_+^2 s_+^2 t_+)^{\frac{1}{3}} < 3t_+.$$

Thus in both cases we have

$$|\omega_2(z)| \leq 10t_+, \quad z \in \mathcal{E}_+^{(2)}. \quad (7.57)$$

Finally, we consider  $z \in \mathcal{E}_+^{(3)}$ . Since  $\operatorname{Im} \omega_2(z) \geq \operatorname{Im} z$ , we can expand (7.47) as

$$F_{\mu_1}(\omega_2(z)) - \omega_2(z) = -\frac{r_+^2}{\omega_2(z)} + O(|\omega_2(z)|^{-3}) = O(|z|^{-1}), \quad (7.58)$$

as  $\operatorname{Im} z \nearrow \infty$ , which in turn implies through (7.48) that

$$\omega_2(z) = z - \frac{r^2}{z + O(|z|^{-1})},$$

as  $\operatorname{Im} z \nearrow \infty$ . Comparison with the Nevanlinna representation of  $\omega_2(z)$  in (7.43) reveals that  $\nu_2(\mathbb{R}) = r^2$ . Using (7.43) we can therefore estimate  $\omega_2(z)$  from above as

$$|\omega_2(z) - z| \leq \frac{r^2}{\operatorname{Im} z}, \quad z \in \mathbb{C}^+. \quad (7.59)$$

In particular, using  $r < r_+ \leq s_+$  we have, for  $i\eta \in \mathcal{E}_+^{(3)}$ , that

$$|\omega_2(i\eta) - i\eta| \leq \frac{r^2}{\eta} \leq \sigma_+ s_+ = t_+, \quad z \in \mathcal{E}_+^{(3)}. \quad (7.60)$$

Combining (7.56), (7.57) and (7.60), we see that there is a numeral  $C$  such that

$$|\omega_2(i\eta) - i\eta| \leq \min \left\{ C\sigma_+ s_+, \frac{r^2}{\eta} \right\} \leq C\sigma_+ s_+, \quad \eta > 0. \quad (7.61)$$

This proves (7.46) and concludes the proof of the lemma.  $\square$

*Proof of Theorem 7.1.* Theorem 7.1 follows by combining (7.16) and (7.46), and adjusting the numerical constants.  $\square$

## 8. PROOF OF THEOREM 4.3 FOR LARGE $\eta$

In this section, we prove Theorem 4.3 for spectral parameters  $z \in \mathbb{C}^+$  with large imaginary parts,  $\eta$ . Here, large  $\eta$  means  $\eta \geq \eta_M$ , for some  $\eta_M \geq 1$  independent of  $N$  to be chosen below.

**8.1. Concentration of  $m_H$  for large  $\eta$ .** In this subsection, we fix an arbitrary  $L > 0$  and a compact interval  $\mathcal{I} \subset \mathbb{R}$ , and consider the domain  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$  introduced in (4.9).

*Proof of (4.11) on  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ .* In this proof, we choose both matrices  $U$  and  $V$  to be either Haar distributed on  $U(N)$  or  $O(N)$ , *i.e.*, we treat the unitary and orthogonal case at once. For simplicity we refer to  $U$  and  $V$  as Haar matrices below.

Our proof consists of two main steps. In the first step, we shall show that

$$\left| m_H(z) - m_A(\omega_B^c(z)) \right| \prec \frac{1}{N\eta^2}, \quad \left| m_H(z) - m_B(\omega_A^c(z)) \right| \prec \frac{1}{N\eta^2}, \quad (8.1)$$

uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ . In the second step, we use the local stability of the system (2.5) with the choice  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$  to conclude (4.11) from (8.1) for large  $\eta$ .

*Step 1: Proof of (8.1).* This step is based on the Gromov-Milman concentration inequality. Let  $\mathfrak{M} \equiv \mathfrak{M}(N)$  stand for the fundamental representation of either  $U(N)$  or  $O(N)$  on  $M_N(\mathbb{C})$ ,

and let  $\mathfrak{M}_1 \equiv \mathfrak{M}_1(N)$  stand for the fundamental representation of either  $SU(N)$  or  $SO(N)$  on  $M_N(\mathbb{C})$ , all endowed with the Riemann metric  $\|ds\|_2$  inherited from  $M_N(\mathbb{C})$  (equipped with the Hilbert-Schmidt norm  $\|\cdot\|_2$ ). We denote by  $\mathbb{P}_{\mathfrak{M}}$ ,  $\mathbb{P}_{\mathfrak{M}_1}$  (the push-forwards of) the Haar measure on  $\mathfrak{M}$ ,  $\mathfrak{M}_1$  respectively. We use the following version of the Gromov-Milman concentration inequality formulated as Corollary 4.4.28 in [1]. If  $g : (\mathfrak{M}(N), \|ds\|_2) \rightarrow \mathbb{C}$  is an  $\mathcal{L}$ -Lipschitz function then

$$\mathbb{P}_{\mathfrak{M}}\left(\left|g(\cdot) - \int_{\mathfrak{M}_1} g(W \cdot) d\mathbb{P}_{\mathfrak{M}_1}(W)\right| > \delta\right) \leq C e^{-c \frac{N\delta^2}{\mathcal{L}^2}}, \quad (8.2)$$

for all  $\delta > 0$ , where  $c > 0$  and  $C$  are numerical constants.

To apply (8.2) with the Haar matrices  $U$  and  $V$  at once, we extend (8.2) to the direct product group  $\mathfrak{M} \times \mathfrak{M}$  by adjusting the constants  $c > 0$  and  $C$ ; see *e.g.*, Theorem 1.11 of [30].

For any deterministic matrix  $Q \in M_{2N}(\mathbb{C})$ , we introduce

$$f(Q, \mathcal{U}, z) := \text{tr} QG(z), \quad z \in \mathbb{C}^+, \quad (8.3)$$

where  $\mathcal{U}$  is given in terms of  $U$  and  $V$  as in (4.14), *i.e.*, is a Haar matrix on  $\mathfrak{M} \times \mathfrak{M}$ . We will view  $f(Q, \mathcal{U}, z)$  as a random variable on  $\mathfrak{M} \times \mathfrak{M}$ . To apply (8.2), we estimate the Lipschitz constant of  $f(Q, \cdot, z) : \mathfrak{M} \times \mathfrak{M} \rightarrow \mathbb{C}$ .

Denote by  $\mathfrak{m} \oplus \mathfrak{m}$  the (fundamental representation of the) Lie algebra of  $\mathfrak{M} \times \mathfrak{M}$ . Note that  $X \in \mathfrak{m} \oplus \mathfrak{m}$  is a blockdiagonal matrix satisfying  $X = -X^*$ . For  $X \in \mathfrak{m} \oplus \mathfrak{m}$  let  $\text{ad}_X : \mathfrak{m} \oplus \mathfrak{m} \rightarrow \mathfrak{m} \oplus \mathfrak{m}$ ,  $Y \mapsto [X, Y]$  with  $[\cdot, \cdot]$  the Lie bracket of  $\mathfrak{m} \oplus \mathfrak{m}$ , *i.e.*, the commutator on  $M_{2N}(\mathbb{C})$ . Let  $B$  be as in (4.14). Then for  $X \in \mathfrak{m} \oplus \mathfrak{m}$  and  $t \in \mathbb{R}$ , we have  $e^{t\text{ad}_X}(UBU^*) = (e^{tX}U)B(e^{tX}U)^*$ , where we used that  $X = -X^*$ . Furthermore, note that

$$\frac{d}{dt} e^{t\text{ad}_X}(UBU^*) = e^{t\text{ad}_X} \text{ad}_X(UBU^*). \quad (8.4)$$

For  $X \in \mathfrak{m} \oplus \mathfrak{m}$ , we then compute, using (8.4) and  $\tilde{B} = UBU^*$ , that

$$\left. \frac{d}{dt} f(Q, e^{tX}U, z) \right|_{t=0} = -\text{tr} QG(\text{ad}_X \tilde{B})G = -\frac{1}{N} \text{Tr} QG(\text{ad}_X \tilde{B})G. \quad (8.5)$$

We thus get the bound

$$\left| \left. \frac{d}{dt} f(Q, e^{tX}U, z) \right|_{t=0} \right| \leq \frac{2}{N} \|B\| \|QG\| \|GX\|_1 \leq \frac{C\|QG\|}{N} \|G\|_2 \|X\|_2, \quad (8.6)$$

where  $\|\cdot\|_1$  denotes the trace norm. We used Schwarz inequality and  $\|B\| \leq C$  by assumption to get the last inequality. Since  $|G(z)|^2 = \frac{\text{Im} G(z)}{\eta}$  and  $\|G(z)\| \leq \eta^{-1}$ , we get from (8.6) that

$$\left| \left. \frac{d}{dt} f(Q, e^{tX}U, z) \right|_{t=0} \right| \leq \frac{C\|Q\|}{\sqrt{N}\eta^2} \|X\|_2, \quad z \in \mathbb{C}^+. \quad (8.7)$$

Thus the Lipschitz constant of  $f(Q)$  is bounded above by  $C\|Q\|/(\sqrt{N}\eta^2)$ . We therefore obtain from (8.2) the concentration inequality

$$\left| f(Q, \mathcal{U}, z) - \int_{\mathfrak{M}_1 \times \mathfrak{M}_1} f(Q, \mathcal{W} \cdot \mathcal{U}, z) d\mathbb{P}_{\mathfrak{M}_1 \times \mathfrak{M}_1}(\mathcal{W}) \right| \prec \frac{\|Q\|}{N\eta^2}, \quad (8.8)$$

where the randomness behind the notation  $\prec$  is provided by the Haar measure on  $\mathfrak{M} \times \mathfrak{M}$ .

We next identify the average appearing on the left side of (8.8). For a function  $g : \mathfrak{M} \times \mathfrak{M} \rightarrow \mathbb{C}$ ,  $\mathcal{U} \mapsto g(\mathcal{U})$ , we introduce the shorthand

$$\tilde{\mathbb{E}}g(\mathcal{U}) := \int_{\mathfrak{M}_1 \times \mathfrak{M}_1} g(\mathcal{W} \cdot \mathcal{U}) d\mathbb{P}_{\mathfrak{M}_1 \times \mathfrak{M}_1}(\mathcal{W}). \quad (8.9)$$

Using the invariance of Haar measure on  $\mathfrak{M}_1 \times \mathfrak{M}_1$ , we are going to compute  $\tilde{\mathbb{E}}\text{tr} G(z)$ . Denote by  $\mathfrak{m}_1 \oplus \mathfrak{m}_1$  the Lie algebra of  $\mathfrak{M}_1 \times \mathfrak{M}_1$ . The following argument is essential due to [33]; see



also [10, 28] for similar arguments. Viewing the Green function as a function (random variable) on  $\mathfrak{M} \times \mathfrak{M}$ ,  $G(\cdot, z) : \mathfrak{M} \times \mathfrak{M} \rightarrow M_N(\mathbb{C})$ , we compute, using (8.4), that

$$\tilde{\mathbb{E}}\left[\frac{d}{dt}G(e^{tX}\mathcal{U}, z)\Big|_{t=0}\right] = -\tilde{\mathbb{E}}\left[G(\mathcal{U}, z)\text{ad}_X(\tilde{B})G(\mathcal{U}, z)\right], \quad (8.10)$$

for any  $X \in \mathfrak{m}_1 \oplus \mathfrak{m}_1$ , where  $\tilde{B} \equiv \mathcal{U}B\mathcal{U}^*$ . On the other hand, by the left-invariance of Haar measure, we also have  $\frac{d}{dt}\tilde{\mathbb{E}}G(e^{tX}\mathcal{U}, z) = 0$ , for all  $t \in \mathbb{R}$  and all  $X \in \mathfrak{m}_1 \oplus \mathfrak{m}_1$ . Thus we get from (8.10) that  $\tilde{\mathbb{E}}[G(\mathcal{U}, z)\text{ad}_X(\tilde{B})G(\mathcal{U}, z)] = 0$ , for any  $X \in \mathfrak{m}_1 \oplus \mathfrak{m}_1$ , *i.e.*, we have

$$\tilde{\mathbb{E}}[G(\mathcal{U}, z)[X, \tilde{B}]G(\mathcal{U}, z)] = 0. \quad (8.11)$$

Such formulas originating from basic symmetries of the model are often called *Ward identities* in physics. Let now  $Y := \hat{e}_i \hat{e}_k^*$ , with  $i \neq k$ ,  $i, k \in \llbracket 1, N \rrbracket$ . We then note that we can decompose  $Y = \frac{1}{2}X_1 + \frac{1}{2i}X_2$ , where  $X_1 := Y - Y^*$  and  $X_2 := (iY + iY^*)$ . Note that  $X_1, X_2 \in \mathfrak{m}_1 \oplus \mathfrak{m}_1$ . Thus we have from (8.11) that

$$\tilde{\mathbb{E}}[G(\mathcal{U}, z)[X_\iota, \tilde{B}]G(\mathcal{U}, z)] = 0, \quad \iota = 1, 2. \quad (8.12)$$

Since  $Y$  is a linear combination of  $X_1$  and  $X_2$ , we conclude by the linearity of the commutator and (8.12) that, for  $i \neq k$ ,

$$\tilde{\mathbb{E}}[G(\mathcal{U}, z)[\hat{e}_i \hat{e}_k^*, \tilde{B}]G(\mathcal{U}, z)] = 0. \quad (8.13)$$

Next, recall the notational convention  $\hat{i} \equiv i + N$ , for  $i \in \llbracket 1, N \rrbracket$ . Using exactly the same argument as above we infer, for  $i \neq k$ ,  $i, k \in \llbracket 1, N \rrbracket$ , that

$$\tilde{\mathbb{E}}[G(\mathcal{U}, z)[\hat{e}_i \hat{e}_k^*, \tilde{B}]G(\mathcal{U}, z)] = 0. \quad (8.14)$$

Thus, taking matrix elements of (8.13) and (8.14), we obtain, for all  $i \in \llbracket 1, N \rrbracket$ ,  $j = i, \hat{i}$ ,

$$\left| \tilde{\mathbb{E}}\left[\tau_1(G(\mathcal{U}, z))(G(\mathcal{U}, z)\tilde{B})_{ji} - \tau_1(\tilde{B}G(\mathcal{U}, z))G_{ji}(\mathcal{U}, z)\right] \right| \leq \frac{C}{N\eta^2}, \quad (8.15)$$

and

$$\left| \tilde{\mathbb{E}}\left[\tau_2(G(\mathcal{U}, z))(G(\mathcal{U}, z)\tilde{B})_{j\hat{i}} - \tau_2(\tilde{B}G(\mathcal{U}, z))G_{j\hat{i}}(\mathcal{U}, z)\right] \right| \leq \frac{C}{N\eta^2}, \quad (8.16)$$

for some constant  $C$  depending only on  $\|B\|$ , where the error terms result from coincidences among indices when using (8.13) and (8.14). Here we also used  $\|G(z)\| \leq \eta^{-1}$ .

Suppressing for simplicity the  $z$ - and  $\mathcal{U}$ -dependences in the notation for the Green function, we next note the identities

$$\begin{aligned} (G\tilde{B})_{ii} &= 1 + zG_{ii} - \bar{\xi}_i G_{\hat{i}\hat{i}}, & (G\tilde{B})_{\hat{i}\hat{i}} &= -\xi_i G_{ii} + zG_{\hat{i}\hat{i}}, \\ (G\tilde{B})_{\hat{i}\hat{i}} &= -\bar{\xi}_i G_{\hat{i}\hat{i}} + zG_{\hat{i}\hat{i}}, & (G\tilde{B})_{ii} &= 1 + zG_{\hat{i}\hat{i}} - \xi_i G_{ii}, \end{aligned} \quad (8.17)$$

for all  $i \in \llbracket 1, N \rrbracket$ , which follow from (4.16). Plugging (8.17) into (8.15) and (8.16) we get

$$\begin{aligned} \left| \tilde{\mathbb{E}}\left[(1 + zG_{ii} - \bar{\xi}_i G_{\hat{i}\hat{i}})\tau_1(G) - G_{ii}\tau_1(\tilde{B}G)\right] \right| &\leq \frac{C}{N\eta^2}, \\ \left| \tilde{\mathbb{E}}\left[(-\bar{\xi}_i G_{\hat{i}\hat{i}} + zG_{\hat{i}\hat{i}})\tau_1(G) - G_{\hat{i}\hat{i}}\tau_1(\tilde{B}G)\right] \right| &\leq \frac{C}{N\eta^2}, \\ \left| \tilde{\mathbb{E}}\left[(-\xi_i G_{ii} + zG_{\hat{i}\hat{i}})\tau_2(G) - G_{\hat{i}\hat{i}}\tau_2(\tilde{B}G)\right] \right| &\leq \frac{C}{N\eta^2}, \\ \left| \tilde{\mathbb{E}}\left[(1 + zG_{\hat{i}\hat{i}} - \xi_i G_{ii})\tau_2(G) - G_{\hat{i}\hat{i}}\tau_2(\tilde{B}G)\right] \right| &\leq \frac{C}{N\eta^2}. \end{aligned} \quad (8.18)$$

Next, by (8.8) we have the concentration inequalities

$$\left| \tau_a(G) - \tilde{\mathbb{E}}[\tau_a(G)] \right| \prec \frac{1}{N\eta^2}, \quad \left| \tau_a(\tilde{B}G) - \tilde{\mathbb{E}}[\tau_a(\tilde{B}G)] \right| \prec \frac{1}{N\eta^2}.$$

For the second estimate we used that  $\tau_a(\tilde{B}G)$  can be brought into the form  $\text{tr } QG$  with a deterministic  $Q$  with the help of (4.16). Hence, we can go back and forth between the tracial quantities  $\tau_a(G)$ ,  $\tau_a(\tilde{B}G)$  and their partial expectations  $\tilde{\mathbb{E}}\tau_a(G)$  and  $\tilde{\mathbb{E}}\tau_a(\tilde{B}G)$ , up to an error  $O_{\prec}(\frac{1}{N\eta^2})$  in the following discussion. Pulling out the expectation of the tracial quantities and combining the first and the third equations in (8.18) we eliminate  $\tilde{\mathbb{E}}G_{ii}$  and get an equation for  $\tilde{\mathbb{E}}G_{ii}$ . After solving for  $\tilde{\mathbb{E}}G_{ii}$ , we may remove the partial expectation  $\tilde{\mathbb{E}}$  from the tracial quantities. We get

$$\begin{aligned} & \left( (z\tau_1(G) - \tau_1(\tilde{B}G))(z\tau_2(G) - \tau_2(\tilde{B}G)) - |\xi_i|^2\tau_1(G)\tau_2(G) \right) \tilde{\mathbb{E}}[G_{ii}] \\ & \quad + \tau_1(G)(z\tau_2(G) - \tau_2(\tilde{B}G)) = O_{\prec}\left(\frac{1}{N\eta^2}\right). \end{aligned}$$

Here we used once more the bound  $\|G\| \leq 1/\eta$ . Dividing the above equation by  $\tau_1(G)\tau_2(G)$  and using the fact  $|\tau_a(G) - \frac{1}{\eta}| \leq O(\eta^{-2})$ , we obtain

$$(\omega_{B,1}^c \omega_{B,2}^c - |\xi_i|^2) \tilde{\mathbb{E}}[G_{ii}] + \omega_{B,2}^c = O_{\prec}\left(\frac{1}{N}\right), \quad (8.19)$$

for all  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ , by choosing  $\eta_M > 0$  sufficiently large. Here we introduced the auxiliary subordination functions

$$\omega_{B,a}^c(z) := z - \frac{\tau_a(\tilde{B}G(z))}{\tau_a(G(z))}, \quad a = 1, 2, \quad z \in \mathbb{C}^+, \quad (8.20)$$

which are defined using the partial traces  $\tau_a$  instead of the full traces as in (4.17).

We further observe that a large  $z$  expansion in  $\mathbb{C}^+$  of the resolvent yields

$$\tau_a(\tilde{B}G(z)) = -\frac{\tau_a(\tilde{B})}{z} + O\left(\frac{1}{|z|^2}\right) = O\left(\frac{1}{|z|^2}\right), \quad (8.21)$$

as  $|z| \rightarrow \infty$ , where we used that  $\tau_a(\tilde{B}) = 0$ . Thus from (8.20) we find that

$$\omega_{B,a}^c(z) = z + O(|z|^{-1}), \quad a = 1, 2, \quad (8.22)$$

as  $|z| \rightarrow \infty$ . Combining (8.19) and (8.22) we find

$$\tilde{\mathbb{E}}[G_{ii}(z)] - \frac{\omega_{B,2}^c(z)}{|\xi_i|^2 - \omega_{B,1}^c(z)\omega_{B,2}^c(z)} = O_{\prec}\left(\frac{1}{N\eta^2}\right), \quad \forall i \in \llbracket 1, N \rrbracket, \quad (8.23)$$

for all  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ , for sufficiently large  $\eta_M > 0$ . Analogously, we have

$$\tilde{\mathbb{E}}[G_{ii}(z)] - \frac{\omega_{B,1}^c(z)}{|\xi_i|^2 - \omega_{B,1}^c(z)\omega_{B,2}^c(z)} = O_{\prec}\left(\frac{1}{N\eta^2}\right), \quad \forall i \in \llbracket 1, N \rrbracket, \quad (8.24)$$

for all  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ . From  $\tau_1(G) = \tau_2(G)$ , see (4.27), and from (8.22), we obtain from (8.23) and (8.24) that

$$\omega_{B,1}^c = \omega_{B,2}^c + O_{\prec}\left(\frac{1}{N}\right) = \omega_B^c + O_{\prec}\left(\frac{1}{N}\right),$$

where  $\omega_B^c$  is defined in (4.17). The second equality follows from the fact that  $\tau_1(G) = \tau_2(G)$  implies that this common value is  $\text{tr } G$ . Hence  $\omega_{B,1}^c \approx \omega_{B,2}^c$  implies  $\tau_1(\tilde{B}G) \approx \tau_2(\tilde{B}G)$ , hence both are close to their average,  $\text{tr } \tilde{B}G$ . We therefore also have

$$\begin{aligned} & \tilde{\mathbb{E}}[G_{ii}(z)] - \frac{\omega_B^c(z)}{|\xi_i|^2 - (\omega_B^c(z))^2} = O_{\prec}\left(\frac{1}{N\eta^2}\right), \\ & \tilde{\mathbb{E}}[G_{ii}(z)] - \frac{\omega_B^c(z)}{|\xi_i|^2 - (\omega_B^c(z))^2} = O_{\prec}\left(\frac{1}{N\eta^2}\right), \quad \forall i \in \llbracket 1, N \rrbracket, \end{aligned} \quad (8.25)$$

uniformly on  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$  by choosing  $\eta_M > 0$  sufficiently large. Averaging (8.25) over  $i$  and using the concentration estimate (8.8) with  $Q = \hat{I}$ , we obtain the first estimate in (8.1).

The second estimate in (8.1) is obtained in the same way by interchanging the rôles of  $A$  and  $B$ . This completes the first step of the argument.

*Step 2: Stability analysis.* We move on to check the stability of the system

$$\Phi_{\mu_A, \mu_B}(\omega_A, \omega_B, z) = 0, \quad (8.26)$$

for  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ ; see (4.28) for the definition of  $\Phi_{\mu_A, \mu_B}$ . First, we will show that  $\omega_A^c(z)$  and  $\omega_B^c(z)$  approximately solve (8.26). Then we will conclude from Lemma A.2 of [5] that  $\omega_A^c(z)$  and  $\omega_B^c(z)$  are close to  $\omega_A(z)$  and  $\omega_B(z)$ .

Using that  $F_{\mu_A}(\omega_B^c(z)) = -1/m_{\mu_A}(\omega_B^c(z))$  and the identity (4.18), we can write

$$F_{\mu_A}(\omega_B^c(z)) - \omega_A^c(z) - \omega_B^c(z) + z = \frac{1}{m_H(z)m_A(\omega_B^c(z))} (m_A(\omega_B^c(z)) - m_H(z)). \quad (8.27)$$

From the resolvent expansion  $G(z) = -1/z + O(1/|z|^2)$  in the large  $|z|$  regime we have that  $|m_H(z) - i\eta^{-1}| \leq C\eta^{-2}$  and  $|\omega_B^c(z) - i\eta| \leq C\eta^{-1}$ ; for the latter estimate we also used  $\text{tr } \tilde{B} = 0$  in (4.17). Thus together with the estimates in (8.1), we get from (8.27) that

$$F_{\mu_A}(\omega_B^c(z)) - \omega_A^c(z) - \omega_B^c(z) + z = O_{\prec} \left( \frac{1}{N} \right),$$

uniformly in  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ , by choosing  $\eta_M > 0$  large enough. Analogously, we also have

$$F_{\mu_B}(\omega_A^c(z)) - \omega_A^c(z) - \omega_B^c(z) + z = O_{\prec} \left( \frac{1}{N} \right),$$

on the same domain. Hence we have

$$\|\Phi(\omega_A^c(z), \omega_B^c(z), z)\|_2 \prec N^{-1}, \quad (8.28)$$

for all  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ .

Next, observe that the deterministic bounds

$$|\omega_A^c(z) - z| \leq \frac{C}{|z|}, \quad |\omega_B^c(z) - z| \leq \frac{C}{|z|} \quad (8.29)$$

hold uniformly for all  $|z| \geq \eta_M$  with sufficiently large  $\eta_M$ . This follows from (4.17), a large  $z$  expansion of  $G(z)$  and  $\text{tr } A = \text{tr } B = 0$ . Consequently, it is easy to check the following deterministic bound also holds uniformly for all  $z$  with  $|z| \geq \eta_M$

$$\|\Phi(\omega_A^c(z), \omega_B^c(z), z)\|_2 \leq \frac{C}{|z|}. \quad (8.30)$$

Then we apply Lemma A.2 of [5]. Thanks to (8.29) and (8.30), the assumptions of Lemma A.2 of [5] are satisfied and we further conclude from (8.28) that

$$|\omega_a^c(z) - \omega_a(z)| \prec \frac{1}{N}, \quad a = A, B, \quad (8.31)$$

uniformly in  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$  by slightly adjusting the value of  $\eta_M$ .

Combining (8.31) with (8.1) and  $|m'_A(\omega)| = O(|\omega|^{-2}) = O(\eta^{-2})$  in the regime where  $\omega \approx i\eta$  and  $\eta$  is large, we find that

$$m_H(z) - m_A(\omega_B(z)) = O_{\prec} \left( \frac{1}{N\eta^2} \right), \quad z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L).$$

Finally, since  $m_{\mu_A \boxplus \mu_B} = m_A(\omega_B)$  we conclude the proof of (4.11) for  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ .  $\square$

**8.2. Green function subordination for large  $\eta$ .** In this subsection, we show the following subordination property for the Green function entries when  $\eta$  is large. Recall from (4.31) the control parameter  $\Lambda_d$ .

**Lemma 8.1.** *Under the conditions and with the notations of Theorem 4.3 there is a (large) constant  $\eta_M$  such that*

$$\Lambda_d(z) \prec \frac{1}{\sqrt{N\eta^4}}, \quad (8.32)$$

uniformly on  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ .

*Proof of Lemma 8.1.* Let  $\eta_M$  be as in Subsection 8.1. From (8.25) and (8.31) we directly get

$$\begin{aligned} \tilde{\mathbb{E}}[G_{ii}(z)] - \frac{\omega_B(z)}{|\xi_i|^2 - \omega_B^2(z)} &= O_{\prec}\left(\frac{1}{N\eta^2}\right), & \tilde{\mathbb{E}}[G_{i\bar{i}}(z)] - \frac{\omega_B(z)}{|\xi_i|^2 - \omega_B^2(z)} &= O_{\prec}\left(\frac{1}{N\eta^2}\right), \\ \tilde{\mathbb{E}}[G_{i\bar{i}}(z)] - \frac{\xi_i}{|\xi_i|^2 - \omega_B^2(z)} &= O_{\prec}\left(\frac{1}{N\eta^2}\right), & \tilde{\mathbb{E}}[G_{ii}(z)] - \frac{\bar{\xi}_i}{|\xi_i|^2 - \omega_B^2(z)} &= O_{\prec}\left(\frac{1}{N\eta^2}\right), \end{aligned} \quad (8.33)$$

for all  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ . Hence, it remains to show the concentration of these entries of the Green function. To this end, we regard, as in Subsection 8.1, the Green function entries as functions of  $\mathcal{U}$ , and use the Gromov-Milman concentration inequality in (8.2). The Lipschitz constant of  $G_{ij}(\cdot, z) : \mathfrak{M} \times \mathfrak{M} \rightarrow \mathbb{C}$ ,  $\mathcal{U} \mapsto G_{ij}(\mathcal{U}, z)$  is estimated by bounding, for  $X \in \mathfrak{m} \oplus \mathfrak{m}$ ,

$$\left| \frac{dG_{ij}(e^{tX}\mathcal{U}, z)}{dt} \Big|_{t=0} \right| = \left| \hat{e}_i^* G(\mathcal{U}, z) \text{ad}_X \tilde{B} G(\mathcal{U}, z) \hat{e}_j \right| \leq \frac{C\|X\|_2}{\eta^2},$$

with a constant  $C$  depending only on  $\|B\|$ , where we first used (8.4) and then Schwarz inequality. Thus by (8.2),

$$\left| G_{ij}(z) - \tilde{\mathbb{E}}[G_{ij}(z)] \right| \prec \frac{1}{\sqrt{N\eta^4}}, \quad z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L).$$

Combining these concentration results with (8.33) we find (8.32). To obtain uniform bounds in  $z \in \mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ , we can apply a simple lattice argument using the Lipschitz continuity of the Green function  $G(z)$  and of the two subordination functions  $\omega_A(z)$  and  $\omega_B(z)$ . See the proof of Theorem 1.8 in Section 3 for a similar argument. The uniform Lipschitz continuity of the subordination functions follows directly from their analyticity on  $\mathcal{S}_{\mathcal{I}}(\eta_M, N^L)$ . This completes the proof of Lemma 8.1.  $\square$

## APPENDIX A.

**A.1. Stochastic domination and large deviation properties.** Recall the stochastic domination in Definition 1.6. The relation  $\prec$  is transitive and it satisfies the following arithmetic rules: if  $X_1 \prec Y_1$  and  $X_2 \prec Y_2$  then  $X_1 + X_2 \prec Y_1 + Y_2$  and  $X_1 X_2 \prec Y_1 Y_2$ . Further assume that  $\Phi(v) \geq N^{-C}$  is deterministic and that  $Y(v)$  is a non-negative random variable satisfying  $\mathbb{E}[Y(v)]^2 \leq N^{C'}$  for all  $v$ . Then  $Y(v) \prec \Phi(v)$ , uniformly in  $v$ , implies  $\mathbb{E}[Y(v)] \prec \Phi(v)$ , uniformly in  $v$ .

Gaussian vectors have well-known large deviation properties. We will use them in the following form whose proof is standard.

**Lemma A.1.** *Let  $X = (x_{ij}) \in M_N(\mathbb{C})$  be a deterministic matrix and let  $\mathbf{y} = (y_i) \in \mathbb{C}^N$  be a deterministic complex vector. For a Gaussian real or complex random vector  $\mathbf{g} = (g_1, \dots, g_N) \in \mathcal{N}_{\mathbb{R}}(0, \sigma^2 I_N)$  or  $\mathcal{N}_{\mathbb{C}}(0, \sigma^2 I_N)$ , we have*

$$|\mathbf{y}^* \mathbf{g}| \prec \sigma \|\mathbf{y}\|_2, \quad |\mathbf{g}^* X \mathbf{g} - \sigma^2 \text{Ntr } X| \prec \sigma^2 \|X\|_2. \quad (\text{A.1})$$

**A.2. Bounds on subordination functions.** Let  $\mu_\alpha, \mu_\beta$  be two  $N$ -independent probability measures on  $\mathbb{R}$  which are compactly supported: there exists a constant  $L < \infty$  such that

$$\text{supp}(\mu_\alpha) \subset [-L, L], \quad \text{supp}(\mu_\beta) \subset [-L, L]. \quad (\text{A.2})$$

Let  $\omega_\alpha, \omega_\beta$  be the subordination functions defined via the system of equations (2.5). The following result is proved in [2].

**Lemma A.2** (Lemma 5.1 and Corollary 5.2 in [2]). *Suppose that neither  $\mu_\alpha$  nor  $\mu_\beta$  is a single point mass and at least of one of them is supported at more than two points. Assume in addition that (A.2) holds. Let  $\mathcal{I} \subset \mathcal{B}_{\mu_\alpha \boxplus \mu_\beta}$  be a compact non-empty interval in the bulk of  $\mu_\alpha \boxplus \mu_\beta$ . Fix any  $0 < \eta_M < \infty$ . Let  $\mu_A, \mu_B$  be ( $N$ -dependent) probability measures on  $\mathbb{R}$ . Then there exist constants  $b_0 > 0, k > 0$  and  $K < \infty, S < \infty$ , which depend only on  $\eta_M, L$  in (A.2), the interval  $\mathcal{I}$  and the measures  $\mu_\alpha$  and  $\mu_\beta$ , such that whenever*

$$d_L(\mu_A, \mu_\beta) + d_L(\mu_B, \mu_\beta) \leq b_0, \quad (\text{A.3})$$

then

$$\begin{aligned} \max_{z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)} |\omega_A(z)| &\leq K, & \max_{z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)} |\omega_B(z)| &\leq K, \\ \min_{z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)} \text{Im} \omega_A(z) &\geq k, & \min_{z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)} \text{Im} \omega_B(z) &\geq k, \\ \max_{z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)} |\omega'_A(z)| &\leq S, & \max_{z \in \mathcal{S}_{\mathcal{I}}(0, \eta_M)} |\omega'_B(z)| &\leq S, \end{aligned} \quad (\text{A.4})$$

for all  $N \geq N_0$  with some sufficiently large  $N_0$  depending only on  $\eta_M, L$  in (A.2), the interval  $\mathcal{I}$  and the measures  $\mu_\alpha$  and  $\mu_\beta$ . Here  $\omega_A, \omega_B$  denote the subordinations functions defined via (2.5) for the choice  $(\mu_1, \mu_2) = (\mu_A, \mu_B)$ .

**A.3. Bounded rank perturbation estimate.** At various places, we use the following perturbation estimate; see Section 3.2 of [3] for proof, for instance.

**Lemma A.3.** *Let  $D \in M_N(\mathbb{C})$  be Hermitian and let  $Q \in M_N(\mathbb{C})$  be arbitrary. Then, for any Hermitian matrix  $R \in M_N(\mathbb{C})$ , we have*

$$|\text{tr}(Q(D + R - z)^{-1}) - \text{tr}(Q(D - z)^{-1})| \leq \frac{\text{rank}(R)\|Q\|}{N\eta}, \quad z = E + i\eta \in \mathbb{C}^+. \quad (\text{A.5})$$

Lemma A.3 also has the following corollary.

**Corollary A.4.** *Let  $Q \in M_{2N}(\mathbb{C})$  be arbitrary matrix. Then there is a numerical constant  $C$  such that, with the notations defined in (5.5), (5.8) and (5.9), we have*

$$\begin{aligned} |\text{tr} QG - \text{tr} QG^{(i)}| &\leq \frac{C\|Q\|}{N\eta}, & |\text{tr} Q\tilde{B}G - \text{tr} Q\tilde{B}^{(i)}G| &\leq \frac{C\|Q\|}{N\eta}, \\ |\text{tr} Q\tilde{B}G - \text{tr} Q\tilde{B}^{(i)}G^{(i)}| &\leq \frac{C\|Q\|}{N\eta}, & |\text{tr} Q\tilde{B}G\tilde{B} - \text{tr} Q\tilde{B}^{(i)}G^{(i)}\tilde{B}^{(i)}| &\leq \frac{C\|Q\|}{N\eta}. \end{aligned} \quad (\text{A.6})$$

*Proof.* The first inequality follows from (A.5) directly since  $H^{(i)}$  is a bounded rank perturbation of  $H$ . Next, we show the second inequality. Note that

$$\text{tr} Q\tilde{B}^{(i)}G - \text{tr} Q\tilde{B}G = \text{tr} Q\tilde{B}^{(i)}G - \text{tr} Q\mathcal{R}_i\tilde{B}^{(i)}\mathcal{R}_iG. \quad (\text{A.7})$$

Denote by  $\hat{\mathbf{r}}_i^u = \ell_i^u(\hat{\mathbf{e}}_i + \mathbf{k}_i^u)$  and  $\hat{\mathbf{r}}_i^v = \ell_i^v(\hat{\mathbf{e}}_i + \mathbf{k}_i^v)$ . By the definition in (5.3) and (5.5), we have  $\mathcal{R}_i = \hat{I} - \hat{\mathbf{r}}_i^u(\hat{\mathbf{r}}_i^u)^* - \hat{\mathbf{r}}_i^v(\hat{\mathbf{r}}_i^v)^*$ . Then it is easy to check the right side of (A.7) is a sum of the terms of the form

$$\frac{\tilde{d}_i}{N}(\hat{\mathbf{r}}_i^a)^* \tilde{B}^{(i)}GQ\hat{\mathbf{r}}_i^b, \quad \frac{\tilde{d}_i}{N}(\hat{\mathbf{r}}_i^a)^*GQ\tilde{B}^{(i)}\hat{\mathbf{r}}_i^b, \quad \frac{\tilde{d}_i}{N}(\hat{\mathbf{r}}_i^a)^*GQ\hat{\mathbf{r}}_i^b, \quad (\text{A.8})$$

or products of some of them, for some  $\tilde{d}_i$  which could be different from one to another, up to the bound  $|\tilde{d}_i| \leq C$ . Here  $a, b = u, v$ . Clearly, the terms in (A.8) are all bounded by  $\frac{C\|Q\|}{N\eta}$ .

This proves the second estimate in (A.6). The third bound in (A.6) follows from the second one and (A.5) immediately. The last one can also be proved analogously.  $\square$

## APPENDIX B.

In this appendix, we bound the terms involving  $\Delta_R^u(i, k)$ , *i.e.*, the terms in (5.65), the last term of (5.78), (5.82) and the last term of (5.87). We summarize the bound in the next lemma.

**Lemma B.1.** *Let  $Q \in M_{2N}(\mathbb{C})$  be arbitrary, with  $\|Q\| \prec 1$ . Let  $X_i = \hat{I}$  or  $\tilde{B}^{(i)}$  and  $X = \hat{I}$  or  $A$ . Suppose that the assumptions in Theorem 5.2 hold. Then,*

$$|\varepsilon_{i2}| \prec \Psi^2, \quad |\varepsilon_{i3}| \prec \Psi^2, \quad (\text{B.1})$$

$$\left| \frac{1}{N} \sum_k^{(i)} \hat{e}_i^* X \Delta_G^u(i, k) \hat{e}_i \hat{e}_k^* X_i G \hat{e}_i \right| \prec \Psi^2, \quad (\text{B.2})$$

$$\left| \frac{1}{N} \sum_k^{(i)} (\mathbf{k}_i^u)^* \Delta_G^u(i, k) \hat{e}_i \hat{e}_k^* X_i G \hat{e}_i \right| \prec \Psi^2, \quad (\text{B.3})$$

$$\left| \frac{1}{N} \sum_k^{(i)} \text{tr} Q X \Delta_G^u(i, k) \hat{e}_k^* X_i G \hat{e}_i \right| \prec \Psi^4. \quad (\text{B.4})$$

*Proof of Lemma B.1.* Recalling (5.46), we see that  $\Delta_R^u(i, k)$  is a sum of terms of the form

$$\tilde{d}_i \tilde{g}_{ik}^u \boldsymbol{\alpha}_i \boldsymbol{\beta}_i^*,$$

for some  $\tilde{d}_i \in \mathbb{C}$  satisfying  $|\tilde{d}_i| \prec 1$ , and  $\boldsymbol{\alpha}_i, \boldsymbol{\beta}_i = \mathbf{e}_i$  or  $\mathbf{h}_i^u$ . Hereafter  $\tilde{d}_i$  can be different from line to line, up to the bound  $|\tilde{d}_i| \prec 1$  uniformly on  $\mathcal{S}_{\mathcal{L}}(\eta_m, \eta_M)$ . By (5.51), we see that  $\Delta_G^u(i, k)$  is a sum of the terms of the form

$$\tilde{d}_i \tilde{g}_{ik}^u G \hat{\boldsymbol{\alpha}}_i \hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G, \quad \tilde{d}_i \tilde{g}_{ik}^u G \mathcal{R}_i \tilde{B}^{(i)} \hat{\boldsymbol{\alpha}}_i \hat{\boldsymbol{\beta}}_i^* G, \quad (\text{B.5})$$

where  $\hat{\boldsymbol{\alpha}}_i, \hat{\boldsymbol{\beta}}_i = \hat{\mathbf{e}}_i$  or  $\mathbf{k}_i^u$ . Then, by the definition in (5.59), we see that  $\varepsilon_{i2}$  is a sum of the terms of the form

$$\frac{1}{N} \tilde{d}_i (\mathbf{k}_i^u)^* \tilde{B}^{(i)} G \hat{\boldsymbol{\alpha}}_i \hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i, \quad \frac{1}{N} \tilde{d}_i (\mathbf{k}_i^u)^* \tilde{B}^{(i)} G \mathcal{R}_i \tilde{B}^{(i)} \hat{\boldsymbol{\alpha}}_i \hat{\boldsymbol{\beta}}_i^* G \hat{\mathbf{e}}_i.$$

Note that using the trivial bound  $\|G\| \leq 1/\eta$ , the terms above are stochastically dominated by

$$\frac{1}{N\eta} |\hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i|, \quad \frac{1}{N\eta} |\hat{\boldsymbol{\beta}}_i^* G \hat{\mathbf{e}}_i|$$

respectively. It is easy to check

$$|\hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{\mathbf{e}}_i| \prec 1, \quad |\hat{\boldsymbol{\beta}}_i^* G \hat{\mathbf{e}}_i| \prec 1, \quad (\text{B.6})$$

for  $\hat{\boldsymbol{\beta}}_i = \hat{\mathbf{e}}_i$  or  $\mathbf{k}_i^u$ . This can be seen from the facts (5.62) and (5.94), and also the bounds (5.31), which hold under the assumption (5.27). From the above discussion, we can see that  $|\varepsilon_{i2}| \prec \frac{1}{N\eta}$ . This proves the first estimate in (B.1). The second estimate on  $\varepsilon_{i3}$  can be verified in the same way. We omit the details.

Now, we prove (B.2). According to (B.5), the left side of (B.2) is a sum of terms of the form

$$\frac{1}{N} \tilde{d}_i \hat{e}_i^* X G \hat{\boldsymbol{\alpha}}_i \hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G \hat{e}_i (\mathbf{k}_i^u)^* X_i G \hat{e}_i, \\ \frac{1}{N} \tilde{d}_i \hat{e}_i^* X G \mathcal{R}_i \tilde{B}^{(i)} \hat{\boldsymbol{\alpha}}_i \hat{\boldsymbol{\beta}}_i^* G \hat{e}_i (\mathbf{k}_i^u)^* X_i G \hat{e}_i.$$

Using (5.77), (B.6) and the bound  $|\hat{\mathbf{e}}_i^* X G \hat{\boldsymbol{\alpha}}_i| \prec \frac{1}{\eta}$  and  $|\hat{\mathbf{e}}_i^* X G \mathcal{R}_i \tilde{B}^{(i)} \hat{\boldsymbol{\alpha}}_i| \prec \frac{1}{\eta}$ , we can get (B.2). Then, (B.3) can be proved similarly to (B.2). Hence, we omit the details. Finally, we show (B.4). According to (B.5),  $\frac{1}{N} \sum_k^{(i)} \text{tr} Q X \Delta_G^u(i, k) \hat{\mathbf{e}}_k^* X_i G \hat{\mathbf{e}}_i$  is a sum of terms of the form

$$\begin{aligned} & \frac{1}{N^2} \tilde{d}_i \hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G Q X G \hat{\boldsymbol{\alpha}}_i (\hat{\mathbf{k}}_i^u)^* X_i G \hat{\mathbf{e}}_i, \\ & \frac{1}{N^2} \tilde{d}_i \hat{\boldsymbol{\beta}}_i^* G Q X G \mathcal{R}_i \tilde{B}^{(i)} \hat{\boldsymbol{\alpha}}_i (\hat{\mathbf{k}}_i^u)^* X_i G \hat{\mathbf{e}}_i. \end{aligned}$$

Then (B.4) follows from (5.77) and the trivial bounds

$$|\hat{\boldsymbol{\beta}}_i^* \tilde{B}^{(i)} \mathcal{R}_i G Q X G \hat{\boldsymbol{\alpha}}_i| \prec \frac{1}{\eta^2}, \quad |\hat{\boldsymbol{\beta}}_i^* G Q X G \mathcal{R}_i \tilde{B}^{(i)} \hat{\boldsymbol{\alpha}}_i| \prec \frac{1}{\eta^2}.$$

Hence, we concluded the proof of Lemma B.1.  $\square$

### APPENDIX C.

In this appendix, we explain how to modify our discussions in Sections 5 and 6 to adapt to the orthogonal setup. Recall our partial randomness decomposition of Haar unitary matrices  $U$  and  $V$  in (5.2). For Haar orthogonal matrices  $U$  and  $V$ , we refer to Appendix A of [3] for an analogous decomposition, with the phases of the  $i$ -th components of  $\mathbf{u}_i$  and  $\mathbf{v}_i$  replaced by the signs of them. We then inherit all the notations introduced in Sections 5 and 6. Under the orthogonal setting, instead of (5.44), we need to use the following integration by parts formula for real Gaussian random variables

$$\int_{\mathbb{R}} g f(g) e^{-\frac{g^2}{2\sigma^2}} dg = \sigma^2 \int_{\mathbb{R}} f'(g) e^{-\frac{g^2}{2\sigma^2}} dg,$$

for differentiable functions  $f : \mathbb{R} \rightarrow \mathbb{R}$ . Consequently, instead of (5.45), here we have, for  $k \neq i$ ,

$$\frac{\partial R_i^a}{\partial g_{ik}^a} = -\frac{(\ell_i^a)^2}{\|\mathbf{g}_i^a\|_2} \mathbf{e}_k (\mathbf{e}_i + \mathbf{h}_i^a)^* - \frac{(\ell_i^a)^2}{\|\mathbf{g}_i^a\|_2} (\mathbf{e}_i + \mathbf{h}_i^a) \mathbf{e}_k^* + 2\Delta_R^a(i, k), \quad a = u, v,$$

where  $\Delta_R^a(i, k)$  is defined in (5.46). Thence we have the following modification of (5.50):

$$\begin{aligned} \frac{\partial G}{\partial g_{ik}^u} &= \text{right side of (5.50)} + c_i^u G (\hat{\mathbf{e}}_i + \mathbf{k}_i^u) \hat{\mathbf{e}}_k^* \tilde{B}^{(i)} \mathcal{R}_i G \\ &\quad + c_i^u G \mathcal{R}_i \tilde{B}^{(i)} (\hat{\mathbf{e}}_i + \mathbf{k}_i^u) \hat{\mathbf{e}}_k^* G. \end{aligned} \quad (\text{C.1})$$

The remaining task is to go through all the discussions in Sections 5 and 6 again, and show that all the estimates which involve the last two terms in (C.1) are negligible at the right order.

To get through the discussions in Section 5 for orthogonal case, it suffices to take the last two terms in (C.1) into the account of the derivation of the equations (5.63) and (5.64), as well as the last two estimates in (5.75).

Using (C.1), we will have the following modification of (5.63):

$$\begin{aligned} \frac{1}{N} \sum_k^{(i)} \frac{\partial (\hat{\mathbf{e}}_k^* \tilde{B}^{(i)} G \hat{\mathbf{e}}_i)}{\partial g_{ik}^u} &= \text{right side of (5.63)} + \frac{c_i^u}{N} \hat{\mathbf{e}}_i^* G \mathcal{R}_i \tilde{B}^{(i)} \hat{I}_1^{(i)} \tilde{B}^{(i)} G (\hat{\mathbf{e}}_i + \mathbf{k}_i^u) \\ &\quad + \frac{c_i^u}{N} \hat{\mathbf{e}}_i^* G \hat{I}_1^{(i)} \tilde{B}^{(i)} G \mathcal{R}_i \tilde{B}^{(i)} (\hat{\mathbf{e}}_i + \mathbf{k}_i^u). \end{aligned} \quad (\text{C.2})$$

Notice that the new terms are qualitatively different from the ones already present in (5.63). In the new terms the summation over  $k$  could be directly performed since  $\mathbf{e}_k$  and  $\mathbf{e}_k^*$  appear directly next to each other, yielding the almost identity  $\hat{I}_1^{(i)}$ . The analogous sums in the old terms, explicitly seen in (5.58), result in a partial trace.

We will show that the last two terms above are of order  $O_{\prec}(\Psi^2)$ . For the first one, note that

$$|\hat{\mathbf{e}}_i^* G \mathcal{R}_i \tilde{B}^{(i)} \hat{I}_1^{(i)} \tilde{B}^{(i)} G (\hat{\mathbf{e}}_i + \mathbf{k}_i^u)| \leq C \|G \hat{\mathbf{e}}_i\|_2 (\|G \hat{\mathbf{e}}_i\|_2 + \|G \mathbf{k}_i^u\|_2), \quad (\text{C.3})$$

for some constant  $C$ . For the last term in (C.2), using (5.13), (5.18) and also  $\mathcal{R}_i^2 = \hat{I}$ , we get

$$\begin{aligned} G\mathcal{R}_i\tilde{B}^{(i)}(\hat{e}_i + \mathbf{k}_i^u) &= -\tilde{\sigma}_i^* G\mathbf{k}_i^u - G\tilde{B}\hat{e}_i = -\tilde{\sigma}_i^* G\mathbf{k}_i^u - \hat{e}_i + G(A - z)\hat{e}_i \\ &= \tilde{\sigma}_i^* G\mathbf{k}_i^u - \hat{e}_i - G\hat{e}_i + \xi_i^* G\hat{e}_i. \end{aligned} \quad (\text{C.4})$$

Thus, applying (C.4), for the last term in (C.2) we have

$$|\hat{e}_i^* G\hat{I}_1^{(i)}\tilde{B}^{(i)}G\mathcal{R}_i\tilde{B}^{(i)}(\hat{e}_i + \mathbf{k}_i^u)| \leq C\|G\hat{e}_i\|_2(\|G\hat{e}_i\|_2 + \|G\hat{e}_i\|_2 + \|G\mathbf{k}_i^u\|_2). \quad (\text{C.5})$$

According to (C.2), (C.3) and (C.5), it suffices to prove

$$\|G\hat{e}_i\|_2 \prec \frac{1}{\sqrt{\eta}}, \quad \|G\hat{e}_i\|_2 \prec \frac{1}{\sqrt{\eta}}, \quad \|G\mathbf{k}_i^u\|_2 \prec \frac{1}{\sqrt{\eta}} \quad (\text{C.6})$$

to get a bound  $O_{\prec}(\Psi^2)$  for the last two terms in (C.2). To show (C.6), we use the identities

$$\|G\hat{e}_i\|_2^2 = \frac{1}{\eta}\text{Im}G_{ii}, \quad \|G\hat{e}_i\|_2^2 = \frac{1}{\eta}\text{Im}G_{ii}, \quad \|G\mathbf{k}_i^u\|_2^2 = \frac{1}{\eta}\text{Im}(\mathbf{k}_i^u)^*G\mathbf{k}_i^u. \quad (\text{C.7})$$

Applying assumption (5.27) and (C.7), we can get the first two estimates in (C.6). Using the last identity of (C.7), (5.84) and (5.85), we can get the last estimate in (C.6). The necessary modification for the proofs of (5.64) and the last three estimates in (5.75) can be done in the same way, we thus omit the details.

For the discussions in Section 6, in the orthogonal case, the averaged analogue of (5.71), *i.e.*, (6.21), still holds. That is because the last two terms in (C.2) and their analog in the equation for  $\frac{1}{N}\sum_k^{(i)}\frac{\partial(\hat{e}_k^*G\hat{e}_i)}{\partial g_{kk}^u}$  are of order  $O_{\prec}(\Psi^2)$ . So the contribution of these additional terms in (6.21) can be absorbed into the last term of (6.21). Thence, the remaining proof is the same as the unitary case. Hence, we completed the necessary modifications for the orthogonal setup.

## REFERENCES

- [1] Anderson, G., Guionnet, A., Zeitouni, O.: *An Introduction to Random Matrices*, Cambridge Stud. Adv. Math. **118**, Cambridge Univ. Press, Cambridge, 2010.
- [2] Bao, Z. G., Erdős, L., Schnelli, K.: *Local stability of the free additive convolution*, J. Funct. Anal. **271(3)**, 672-719 (2016).
- [3] Bao, Z. G., Erdős, L., Schnelli, K.: *Local law of addition of random matrices on optimal scale*, Comm. Math. Phys. **349(3)**, 947-990 (2017).
- [4] Bao, Z. G., Erdős, L., Schnelli, K.: *Convergence rate for spectral distribution of addition of random matrices*, Adv. Math. **319**, 251-291 (2017).
- [5] Bao, Z. G., Erdős, L., Schnelli, K.: *Spectral rigidity for addition of random matrices at the regular edge*, arXiv:1708.01597, (2017).
- [6] Belinschi, S.: *A note on regularity for free convolutions*, Ann. Inst. Henri Poincaré Probab. Stat. **42(5)**, 635-648 (2006).
- [7] Belinschi, S.: *The Lebesgue decomposition of the free additive convolution of two probability distributions*, Probab. Theory Related Fields **142(1-2)**, 125-150 (2008).
- [8] Belinschi, S., Bercovici, H.: *A new approach to subordination results in free probability*, J. Anal. Math. **101(1)**, 357-365 (2007).
- [9] Benaych-Georges, F.: *Exponential bounds for the support convergence in the single ring theorem*, J. Funct. Anal. **268**, 3492-3507 (2015).
- [10] Benaych-Georges, F.: *Local single ring theorem*, arXiv:1501.07840, Ann. Probab. (appeared online).
- [11] Bercovici, H., Voiculescu, D.: *Free convolution of measures with unbounded support*, Indiana Univ. Math. J. **42**, 733-773 (1993).
- [12] Biane, P.: *Process with free increments*, Math. Z. **227(1)**, 143-174 (1998).
- [13] Bordenave, C., Chafaï, D.: *Around the circular law*, Probability Surveys **9**, 1-89 (2012).
- [14] Bourgade, P., Yau, H.-T., Yin, J.: *Local circular law for random matrices*, Probab. Theory Related Fields **159(3-4)** 545-595 (2014).
- [15] Bourgade, P., Yau, H.-T., Yin, J.: *The local circular law II: the edge case*, Probab. Theory Related Fields **159(3-4)** 619-660 (2014).
- [16] Chistyakov, G. P., Götze, F.: *The arithmetic of distributions in free probability theory*, Cent. Euro. J. Math. **9**, 997-1050 (2011).
- [17] Diaconis, P., Shahshahani, M.: *The subgroup algorithm for generating uniform random variables*, Probab. Engrg. Inform. Sci. **1(01)**, 15-32 (1987).



- [18] Erdős, L.: *Random matrices, log-gases and Hölder regularity*. Proceedings of ICM 2014, Seoul, Vol. III. 213-236 (2015).
- [19] Erdős, L., Yau, H.-T., Yin, J.: *Universality for generalized Wigner matrices with Bernoulli distribution*, J. Comb. **2(1)**, 15-85 (2011).
- [20] Erdős, L., Knowles, A., Yau, H.-T.: *Averaging fluctuations in resolvents of random band matrices*, Ann. Henri Poincaré **14**, 1837-1926 (2013).
- [21] Erdős, L., Schlein, B., Yau, H.-T.: *Local semicircle law and complete delocalization for Wigner random matrices*. Comm. Math. Phys. **287**, 641-655 (2009).
- [22] Erdős, L., Yau, H.-T., Yin, J.: *Bulk universality for generalized Wigner matrices*, Probab. Theory Related Fields **154(1-2)**: 341-407 (2012).
- [23] Feinberg, J., Zee, A.: *Non-Gaussian non-Hermitian random matrix theory: phase transition and addition formalism*, Nuclear Phys. B **501**, 643-669 (1997).
- [24] Girko, V. L.: *The circular law*, Teor. Veroyatnost. i Primenen. **29(4)**, 669-679 (1984).
- [25] Guionnet, A., Zeitouni, O.: *Support convergence in the single ring theorem*, Probab. Theory Related Fields **154** (3-4): 661-675 (2012).
- [26] Guionnet, A., Krishnapur, M., Zeitouni, O.: *The single ring theorem*, Ann. of Math. (2) **174**, 1189-1217 (2011).
- [27] Haagerup, U., Larsen, F.: *Brown's spectral distribution measure for  $R$ -diagonal elements in finite von Neumann algebras*, J. Funct. Anal. **176(2)**, 331-367 (2000).
- [28] Kargin, V.: *Subordination for the sum of two random matrices*, Ann. Probab. **43(4)**, 2119-2150 (2015).
- [29] Lee, J. O., Schnelli, K.: *Local law and Tracy-Widom limit for sparse random matrices*, arXiv:1605.08767 (2016).
- [30] Ledoux, M.: *The Concentration of Measure Phenomenon*, Providence, RI: American Mathematical Society, 2001.
- [31] Maassen, H.: *Addition of freely independent random variables*, J. Func. Anal. **106(2)**, 409-438 (2000).
- [32] Mezzadri, F.: *How to generate random matrices from the classical compact groups*, Notices Amer. Math. Soc. **54(5)**, 592-604 (2007).
- [33] Pastur, L., Vasilchuk, V.: *On the law of addition of random matrices*, Comm. Math. Phys. **214.2**, 249-286 (2000).
- [34] Rudelson, M., Vershynin, R.: *Invertibility of random matrices: unitary and orthogonal perturbations*, J. Amer. Math. Soc. **27(2)**, 293-338 (2014).
- [35] Tao, T., Vu, V.: *Random matrices: universality of ESDs and the circular law*, with an appendix by Krishnapur, M., Ann. Probab., **38(5)**, 2023-2065 (2010).
- [36] Tao, T., Vu, V.: *Random matrices: universality of local spectral statistics of non-Hermitian matrices*, Ann. Probab. **43(2)**, 782-874 (2015).
- [37] Voiculescu, D.: *The analogues of entropy and of Fisher's information theory in free probability theory, I*, Comm. Math. Phys. **155**, 71-92 (1993).
- [38] Yin, J.: *The local circular law III: general case*, Probab. Theory Related Fields **160(3-4)**, 679-732 (2014).