

471 A Proofs

472 A.1 Proof of Theorem 4.1

473 *Proof.* By the definition of $\langle \nu_{\text{DE}}, f_{\text{DE}}^{(1)}, f_{\text{DE}}^{(2)} \rangle$, we have:

$$\begin{aligned} \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y}) &= (2\pi)^{-\frac{d}{2}} D^2 \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\|\omega\|^2 + 2\omega^\top \mathbf{A}\omega \right. \\ &\quad \left. + \omega^\top (\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y}) + \mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + \mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y}\right) d\omega \\ &= (2\pi)^{-\frac{d}{2}} D^2 \exp\left(\mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + \mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y}\right) \\ &\quad \times \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\omega^\top (\mathbf{I}_d - 4\mathbf{A})\omega + \omega^\top (\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y})\right) d\omega. \end{aligned}$$

474 Since $8\mathbf{A} \prec \mathbf{I}_d$, we have $4\mathbf{A} \prec 0.5\mathbf{I}_d \prec \mathbf{I}_d$, meaning that $\mathbf{I}_d - 4\mathbf{A}$ is positive definite and invertible.
475 The following identity is straightforward to check:

$$\begin{aligned} -\frac{1}{2}\omega^\top (\mathbf{I}_d - 4\mathbf{A})\omega + \omega^\top (\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y}) &= -\frac{1}{2}(\omega - \mu)^\top \Sigma^{-1}(\omega - \mu) + \frac{1}{2}\mu^\top \Sigma^{-1}\mu, \\ \Sigma &= (\mathbf{I}_d - 4\mathbf{A})^{-1}, \quad \mu = \Sigma(\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y}). \end{aligned}$$

476 Therefore, we have:

$$\begin{aligned} \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y}) &= (2\pi)^{-d/2} D^2 \exp\left(\mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + \mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y} + \frac{1}{2}\mu^\top \Sigma^{-1}\mu\right) \\ &\quad \times \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}(\omega - \mu)^\top \Sigma^{-1}(\omega - \mu)\right) d\omega. \end{aligned}$$

477 Next, we use the fact that the integral of the multivariate Gaussian distribution with mean μ and
478 variance Σ is 1:

$$(2\pi)^{-d/2} \det(\Sigma)^{-1/2} \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}(\omega - \mu)^\top \Sigma^{-1}(\omega - \mu)\right) d\omega = 1.$$

479 From that we conclude:

$$\begin{aligned} \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y}) &= D^2 \det(\Sigma)^{1/2} \exp\left(\mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + \mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y} + \frac{1}{2}\mu^\top \Sigma^{-1}\mu\right) \\ &= D^2 \det(\mathbf{I}_d - 4\mathbf{A})^{-1/2} \exp\left(\mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + \mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y} \right. \\ &\quad \left. + \frac{1}{2}(\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y})^\top (\mathbf{I}_d - 4\mathbf{A})^{-1}(\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y})\right) \\ &= D^2 \det(\mathbf{I}_d - 4\mathbf{A})^{-1/2} \exp\left(\mathbf{x}^\top \left(\mathbf{C}^{(1)} + \frac{1}{2}(\mathbf{B}^{(1)})^\top (\mathbf{I}_d - 4\mathbf{A})^{-1}\mathbf{B}^{(1)}\right) \mathbf{x} \right. \\ &\quad \left. + \mathbf{y}^\top \left(\mathbf{C}^{(2)} + \frac{1}{2}(\mathbf{B}^{(2)})^\top (\mathbf{I}_d - 4\mathbf{A})^{-1}\mathbf{B}^{(2)}\right) \mathbf{y} + \mathbf{x}^\top (\mathbf{B}^{(1)})^\top (\mathbf{I}_d - 4\mathbf{A})^{-1}\mathbf{B}^{(2)}\mathbf{y}\right). \end{aligned}$$

480 Based on this expression, we conclude that, indeed, $\mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y}) = K^{(0)}(\mathbf{x}, \mathbf{y})$ for all
481 $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ if the conditions from theorem's statement are satisfied.

482 Next, we calculate expression for the variance. For any random variable Z , $\text{Var } Z = \mathbb{E}Z^2 - (\mathbb{E}Z)^2$.

483 In particular, if $Z = f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y})$, $\omega \sim \nu_{\text{DE}}$, we get:

$$\begin{aligned} \text{Var}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y}) &= \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x})^2 f_{\text{DE}}^{(2)}(\omega, \mathbf{y})^2 - \left(\mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x}) f_{\text{DE}}^{(2)}(\omega, \mathbf{y})\right)^2 \\ &= \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\omega, \mathbf{x})^2 f_{\text{DE}}^{(2)}(\omega, \mathbf{y})^2 - K^{(0)}(\mathbf{x}, \mathbf{y})^2. \end{aligned}$$

484 We have:

$$\begin{aligned} \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x})^2 f_{\text{DE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y})^2 &= (2\pi)^{\frac{d}{2}} D^4 \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\|\boldsymbol{\omega}\|^2 + 4\boldsymbol{\omega}^\top \mathbf{A}\boldsymbol{\omega} + 2\boldsymbol{\omega}^\top (\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y})\right. \\ &\quad \left.+ 2\mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + 2\mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y}\right) d\boldsymbol{\omega} = (2\pi)^{\frac{d}{2}} D^4 \exp\left(2\mathbf{x}^\top \mathbf{C}^{(1)}\mathbf{x} + 2\mathbf{y}^\top \mathbf{C}^{(2)}\mathbf{y}\right) \\ &\quad \times \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2}\boldsymbol{\omega}^\top (\mathbf{I}_d - 8\mathbf{A})\boldsymbol{\omega} + 2\boldsymbol{\omega}^\top (\mathbf{B}^{(1)}\mathbf{x} + \mathbf{B}^{(2)}\mathbf{y})\right) d\boldsymbol{\omega}. \end{aligned}$$

485 Evaluation of the integral above can be done in the same way as calculation of
 486 $\mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}) f_{\text{DE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y})$, noticing that $\mathbf{I}_d - 8\mathbf{A}$ is positive definite and invertible. The result
 487 is as follows:

$$\begin{aligned} \mathbb{E}_{\nu_{\text{DE}}} f_{\text{DE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x})^2 f_{\text{DE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y})^2 &= D^4 \det(\mathbf{I}_d - 8\mathbf{A})^{-1/2} \\ &\quad \times \exp\left(2\mathbf{x}^\top \left(\mathbf{C}^{(1)} + (\mathbf{B}^{(1)})^\top (\mathbf{I}_d - 8\mathbf{A})^{-1} \mathbf{B}^{(1)}\right) \mathbf{x}\right. \\ &\quad \left.+ 2\mathbf{y}^\top \left(\mathbf{C}^{(2)} + (\mathbf{B}^{(2)})^\top (\mathbf{I}_d - 8\mathbf{A})^{-1} \mathbf{B}^{(2)}\right) \mathbf{y} + 4\mathbf{x}^\top (\mathbf{B}^{(1)})^\top (\mathbf{I}_d - 8\mathbf{A})^{-1} \mathbf{B}^{(2)} \mathbf{y}\right). \end{aligned}$$

488 We conclude that the variance expression given in the theorem's statement is correct. \square

489 A.2 Important lemma

490 Below, we prove an important lemma which is used in the subsequent proofs:

491 **Lemma A.1.** Consider a function $f : (-\infty, \frac{1}{8})$ defined as

$$f(A) = \log(1 - 4A) - \frac{1}{2} \log(1 - 8A) + \frac{\phi}{1 - 8A} \quad (14)$$

492 where $\phi \geq 0$. Then, the minimum of f on $(-\infty, \frac{1}{8})$ is achieved at

$$A^* = \frac{1}{16} \left(1 - 2\phi - \sqrt{(2\phi + 1)^2 + 8\phi}\right). \quad (15)$$

493 *Proof.* Set $\gamma = (1 - 8A)^{-1} \in (0, +\infty)$. Note that there is a one-to-one correspondence between
 494 $\gamma \in (0, +\infty)$ and $A \in (-\infty, \frac{1}{8})$. Hence, we can substitute $\gamma^{-1} = 1 - 8A$ and $1 - 4A =$
 495 $((1 - 8A) + 1)/2 = (\gamma^{-1} + 1)/2 = \frac{1+\gamma}{2\gamma}$ in (14) and equivalently perform minimization with respect
 496 to γ :

$$\min_{\gamma \in (0, +\infty)} h(\gamma) = \log\left(\frac{\gamma + 1}{2\gamma}\right) + \frac{1}{2} \log \gamma + \phi \gamma = \log(\gamma + 1) - \frac{1}{2} \log \gamma - \log 2 + \phi \gamma.$$

497 For $h(\cdot)$'s derivative, we have:

$$h'(\gamma) = \phi + \frac{1}{\gamma + 1} - \frac{1}{2\gamma} = \phi + \frac{\gamma - 1}{2\gamma(\gamma + 1)} \quad (16)$$

$$= \frac{2\phi\gamma(\gamma + 1) + \gamma - 1}{2\gamma(\gamma + 1)} = \frac{2\phi\gamma^2 + (2\phi + 1)\gamma - 1}{2\gamma(\gamma + 1)}. \quad (17)$$

498 Based on (16), we see that $h'(\gamma) \rightarrow -\infty$ as $\gamma \rightarrow 0$ and $h'(\gamma) > \phi \geq 0$ for all $\gamma > 1$. Hence,
 499 we conclude that $h(\cdot)$ is bounded from below on $(0, +\infty)$ and the global minimum γ^* on $(0, +\infty)$
 500 exists and it is one of the points satisfying $h'(\gamma^*) = 0$. Hence, it's one of the positive roots of the
 501 polynomial in numerator of (17).

502 If $\phi = 0$, there is a single root $\gamma^* = 1$ of the polynomial in the numerator of (17), hence it is a global
 503 minimum of $h(\cdot)$. If $\phi > 0$, then there are two roots of the polynomial in the numerator of (17):

$$\gamma_-^* = \frac{-(2\phi + 1) - \sqrt{(2\phi + 1)^2 + 8\phi}}{4\phi},$$

$$\gamma_+^* = \frac{-(2\phi + 1) + \sqrt{(2\phi + 1)^2 + 8\phi}}{4\phi}. \quad (18)$$

504 Note that, if $\phi > 0$, then $2\phi + 1 > 0$ and $(2\phi + 1)^2 + 8\phi \geq (2\phi + 1)^2$. Hence, $\gamma_-^* < 0$ and
 505 $\gamma_+^* > 0$. We conclude that $\gamma^* = \gamma_+^*$ is the minimum of $h(\cdot)$ on $(0, +\infty)$. We multiply numerator and
 506 denominator of (18)'s right hand side by $(2\phi + 1) + \sqrt{(2\phi + 1)^2 + 8\phi} > 0$:

$$\gamma^* = \gamma_+^* = \frac{((2\phi + 1)^2 + 8\phi) - (2\phi + 1)^2}{4\phi \left((2\phi + 1) + \sqrt{(2\phi + 1)^2 + 8\phi} \right)} = \frac{2}{2\phi + 1 + \sqrt{(2\phi + 1)^2 + 8\phi}}. \quad (19)$$

507 Note that the right hand side of (19) is equivalent to (18) when $\phi > 0$ but also holds for the case
 508 when $\phi = 0$ (i.e. when $\gamma^* = 1$). We conclude that $f(\cdot)$ is minimized at $\mathbf{A}^* = \frac{1}{8}(1 - (\gamma^*)^{-1})$ since
 509 $\gamma^* = (1 - 8A^*)^{-1}$. It's easy to see that (15) follows from (19) directly. \square

510 A.3 Proof of Theorem 4.2

511 *Proof.* With $\mathbf{A} = A\mathbf{I}_d$, the conditions from Theorem 4.1 read as

$$8A < 1, \quad \frac{1}{1 - 4A}(\mathbf{B}^{(1)})^\top \mathbf{B}^{(2)} = \mathbf{I}_d, \quad \mathbf{C}^{(k)} = -\frac{1}{2(1 - 4A)}(\mathbf{B}^{(k)})^\top \mathbf{B}^{(k)}, \quad D = (1 - 4A)^{d/4} \quad (20)$$

512 for $k \in \{1, 2\}$. And the variance expression (9) for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ transforms into

$$\begin{aligned} \text{Var}_{\nu_{\text{ADE}}} f_{\text{ADE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}) f_{\text{ADE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y}) &= D^4 (1 - 8A)^{-d/2} \exp \left(2\mathbf{x}^\top \left(\mathbf{C}^{(1)} + \frac{1}{1 - 8A}(\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)} \right) \mathbf{x} \right. \\ &\quad \left. + 2\mathbf{y}^\top \left(\mathbf{C}^{(2)} + \frac{1}{1 - 8A}(\mathbf{B}^{(2)})^\top \mathbf{B}^{(2)} \right) \mathbf{y} + \frac{4}{1 - 8A} \mathbf{x}^\top (\mathbf{B}^{(1)})^\top \mathbf{B}^{(2)} \mathbf{y} \right) - K^{(0)}(\mathbf{x}, \mathbf{y})^2. \end{aligned}$$

513 We express $\mathbf{C}^{(k)}$ through $A, \mathbf{B}^{(k)}$ and D through A using (20) in the equation above:

$$\begin{aligned} \text{Var}_{\nu_{\text{ADE}}} f_{\text{ADE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}) f_{\text{ADE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y}) &= \left(\frac{1 - 4A}{\sqrt{1 - 8A}} \right)^d \exp \left(\left(\frac{2}{1 - 8A} - \frac{1}{1 - 4A} \right) \mathbf{x}^\top (\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)} \mathbf{x} \right. \\ &\quad \left. + \left(\frac{2}{1 - 8A} - \frac{1}{1 - 4A} \right) \mathbf{y}^\top (\mathbf{B}^{(2)})^\top \mathbf{B}^{(2)} \mathbf{y} + \frac{4}{1 - 8A} \mathbf{x}^\top (\mathbf{B}^{(1)})^\top \mathbf{B}^{(2)} \mathbf{y} \right) - K^{(0)}(\mathbf{x}, \mathbf{y})^2. \end{aligned}$$

514 Since $\frac{1}{1 - 4A}(\mathbf{B}^{(1)})^\top \mathbf{B}^{(2)}$ is a full-rank matrix \mathbf{I}_d (20), both $\mathbf{B}^{(1)}$ and $\mathbf{B}^{(2)}$ are full-rank. Hence, we
 515 can express $\mathbf{B}^{(2)} = (1 - 4A)(\mathbf{B}^{(1)})^{-\top}$. Also, note that

$$\frac{2}{1 - 8A} - \frac{1}{1 - 4A} = \frac{2 - 8A - 1 + 8A}{(1 - 8A)(1 - 4A)} = (1 - 8A)^{-1}(1 - 4A)^{-1}.$$

516 We rewrite the expression for the variance using the identity above and the formula for $\mathbf{B}^{(2)}$:

$$\begin{aligned} \text{Var}_{\nu_{\text{ADE}}} f_{\text{ADE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}) f_{\text{ADE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y}) &= \left(\frac{1 - 4A}{\sqrt{1 - 8A}} \right)^d \exp \left((1 - 8A)^{-1}(1 - 4A)^{-1} \mathbf{x}^\top (\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)} \mathbf{x} \right. \\ &\quad \left. + (1 - 8A)^{-1}(1 - 4A) \mathbf{y}^\top ((\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)})^{-1} \mathbf{y} + 4(1 - 8A)^{-1}(1 - 4A) \mathbf{x}^\top \mathbf{y} \right) - K^{(0)}(\mathbf{x}, \mathbf{y})^2. \end{aligned}$$

517 We use the expression above to rewrite (8) for $\langle \nu, f^{(1)}, f^{(2)} \rangle = \langle \nu_{\text{ADE}}, f_{\text{ADE}}^{(1)}, f_{\text{ADE}}^{(2)} \rangle$ as follows:

$$\begin{aligned} \bar{\mathcal{L}}(\boldsymbol{\theta}_{\text{ADE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{ADE}}) &= L^{-2} \sum_{1 \leq i, j \leq L} \log(\text{Var}_{\nu_{\text{ADE}}} f_{\text{ADE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}^{(i)}) f_{\text{ADE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y}^{(j)}) + K^{(0)}(\mathbf{x}^{(i)}, \mathbf{y}^{(j)})) \\ &= d \log(1 - 4A) - \frac{d}{2} \log(1 - 8A) + (1 - 8A)^{-1}(1 - 4A)^{-1} L^{-1} \sum_{i=1}^L (\mathbf{x}^{(i)})^\top (\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)} \mathbf{x}^{(i)} \\ &\quad + (1 - 8A)^{-1}(1 - 4A) L^{-1} \sum_{j=1}^L (\mathbf{y}^{(j)})^\top (\mathbf{B}^{(1)})^{-1} (\mathbf{B}^{(1)})^{-\top} \mathbf{y}^{(j)} \end{aligned}$$

$$+4(1-8A)^{-1}(1-4A)L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)})^\top \mathbf{y}^{(j)}. \quad (21)$$

518 Denote $\mathbf{E} = (\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)}$. Then (21) becomes:

$$\begin{aligned} \bar{\mathcal{L}}(\boldsymbol{\theta}_{\text{ADE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{ADE}}) &= d \log(1-4A) - \frac{d}{2} \log(1-8A) \\ &\quad + (1-8A)^{-1}(1-4A)^{-1} L^{-1} \sum_{i=1}^L (\mathbf{x}^{(i)})^\top \mathbf{E} \mathbf{x}^{(i)} \\ &\quad + (1-8A)^{-1}(1-4A)L^{-1} \sum_{j=1}^L (\mathbf{y}^{(j)})^\top \mathbf{E}^{-1} \mathbf{y}^{(j)} + 4(1-8A)^{-1}(1-4A)L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)})^\top \mathbf{y}^{(j)}. \end{aligned} \quad (22)$$

519 We next prove the following lemma:

520 **Lemma A.2.** Let $\mathbf{B}^{(1)*} = \sqrt{1-4A} \boldsymbol{\Sigma}^{1/2} \mathbf{U}^\top (\boldsymbol{\Lambda}^{(1)})^{-1/2} (\mathbf{Q}^{(1)})^\top$. When A ($8A < 1$) is fixed,
521 $\mathbf{E} = \mathbf{E}^* = (\mathbf{B}^{(1)*})^\top \mathbf{B}^{(1)*}$ minimizes the right hand side of (22) with respect to \mathbf{E} .

522 *Proof.* We have:

$$\begin{aligned} L^{-1} \sum_{i=1}^L (\mathbf{x}^{(i)})^\top \mathbf{E} \mathbf{x}^{(i)} &= L^{-1} \sum_{i=1}^L \text{Trace}((\mathbf{x}^{(i)})^\top \mathbf{E} \mathbf{x}^{(i)}) = L^{-1} \sum_{i=1}^L \text{Trace}(\mathbf{E} \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top) \\ &= \text{Trace} \left(\mathbf{E} \left(L^{-1} \sum_{i=1}^L \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top \right) \right) = \text{Trace}(\mathbf{E} \mathbf{M}^{(1)}) \end{aligned}$$

523 where we use the cyclic property of trace $\text{Trace}(\cdot)$ and linearity of trace. Analogously, we obtain
524 $L^{-1} \sum_{j=1}^L (\mathbf{y}^{(j)})^\top \mathbf{E}^{-1} \mathbf{y}^{(j)} = \text{Trace}(\mathbf{E}^{-1} \mathbf{M}^{(2)})$. Assuming that A is fixed, optimization of (22)
525 with respect to \mathbf{E} reduces to the following minimization problem:

$$\min_{\mathbf{E} \in \mathbb{S}_d, \mathbf{E} \succ 0} \mathcal{F}(\mathbf{E}) = \beta_1 \text{Trace}(\mathbf{E} \mathbf{M}^{(1)}) + \beta_2 \text{Trace}(\mathbf{E}^{-1} \mathbf{M}^{(2)}) \quad (23)$$

526 where $\beta_1 = (1-8A)^{-1}(1-4A)^{-1}$, $\beta_2 = (1-8A)^{-1}(1-4A)$ and the constraint $\mathbf{E} \in \mathbb{S}_d, \mathbf{E} \succ 0$
527 follows from the fact that $\mathbf{E} = (\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)}$ and \mathbf{E} is invertible. We have $1-8A > 0$ and
528 $1-4A = (1-8A)/2 + 1/2 > 0$. Hence, $\beta_1, \beta_2 > 0$. For any $\mathbf{E} \succ 0$ and any $\boldsymbol{\Delta} \in \mathbb{S}_d$ there is $t \in \mathbb{R}$
529 small enough such that $\mathbf{E} + t\boldsymbol{\Delta}$ is invertible and the following Neumann series is convergent:

$$(\mathbf{E} + t\boldsymbol{\Delta})^{-1} = \mathbf{E}^{-1} (\mathbf{I}_d + t\boldsymbol{\Delta} \mathbf{E}^{-1})^{-1} = \sum_{l=0}^{\infty} (-t)^l \mathbf{E}^{-1} (\boldsymbol{\Delta} \mathbf{E}^{-1})^l$$

530 We further deduce:

$$\begin{aligned} \text{Trace}((\mathbf{E} + t\boldsymbol{\Delta})^{-1} \mathbf{M}^{(2)}) &= \text{Trace} \left(\left(\sum_{l=0}^{\infty} (-t)^l \mathbf{E}^{-1} (\boldsymbol{\Delta} \mathbf{E}^{-1})^l \right) \mathbf{M}^{(2)} \right) = \\ &= \sum_{l=0}^{\infty} (-t)^l \text{Trace} \left(\mathbf{E}^{-1} (\boldsymbol{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right) \end{aligned}$$

531 and, therefore,

$$\begin{aligned} \mathcal{F}(\mathbf{E} + t\boldsymbol{\Delta}) &= \beta_1 \text{Trace}((\mathbf{E} + t\boldsymbol{\Delta}) \mathbf{M}^{(1)}) + \beta_2 \sum_{l=0}^{\infty} (-t)^l \text{Trace} \left(\mathbf{E}^{-1} (\boldsymbol{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right) \\ &= \beta_1 \text{Trace}(\mathbf{E} \mathbf{M}^{(1)}) + t\beta_1 \text{Trace}(\boldsymbol{\Delta} \mathbf{M}^{(1)}) + \beta_2 \sum_{l=0}^{\infty} (-t)^l \text{Trace} \left(\mathbf{E}^{-1} (\boldsymbol{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right). \end{aligned} \quad (24)$$

Further, we have:

$$\begin{aligned}\frac{\partial}{\partial t}\mathcal{F}(\mathbf{E} + t\mathbf{\Delta}) &= \beta_1 \text{Trace}(\mathbf{\Delta}\mathbf{M}^{(1)}) + \beta_2 \sum_{l=1}^{\infty} (-1)^l l t^{l-1} \text{Trace} \left(\mathbf{E}^{-1} (\mathbf{\Delta}\mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right), \\ \frac{\partial^2}{(\partial t)^2}\mathcal{F}(\mathbf{E} + t\mathbf{\Delta}) &= \beta_2 \sum_{l=2}^{\infty} (-1)^l l(l-1) t^{l-2} \text{Trace} \left(\mathbf{E}^{-1} (\mathbf{\Delta}\mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right), \\ \left. \frac{\partial^2}{(\partial t)^2}\mathcal{F}(\mathbf{E} + t\mathbf{\Delta}) \right|_{t=0} &= 2\beta_2 \text{Trace} \left(\mathbf{E}^{-1} \mathbf{\Delta} \mathbf{E}^{-1} \mathbf{\Delta} \mathbf{E}^{-1} \mathbf{M}^{(2)} \right).\end{aligned}\quad (25)$$

We replace $\mathbf{M}^{(2)} = \mathbf{Q}^{(2)}(\mathbf{\Lambda}^{(2)})^{1/2}(\mathbf{\Lambda}^{(2)})^{1/2}(\mathbf{Q}^{(2)})^\top$ and apply the cyclic property of trace in (25):

$$\begin{aligned}\left. \frac{\partial^2}{(\partial t)^2}\mathcal{F}(\mathbf{E} + t\mathbf{\Delta}) \right|_{t=0} &= 2\beta_2 \text{Trace} \left((\mathbf{\Lambda}^{(2)})^{1/2} (\mathbf{Q}^{(2)})^\top \mathbf{E}^{-1} \mathbf{\Delta} \mathbf{E}^{-1} \mathbf{\Delta} \mathbf{E}^{-1} \mathbf{Q}^{(2)} (\mathbf{\Lambda}^{(2)})^{1/2} \right) \\ &= 2\beta_2 \text{Trace} (\mathbf{T} \mathbf{E}^{-1} \mathbf{T}^\top)\end{aligned}$$

where $\mathbf{T} = (\mathbf{\Lambda}^{(2)})^{1/2} (\mathbf{Q}^{(2)})^\top \mathbf{E}^{-1} \mathbf{\Delta}$. Since \mathbf{E} is positive definite, \mathbf{E}^{-1} is also positive definite and $\mathbf{T} \mathbf{E}^{-1} \mathbf{T}^\top$ is at least positive semidefinite. Hence, $\text{Trace} (\mathbf{T} \mathbf{E}^{-1} \mathbf{T}^\top) \geq 0$ and also $\left. \frac{\partial^2}{(\partial t)^2}\mathcal{F}(\mathbf{E} + t\mathbf{\Delta}) \right|_{t=0} \geq 0$. We conclude that $\mathcal{F}(\mathbf{E})$ is a convex function on $\{\mathbf{E} \in \mathbb{S}_d \mid \mathbf{E} \succ 0\}$. Since $\{\mathbf{E} \in \mathbb{S}_d \mid \mathbf{E} \succ 0\}$ is an open set, (every) global minimum \mathbf{E} of (23) satisfies two conditions

$$1) \mathbf{E} \succ 0, \quad \text{and} \quad 2) \nabla \mathcal{F}(\mathbf{E}) = \mathbf{0}_{d \times d} \quad (26)$$

Set $t = 1$ and assume that $\mathbf{\Delta} \in \mathbb{S}_d$ is small enough by norm so that $\mathbf{E} + \mathbf{\Delta}$ is invertible and the Neumann series for $(\mathbf{I}_d + \mathbf{\Delta} \mathbf{E}^{-1})^{-1}$ is convergent. Then, (24) holds for $t = 1$:

$$\begin{aligned}\mathcal{F}(\mathbf{E} + \mathbf{\Delta}) &= \beta_1 \text{Trace}(\mathbf{E} \mathbf{M}^{(1)}) + \beta_1 \text{Trace}(\mathbf{\Delta} \mathbf{M}^{(1)}) + \beta_2 \sum_{l=0}^{\infty} (-1)^l \text{Trace} \left(\mathbf{E}^{-1} (\mathbf{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right) \\ &= \mathcal{F}(\mathbf{E}) + \beta_1 \text{Trace}(\mathbf{\Delta} \mathbf{M}^{(1)}) + \beta_2 \sum_{l=1}^{\infty} (-1)^l \text{Trace} \left(\mathbf{E}^{-1} (\mathbf{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right) \\ &= \mathcal{F}(\mathbf{E}) + \beta_1 \text{Trace}(\mathbf{\Delta} \mathbf{M}^{(1)}) - \beta_2 \text{Trace} \left(\mathbf{E}^{-1} \mathbf{\Delta} \mathbf{E}^{-1} \mathbf{M}^{(2)} \right) \\ &\quad + \beta_2 \sum_{l=2}^{\infty} (-1)^l \text{Trace} \left(\mathbf{E}^{-1} (\mathbf{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right).\end{aligned}$$

Clearly, $\beta_2 \sum_{l=2}^{\infty} (-1)^l \text{Trace} \left(\mathbf{E}^{-1} (\mathbf{\Delta} \mathbf{E}^{-1})^l \mathbf{M}^{(2)} \right) = o(\|\mathbf{\Delta}\|)$ where $\|\cdot\|$ is an L_2 -norm. Also, using the cyclic property of trace, we get:

$$\text{Trace} \left(\mathbf{E}^{-1} \mathbf{\Delta} \mathbf{E}^{-1} \mathbf{M}^{(2)} \right) = \text{Trace} \left(\mathbf{\Delta} \mathbf{E}^{-1} \mathbf{M}^{(2)} \mathbf{E}^{-1} \right).$$

Therefore, we have:

$$\mathcal{F}(\mathbf{E} + \mathbf{\Delta}) = \text{Trace} \left(\mathbf{\Delta} \left(\beta_1 \mathbf{M}^{(1)} - \beta_2 \mathbf{E}^{-1} \mathbf{M}^{(2)} \mathbf{E}^{-1} \right) \right) + o(\|\mathbf{\Delta}\|). \quad (27)$$

Since $\mathbf{\Delta}, \mathbf{E}^{-1}, \mathbf{M}^{(1)}, \mathbf{M}^{(2)} \in \mathbb{S}_d$, from (27) it follows that

$$\nabla \mathcal{F}(\mathbf{E}) = \beta_1 \mathbf{M}^{(1)} - \beta_2 \mathbf{E}^{-1} \mathbf{M}^{(2)} \mathbf{E}^{-1}. \quad (28)$$

Let $\mathbf{E}^* = (\mathbf{B}^{(1)*})^\top \mathbf{B}^{(1)*} \succeq 0$. Note that

$$\sqrt[4]{\frac{\beta_2}{\beta_1}} = \sqrt[4]{\frac{(1-8A)^{-1}(1-4A)}{(1-8A)^{-1}(1-4A)^{-1}}} = \sqrt{1-4A}.$$

Since $\sqrt{\beta_2/\beta_1} \neq 0$, $\mathbf{\Sigma}, \mathbf{U}, \mathbf{\Lambda}^{-1/2}, \mathbf{Q}^{(1)}$ are full-rank, \mathbf{E}^* is also full-rank, therefore $\mathbf{E}^* \succ 0$ and it satisfies condition 1 from (26). Observe that

$$\mathbf{E}^* \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{1/2} = \sqrt{\beta_2/\beta_1} \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top (\mathbf{\Lambda}^{(1)})^{-1/2} (\mathbf{Q}^{(1)})^\top \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{1/2}$$

$$\begin{aligned}
&= \sqrt{\beta_2/\beta_1} \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top \\
&= \sqrt{\beta_2/\beta_1} \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top) \mathbf{V} \mathbf{U}^\top \\
&= \sqrt{\beta_2/\beta_1} \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} ((\mathbf{\Lambda}^{(1)})^{\frac{1}{2}} (\mathbf{Q}^{(1)})^\top \mathbf{Q}^{(2)} (\mathbf{\Lambda}^{(2)})^{\frac{1}{2}}) \mathbf{V} \mathbf{U}^\top \\
&= \sqrt{\beta_2/\beta_1} \mathbf{Q}^{(2)} (\mathbf{\Lambda}^{(2)})^{1/2} \mathbf{V} \mathbf{U}^\top
\end{aligned}$$

547 where we use definitions of \mathbf{E}^* , \mathbf{U} , $\mathbf{\Sigma}$, \mathbf{V} and orthogonality of $\mathbf{Q}^{(1)}$, $\mathbf{Q}^{(2)}$, \mathbf{U} , \mathbf{V} . Hence, we deduce
548 that

$$\beta_1 \mathbf{E}^* \mathbf{M}^{(1)} \mathbf{E}^* = \beta_1 \mathbf{E}^* \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{1/2} \left((\mathbf{\Lambda}^{(1)})^{1/2} (\mathbf{Q}^{(1)})^\top \mathbf{E}^* \right) \quad (29)$$

$$\begin{aligned}
&= \beta_1 \frac{\beta_2}{\beta_1} \mathbf{Q}^{(2)} (\mathbf{\Lambda}^{(2)})^{1/2} \left((\mathbf{\Lambda}^{(2)})^{1/2} (\mathbf{Q}^{(2)})^\top \right) \\
&= \beta_2 \mathbf{M}^{(2)} \quad (30)
\end{aligned}$$

549 by the definition of $\mathbf{Q}^{(1)}$, $\mathbf{\Lambda}^{(1)}$, $\mathbf{Q}^{(2)}$, $\mathbf{\Lambda}^{(2)}$ and due to orthogonality of \mathbf{V} , \mathbf{U} . By left- and right-
550 multiplication of (30) by $(\mathbf{E}^*)^{-1}$ we deduce that

$$\beta_1 \mathbf{M}^{(1)} = \beta_2 (\mathbf{E}^*)^{-1} \mathbf{M}^{(2)} (\mathbf{E}^*)^{-1}$$

551 or, in other words, $\nabla \mathcal{F}(\mathbf{E}^*) = \mathbf{0}_{d \times d}$ and the condition 2 from (26) is also satisfied. We conclude that
552 the global minimum of (23) is achieved at \mathbf{E}^* . \square

553 According to Lemma A.2, $\mathbf{B}^{(1)} = \mathbf{B}^{(1)*}$ is a global minimum of (21)'s right hand side when A is
554 fixed. Indeed, if there is $\mathbf{B}^{(1)}$ which leads to a smaller value of (21), $\mathbf{E} = (\mathbf{B}^{(1)})^\top \mathbf{B}^{(1)}$ would lead
555 to a smaller value of (22)'s right hand side. Also, this \mathbf{E} is positive definite by definition (note that
556 $\mathbf{B}^{(1)}$ is nonsingular), leading to contradiction with Lemma A.2.

557 Substituting \mathbf{E}^* instead of \mathbf{E} in (22) corresponds to the minimum value of $\bar{\mathcal{L}}(\boldsymbol{\theta}_{\text{AGE}}; \alpha, \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{AGE}})$
558 for a fixed A . Our next step is to minimize this expression with respect to A . Denote $\mathbf{F} =$
559 $\mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top (\mathbf{\Lambda}^{(1)})^{-1/2} (\mathbf{Q}^{(1)})^\top$. Then $\mathbf{E}^* = (1 - 4A) \mathbf{F}$ where \mathbf{F} doesn't depend on A .
560 We substitute \mathbf{E}^* into (22) and get:

$$\begin{aligned}
&d \log(1 - 4A) - \frac{d}{2} \log(1 - 8A) + (1 - 8A)^{-1} (1 - 4A)^{-1} \text{Trace}((1 - 4A) \mathbf{F} \mathbf{M}^{(1)}) \\
&\quad + (1 - 8A)^{-1} (1 - 4A) \text{Trace}((1 - 4A)^{-1} \mathbf{F}^{-1} \mathbf{M}^{(2)}) \\
&\quad + 4(1 - 8A)^{-1} (1 - 4A) L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)})^\top \mathbf{y}^{(j)} \\
&= d \log(1 - 4A) - \frac{d}{2} \log(1 - 8A) + (1 - 8A)^{-1} \text{Trace}(\mathbf{F} \mathbf{M}^{(1)}) + (1 - 8A)^{-1} \text{Trace}(\mathbf{F}^{-1} \mathbf{M}^{(2)}) \\
&\quad + 2(1 + (1 - 8A)^{-1}) d\mu^{(3)} \quad (31)
\end{aligned}$$

561 where we also replace

$$L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)})^\top \mathbf{y}^{(j)} = L^{-2} \left(\sum_{i=1}^L \mathbf{x}^{(i)} \right)^\top \left(\sum_{j=1}^L \mathbf{y}^{(j)} \right) = d\mu^{(3)}$$

562 and

$$(1 - 8A)^{-1} (1 - 4A) = \frac{(1 - 8A) + 1}{2(1 - 8A)} = \frac{1}{2} (1 + (1 - 8A)^{-1})$$

563 Based on (30) and since $\mathbf{F} = \sqrt{\beta_1/\beta_2} \mathbf{E}^*$, we conclude that $\mathbf{F} \mathbf{M}^{(1)} \mathbf{F} = \mathbf{M}^{(2)}$, or $\mathbf{M}^{(1)} \mathbf{F} =$
564 $\mathbf{F}^{-1} \mathbf{M}^{(2)}$. Using the cyclic property of trace, we get:

$$\text{Trace}(\mathbf{F} \mathbf{M}^{(1)}) = \text{Trace}(\mathbf{M}^{(1)} \mathbf{F}) = \text{Trace}(\mathbf{F}^{-1} \mathbf{M}^{(2)}).$$

565 By the definition of \mathbf{F} , $\mathbf{\Lambda}^{(1)}$, $\mathbf{Q}^{(1)}$ and using the cyclic property and orthogonality of $\mathbf{Q}^{(1)}$, \mathbf{U} , we
566 have:

$$\text{Trace}(\mathbf{F} \mathbf{M}^{(1)}) = \text{Trace} \left(\mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{Q}^{(1)} \right)^\top \left(\mathbf{Q}^{(1)} \mathbf{\Lambda}^{(1)} (\mathbf{Q}^{(1)})^\top \right)$$

$$\begin{aligned}
&= \text{Trace} \left(\mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top (\mathbf{\Lambda}^{(1)})^{1/2} (\mathbf{Q}^{(1)})^\top \right) \\
&= \text{Trace} \left(\mathbf{\Sigma} \mathbf{U}^\top (\mathbf{\Lambda}^{(1)})^{1/2} (\mathbf{Q}^{(1)})^\top \mathbf{Q}^{(1)} (\mathbf{\Lambda}^{(1)})^{-1/2} \mathbf{U} \right) \\
&= \text{Trace}(\mathbf{\Sigma}) = \sum_{l=1}^d \Sigma_{l,l}.
\end{aligned}$$

Hence, (31) finally becomes:

$$\begin{aligned}
&d \log(1 - 4A) - \frac{d}{2} \log(1 - 8A) + 2(1 - 8A)^{-1} \sum_{l=1}^d \Sigma_{l,l} + 2(1 + (1 - 8A)^{-1}) d\mu^{(3)} \\
&= d \left(\log(1 - 4A) - \frac{1}{2} \log(1 - 8A) + 2(1 - 8A)^{-1} \left(d^{-1} \sum_{l=1}^d \Sigma_{l,l} + \mu^{(3)} \right) + 2\mu^{(3)} \right). \quad (32)
\end{aligned}$$

Next, we use Lemma A.1 ($\phi = d^{-1} \sum_{l=1}^d \Sigma_{l,l} + \mu^{(3)} \geq 0$) for deriving expression for A which minimizes (32). This expression coincides with the one in Theorem's statement. The expressions for $\mathbf{B}^{(2)}$, $\mathbf{C}^{(1)}$, $\mathbf{C}^{(2)}$ follow directly from (20), optimal $\mathbf{B}^{(1)} = \mathbf{B}^{(1)*}$ and A . (10) follows from (32). The proof is concluded. \square

A.4 Proof of Theorem 4.3

Proof. With $\mathbf{B}^{(1)} = \mathbf{B}^{(2)} = \mathbf{B}$ and $\mathbf{C}^{(1)} = \mathbf{C}^{(2)} = \mathbf{C}$, the conditions from Theorem 4.1 read as

$$8\mathbf{A} \prec \mathbf{I}_d, \quad \mathbf{B}^\top (\mathbf{I}_d - 4\mathbf{A})^{-1} \mathbf{B} = \mathbf{I}_d, \quad \mathbf{C} = -\frac{1}{2} \mathbf{B}^\top (\mathbf{I}_d - 4\mathbf{A})^{-1} \mathbf{B} = -\frac{1}{2} \mathbf{I}_d, \quad D = \det(\mathbf{I}_d - 4\mathbf{A})^{1/4}. \quad (33)$$

Denote $\mathbf{Q} = (\mathbf{I}_d - 4\mathbf{A})^{-1/2} \mathbf{B} \in \mathbb{R}^{d \times d}$. Then, according to (33), $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}_d$, that is $\mathbf{Q} \in \mathbb{O}_d$. We rewrite (9) using (33) and then substitute $\mathbf{B} = (\mathbf{I}_d - 4\mathbf{A})^{1/2} \mathbf{Q}$:

$$\begin{aligned}
&\text{Var}_{\nu_{\text{SDE}}} f_{\text{SDE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}) f_{\text{SDE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y}) = \det(\mathbf{I}_d - 4\mathbf{A}) \det(\mathbf{I}_d - 8\mathbf{A})^{-1/2} \exp \left(-\|\mathbf{x}\|^2 \right. \\
&\quad \left. + 2\mathbf{x}^\top \mathbf{B}^\top (\mathbf{I}_d - 8\mathbf{A})^{-1} \mathbf{B} \mathbf{x} - \|\mathbf{y}\|^2 + 2\mathbf{y}^\top \mathbf{B}^\top (\mathbf{I}_d - 8\mathbf{A})^{-1} \mathbf{B} \mathbf{y} + 4\mathbf{x}^\top \mathbf{B}^\top (\mathbf{I}_d - 8\mathbf{A})^{-1} \mathbf{B} \mathbf{y} \right) \\
&\quad - K^{(0)}(\mathbf{x}, \mathbf{y})^2 = \det(\mathbf{I}_d - 4\mathbf{A})^{1/4} \det(\mathbf{I}_d - 8\mathbf{A})^{-1/2} \exp \left(-\|\mathbf{x}\|^2 - 2\mathbf{x}^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{x} - \|\mathbf{y}\|^2 \right. \\
&\quad \left. - 2\mathbf{y}^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{y} - 4\mathbf{x}^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{y} \right) - K^{(0)}(\mathbf{x}, \mathbf{y})^2 \quad (34)
\end{aligned}$$

where we denote:

$$\begin{aligned}
\mathbf{E} &= -(\mathbf{I}_d - 4\mathbf{A})^{1/2} (\mathbf{I}_d - 8\mathbf{A})^{-1} (\mathbf{I}_d - 4\mathbf{A})^{1/2} = -(\mathbf{I}_d - 4\mathbf{A})(\mathbf{I}_d - 8\mathbf{A})^{-1} \\
&= -\frac{1}{2} ((\mathbf{I}_d - 8\mathbf{A}) + \mathbf{I}_d) (\mathbf{I}_d - 8\mathbf{A})^{-1} = -\frac{1}{2} \mathbf{I}_d - \frac{1}{2} (\mathbf{I}_d - 8\mathbf{A})^{-1} \quad (35)
\end{aligned}$$

which is in \mathbb{D}_d since $\mathbf{A} \in \mathbb{D}_d$. Next, we observe:

$$2\mathbf{x}^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{x} + 2\mathbf{y}^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{y} + 4\mathbf{x}^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{y} = 2(\mathbf{x} + \mathbf{y})^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} (\mathbf{x} + \mathbf{y})$$

We plug this into (34) and use the resulting expression to rewrite (8) for $\langle \nu, f^{(1)}, f^{(2)} \rangle = \langle \nu_{\text{SDE}}, f_{\text{SDE}}^{(1)}, f_{\text{SDE}}^{(2)} \rangle$ as follows:

$$\begin{aligned}
\bar{\mathcal{L}}(\boldsymbol{\theta}_{\text{SDE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{SDE}}) &= L^{-2} \sum_{1 \leq i, j \leq L} \log(\text{Var}_{\nu_{\text{SDE}}} f_{\text{SDE}}^{(1)}(\boldsymbol{\omega}, \mathbf{x}^{(i)}) f_{\text{SDE}}^{(2)}(\boldsymbol{\omega}, \mathbf{y}^{(j)}) + K^{(0)}(\mathbf{x}^{(i)}, \mathbf{y}^{(j)})) \\
&= \log \det(\mathbf{I}_d - 4\mathbf{A}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\
&\quad - 2L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)})^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)}). \quad (36)
\end{aligned}$$

580 Using linearity and cyclic property of trace, we deduce that

$$\begin{aligned}
& L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)})^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)}) = \\
& = L^{-2} \sum_{1 \leq i, j \leq L} \text{Trace} \left((\mathbf{x}^{(i)} + \mathbf{y}^{(j)})^\top \mathbf{Q}^\top \mathbf{E} \mathbf{Q} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)}) \right) \\
& = L^{-2} \sum_{1 \leq i, j \leq L} \text{Trace} \left(\mathbf{Q}^\top \mathbf{E} \mathbf{Q} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)}) (\mathbf{x}^{(i)} + \mathbf{y}^{(j)})^\top \right) \\
& = \text{Trace} \left(\mathbf{Q}^\top \mathbf{E} \mathbf{Q} \left(L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)}) (\mathbf{x}^{(i)} + \mathbf{y}^{(j)})^\top \right) \right)
\end{aligned}$$

581 Observe that

$$\begin{aligned}
& L^{-2} \sum_{1 \leq i, j \leq L} (\mathbf{x}^{(i)} + \mathbf{y}^{(j)}) (\mathbf{x}^{(i)} + \mathbf{y}^{(j)})^\top = \\
& = L^{-2} \sum_{1 \leq i, j \leq L} \left(\mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top + \mathbf{x}^{(i)} (\mathbf{y}^{(j)})^\top + \mathbf{y}^{(j)} (\mathbf{x}^{(i)})^\top + \mathbf{y}^{(j)} (\mathbf{y}^{(j)})^\top \right) \\
& = L^{-1} \sum_{i=1}^L \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^\top + \left(L^{-1} \sum_{i=1}^L \mathbf{x}^{(i)} \right) \left(L^{-1} \sum_{j=1}^L \mathbf{y}^{(j)} \right)^\top + \left(L^{-1} \sum_{j=1}^L \mathbf{y}^{(j)} \right) \left(L^{-1} \sum_{i=1}^L \mathbf{x}^{(i)} \right)^\top \\
& \quad + L^{-1} \sum_{j=1}^L \mathbf{y}^{(j)} (\mathbf{y}^{(j)})^\top = \mathbf{M}^{(1)} + \boldsymbol{\mu}^{(4)} (\boldsymbol{\mu}^{(5)})^\top + \boldsymbol{\mu}^{(5)} (\boldsymbol{\mu}^{(4)})^\top + \mathbf{M}^{(2)}.
\end{aligned}$$

582 Denote $\mathbf{N} = \mathbf{M}^{(1)} + \boldsymbol{\mu}^{(4)} (\boldsymbol{\mu}^{(5)})^\top + \boldsymbol{\mu}^{(5)} (\boldsymbol{\mu}^{(4)})^\top + \mathbf{M}^{(2)}$. We conclude that

$$\begin{aligned}
\bar{\mathcal{L}}(\boldsymbol{\theta}_{\text{SDE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{SDE}}) &= \log \det(\mathbf{I}_d - 4\mathbf{A}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 \\
&\quad - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 - 2 \text{Trace}(\mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{N}). \tag{37}
\end{aligned}$$

583 With \mathbf{A} fixed, we minimize the right hand side of (37) with respect to \mathbf{Q} which is equivalent
584 to minimizing $\bar{\mathcal{L}}(\boldsymbol{\theta}_{\text{SDE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{SDE}})$ with respect to \mathbf{B} with fixed \mathbf{A} , since there is a one-to-one
585 correspondence between \mathbf{B} and \mathbf{Q} . This is equivalent to maximizing, again using the cyclic property
586 of trace,

$$\text{Trace}(\mathbf{Q}^\top \mathbf{E} \mathbf{Q} \mathbf{N}) = \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top) \tag{38}$$

587 with respect to \mathbf{Q} . We prove the following lemma first:

588 **Lemma A.3.** Suppose that diagonal entries of \mathbf{E} are all distinct, and the same holds for $\boldsymbol{\Lambda}^{(3)}$. Let
589 $\boldsymbol{\Pi} \in \{0, 1\}^{d \times d}$ be a permutation matrix sorting diagonal entries of \mathbf{E} (i.e. by applying $\boldsymbol{\Pi} \mathbf{E} \boldsymbol{\Pi}^\top$) in a
590 descending order corresponding to a permutation $\pi \in \mathbb{N}^d$. Set $\mathbf{Q}^* = \boldsymbol{\Pi}^\top (\mathbf{Q}^{(3)})^\top \in \mathbb{O}_d$. Then we
591 have:

$$\text{Trace}(\mathbf{E} \mathbf{Q}^* \mathbf{N} (\mathbf{Q}^*)^\top) = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \boldsymbol{\Lambda}_{l,l}^{(3)} \tag{39}$$

$$= \sup_{\mathbf{Q} \in \mathbb{O}_d} \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top) \tag{40}$$

592 *Proof.* First of all, we have:

$$\text{Trace}(\mathbf{E} \mathbf{Q}^* \mathbf{N} (\mathbf{Q}^*)^\top) = \text{Trace}(\mathbf{E} \boldsymbol{\Pi}^\top (\mathbf{Q}^{(3)})^\top \mathbf{N} \mathbf{Q}^{(3)} \boldsymbol{\Pi}) = \text{Trace}(\mathbf{E} \boldsymbol{\Pi}^\top \boldsymbol{\Lambda}^{(3)} \boldsymbol{\Pi}) \tag{41}$$

$$= \text{Trace} \left(\Pi \mathbf{E} \Pi^\top \mathbf{\Lambda}^{(3)} \right) = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)}, \quad (42)$$

593 i.e. (39) is satisfied.

594 Optimization for finding $\sup_{\mathbf{Q} \in \mathbb{O}_d} \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top)$ is a well-studied problem [8]. By the definition,
 595 $\mathbf{\Lambda}^{(3)}$ has eigenvalues of \mathbf{N} on the main diagonal and $\mathbf{E} \in \mathbb{D}_d$ hence it contains its eigenvalues on its
 596 main diagonal. Then, as proven in [8], \mathbf{Q}^* is indeed a global maximum of this problem in the case of
 597 distinct eigenvalues for \mathbf{E} and \mathbf{N} . That is, (40) is proven. \square

598 Next, we prove a generalization of Lemma A.3 when diagonal entries of \mathbf{E} and $\mathbf{\Lambda}^{(3)}$ are not necessarily
 599 distinct:

600 **Lemma A.4.** *Let $\Pi \in \{0, 1\}^{d \times d}$ be a permutation matrix sorting diagonal entries of \mathbf{E} (i.e. by*
 601 *applying $\Pi \mathbf{E} \Pi^\top$) in **any non-ascending** order corresponding to a permutation $\pi \in \mathbb{N}^d$. Set*
 602 *$\mathbf{Q}^* = \Pi^\top (\mathbf{Q}^{(3)})^\top \in \mathbb{O}_d$. Then we have:*

$$\text{Trace}(\mathbf{E} \mathbf{Q}^* \mathbf{N} (\mathbf{Q}^*)^\top) = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} \quad (43)$$

$$= \sup_{\mathbf{Q} \in \mathbb{O}_d} \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top) \quad (44)$$

603 *Proof.* In the same way as (41-42), we show that $\text{Trace}(\mathbf{E} \mathbf{Q}^* \mathbf{N} (\mathbf{Q}^*)^\top) = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)}$, i.e.
 604 (43) is satisfied. Next we prove that for any $\mathbf{Q} \in \mathbb{O}_d$,

$$\text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top) \leq \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)}. \quad (45)$$

605 which would imply (44).

606 Our proof is by contradiction. First of all, we can assume that $\mathbf{E}, \mathbf{\Lambda}^{(3)}$ are nonzero matrices since
 607 otherwise we have (44) trivially. Since $\text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top)$ is a continuous function of \mathbf{Q} and \mathbb{O}_d is
 608 compact, $\sup_{\mathbf{Q} \in \mathbb{O}_d} \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top)$ is finite. Suppose that there is $\delta > 0$ such that

$$\delta = \sup_{\mathbf{Q} \in \mathbb{O}_d} \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top) - \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)}. \quad (46)$$

609 Let $\tilde{\mathbf{E}}, \tilde{\mathbf{\Lambda}}^{(3)} \in \mathbb{D}_d$ be matrices with all distinct values on the diagonal such that

$$\|\tilde{\mathbf{E}} - \mathbf{E}\|_F \leq \min \left(\|\mathbf{E}\|_F, \frac{\delta}{12\|\mathbf{\Lambda}^{(3)}\|_F} \right), \quad \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F \leq \frac{\delta}{12\|\mathbf{E}\|_F} \quad (47)$$

610 where $\|\cdot\|_F$ denotes Frobenius norm and $\|\mathbf{E}\|_F, \|\mathbf{\Lambda}^{(3)}\|_F \neq 0$ since these are nonzero matrices.

611 Further, we assume that diagonal entries of $\tilde{\mathbf{\Lambda}}^{(3)}$ are sorted in a descending order and, in addition
 612 to $\mathbf{\Lambda}^{(3)}$, π also sorts entries of $\tilde{\mathbf{E}}$ in a non-ascending (descending) order. Clearly, such $\tilde{\mathbf{E}}, \tilde{\mathbf{\Lambda}}^{(3)}$ can
 613 be obtained by small perturbations of $\mathbf{E}, \mathbf{\Lambda}^{(3)}$. Also, denote $\tilde{\mathbf{N}} = \mathbf{Q}^{(3)} \tilde{\mathbf{\Lambda}}^{(3)} (\mathbf{Q}^{(3)})^\top$. Since \mathbb{O}_d is
 614 a compact closed set and $\text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top)$ is a continuous function of \mathbf{Q} , there exists $\mathbf{Q}^{**} \in \mathbb{O}_d$
 615 such that

$$\text{Trace}(\mathbf{E} \mathbf{Q}^{**} \mathbf{N} (\mathbf{Q}^{**})^\top) = \sup_{\mathbf{Q} \in \mathbb{O}_d} \text{Trace}(\mathbf{E} \mathbf{Q} \mathbf{N} \mathbf{Q}^\top). \quad (48)$$

616 By the definition of $\tilde{\mathbf{E}}, \tilde{\mathbf{\Lambda}}^{(3)}, \tilde{\mathbf{N}}$, we have:

$$\begin{aligned} & \text{Trace}(\mathbf{E} \mathbf{Q}^{**} \mathbf{N} (\mathbf{Q}^{**})^\top) - \text{Trace}(\tilde{\mathbf{E}} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top) \\ &= \left(\text{Trace}(\mathbf{E} \mathbf{Q}^{**} \mathbf{N} (\mathbf{Q}^{**})^\top) - \text{Trace}(\mathbf{E} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top) \right) + \left(\text{Trace}(\mathbf{E} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top) - \text{Trace}(\tilde{\mathbf{E}} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top) \right) \end{aligned}$$

$$\begin{aligned}
& -\text{Trace} \left(\tilde{\mathbf{E}} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top \right) \\
& = \text{Trace} \left(\mathbf{E} \mathbf{Q}^{**} (\mathbf{N} - \tilde{\mathbf{N}}) (\mathbf{Q}^{**})^\top \right) + \text{Trace} \left((\mathbf{E} - \tilde{\mathbf{E}}) \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top \right).
\end{aligned}$$

Next, we apply Cauchy-Schwarz inequality to both terms:

$$\begin{aligned}
\text{Trace} \left(\mathbf{E} \mathbf{Q}^{**} (\mathbf{N} - \tilde{\mathbf{N}}) (\mathbf{Q}^{**})^\top \right) & \leq \|(\mathbf{Q}^{**})^\top \mathbf{E}\|_F \|(\mathbf{N} - \tilde{\mathbf{N}}) (\mathbf{Q}^{**})^\top\|_F = \|\mathbf{E}\|_F \|\mathbf{N} - \tilde{\mathbf{N}}\|_F, \\
\text{Trace} \left((\mathbf{E} - \tilde{\mathbf{E}}) \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top \right) & \leq \|(\mathbf{Q}^{**})^\top (\mathbf{E} - \tilde{\mathbf{E}})\|_F \|\tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top\|_F = \|\mathbf{E} - \tilde{\mathbf{E}}\|_F \|\tilde{\mathbf{N}}\|_F
\end{aligned}$$

where we use invariance of the Frobenius norm under multiplications by orthogonal matrices. Using this invariance again, we deduce that

$$\begin{aligned}
\|\mathbf{N} - \tilde{\mathbf{N}}\|_F & = \|\mathbf{Q}^{(3)} (\mathbf{\Lambda}^{(3)} - \tilde{\mathbf{\Lambda}}^{(3)}) (\mathbf{Q}^{(3)})^\top\|_F = \|\mathbf{\Lambda}^{(3)} - \tilde{\mathbf{\Lambda}}^{(3)}\|_F, \\
\|\tilde{\mathbf{N}}\|_F & = \|\mathbf{Q}^{(3)} \tilde{\mathbf{\Lambda}}^{(3)} (\mathbf{Q}^{(3)})^\top\|_F = \|\tilde{\mathbf{\Lambda}}^{(3)}\|_F.
\end{aligned}$$

We conclude that

$$\text{Trace} \left(\mathbf{E} \mathbf{Q}^{**} \mathbf{N} (\mathbf{Q}^{**})^\top \right) \leq \text{Trace} \left(\tilde{\mathbf{E}} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top \right) + \|\mathbf{E}\|_F \|\mathbf{\Lambda}^{(3)} - \tilde{\mathbf{\Lambda}}^{(3)}\|_F + \|\mathbf{E} - \tilde{\mathbf{E}}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)}\|_F. \quad (49)$$

Next, we apply Lemma A.3 to $\mathbf{E} = \hat{\mathbf{E}}$, $\mathbf{\Lambda}^{(3)} = \hat{\mathbf{\Lambda}}^{(3)}$ and deduce that

$$\begin{aligned}
\text{Trace} \left(\tilde{\mathbf{E}} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top \right) & \leq \sum_{l=1}^d \tilde{\mathbf{E}}_{\pi_l, \pi_l} \tilde{\mathbf{\Lambda}}_{l,l}^{(3)} = \sum_{l=1}^d \left(\mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \left(\tilde{\mathbf{E}}_{\pi_l, \pi_l} - \mathbf{E}_{\pi_l, \pi_l} \right) \mathbf{\Lambda}_{l,l}^{(3)} \right. \\
& \quad \left. + \tilde{\mathbf{E}}_{\pi_l, \pi_l} \left(\tilde{\mathbf{\Lambda}}_{l,l}^{(3)} - \mathbf{\Lambda}_{l,l}^{(3)} \right) \right) \\
& = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \sum_{l=1}^d \left(\tilde{\mathbf{E}}_{\pi_l, \pi_l} - \mathbf{E}_{\pi_l, \pi_l} \right) \mathbf{\Lambda}_{l,l}^{(3)} + \sum_{l=1}^d \tilde{\mathbf{E}}_{\pi_l, \pi_l} \left(\tilde{\mathbf{\Lambda}}_{l,l}^{(3)} - \mathbf{\Lambda}_{l,l}^{(3)} \right) \\
& = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \text{Trace} \left(\mathbf{\Pi} (\tilde{\mathbf{E}} - \mathbf{E}) \mathbf{\Pi}^\top \tilde{\mathbf{\Lambda}}^{(3)} \right) + \text{Trace} \left(\mathbf{\Pi} \tilde{\mathbf{E}} \mathbf{\Pi}^\top (\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}) \right).
\end{aligned}$$

We apply Cauchy-Schwarz inequality again to the second and the third term:

$$\begin{aligned}
\text{Trace} \left(\mathbf{\Pi} (\tilde{\mathbf{E}} - \mathbf{E}) \mathbf{\Pi}^\top \tilde{\mathbf{\Lambda}}^{(3)} \right) & \leq \|(\tilde{\mathbf{E}} - \mathbf{E}) \mathbf{\Pi}^\top\|_F \|\mathbf{\Pi}^\top \tilde{\mathbf{\Lambda}}^{(3)}\|_F = \|\tilde{\mathbf{E}} - \mathbf{E}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)}\|_F, \\
\text{Trace} \left(\mathbf{\Pi} \tilde{\mathbf{E}} \mathbf{\Pi}^\top (\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}) \right) & \leq \|\tilde{\mathbf{E}} \mathbf{\Pi}^\top\|_F \|\mathbf{\Pi}^\top (\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)})\|_F = \|\tilde{\mathbf{E}}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F
\end{aligned}$$

where we use invariance of the Frobenius norm under column and row permutations. We conclude that

$$\text{Trace} \left(\tilde{\mathbf{E}} \mathbf{Q}^{**} \tilde{\mathbf{N}} (\mathbf{Q}^{**})^\top \right) \leq \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \|\tilde{\mathbf{E}} - \mathbf{E}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)}\|_F + \|\tilde{\mathbf{E}}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F.$$

We combine this inequality with (49) and obtain:

$$\text{Trace} \left(\mathbf{E} \mathbf{Q}^{**} \mathbf{N} (\mathbf{Q}^{**})^\top \right) \leq \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + 2\|\mathbf{E}\|_F \|\mathbf{\Lambda}^{(3)} - \tilde{\mathbf{\Lambda}}^{(3)}\|_F + 2\|\mathbf{E} - \tilde{\mathbf{E}}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)}\|_F. \quad (50)$$

Next, we use triangle inequality and deduce that

$$\|\tilde{\mathbf{\Lambda}}^{(3)}\|_F \leq \|\mathbf{\Lambda}^{(3)}\|_F + \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F.$$

Hence, we continue (50):

$$\text{Trace} \left(\mathbf{E} \mathbf{Q}^{**} \mathbf{N} (\mathbf{Q}^{**})^\top \right) \leq \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + 2\|\mathbf{E}\|_F \|\mathbf{\Lambda}^{(3)}\|_F$$

$$\begin{aligned}
& -\tilde{\mathbf{\Lambda}}^{(3)}\|_F + 2\|\mathbf{E} - \tilde{\mathbf{E}}\|_F \left(\|\mathbf{\Lambda}^{(3)}\|_F + \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F \right) \\
& = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + 2\|\mathbf{E}\|_F \|\mathbf{\Lambda}^{(3)} - \tilde{\mathbf{\Lambda}}^{(3)}\|_F + 2\|\mathbf{\Lambda}^{(3)}\|_F \|\mathbf{E} - \tilde{\mathbf{E}}\|_F + 2\|\mathbf{E} - \tilde{\mathbf{E}}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F \\
& \leq \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + 2\|\mathbf{E}\|_F \|\mathbf{\Lambda}^{(3)} - \tilde{\mathbf{\Lambda}}^{(3)}\|_F + 2\|\mathbf{\Lambda}^{(3)}\|_F \|\mathbf{E} - \tilde{\mathbf{E}}\|_F + 2\|\mathbf{E}\|_F \|\tilde{\mathbf{\Lambda}}^{(3)} - \mathbf{\Lambda}^{(3)}\|_F
\end{aligned}$$

628 where in the last transition we use $\|\mathbf{E} - \tilde{\mathbf{E}}\|_F \leq \|\mathbf{E}\|_F$ which is according to (47). We continue this
629 chain of inequalities using (47) again:

$$\begin{aligned}
\text{Trace}(\mathbf{E}\mathbf{Q}^{**}\mathbf{N}(\mathbf{Q}^{**})^\top) & \leq \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \frac{2}{12}\delta + \frac{2}{12}\delta + \frac{2}{12}\delta = \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \frac{\delta}{2} \\
& < \sum_{l=1}^d \mathbf{E}_{\pi_l, \pi_l} \mathbf{\Lambda}_{l,l}^{(3)} + \delta.
\end{aligned}$$

630 This is a contradiction with (46) taking into account \mathbf{Q}^{**} 's definition (48). Hence, (44) is proven. \square

631 Let \mathbf{Q}^* be defined as in Lemma A.4's statement. Further, we denote $\pi(\mathbf{A}) = \pi$, $\Pi(\mathbf{A}) = \Pi$ where
632 π, Π are defined as in Lemma A.4's statement. That is, $\pi(\mathbf{A})$ denotes some permutation which sorts
633 diagonal entries of \mathbf{E} in a non-ascending order. It's a function of \mathbf{A} since \mathbf{E} is a function of \mathbf{A} defined
634 in (35). In fact, based on (35), we see that $\pi(\mathbf{A})$ is some permutation which sorts diagonal entries of
635 \mathbf{A} in a non-descending order. $\Pi(\mathbf{A})$ denotes a permutation matrix corresponding to $\pi(\mathbf{A})$. That is,
636 diagonal entries of $\Pi(\mathbf{A})\mathbf{A}\Pi(\mathbf{A})^\top$ are sorted in a non-descending order.

637 Let $\mathcal{G}(\mathbf{A})$ denote the right hand side of (37) where we substitute $\mathbf{Q} = \mathbf{Q}^*$. That is, $\mathcal{G}(\mathbf{A})$ is an
638 optimal value of $\bar{\mathcal{L}}(\theta_{\text{SDE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{SDE}})$ with \mathbf{A} fixed:

$$\begin{aligned}
\mathcal{G}(\mathbf{A}) & = \log \det(\mathbf{I}_d - 4\mathbf{A}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\
& \quad - 2 \sum_{l=1}^d \mathbf{E}_{\pi(\mathbf{A})_l, \pi(\mathbf{A})_l} \mathbf{\Lambda}_{l,l}^{(3)} \\
& = \log \det(\mathbf{I}_d - 4\mathbf{A}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\
& \quad + \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{\pi(\mathbf{A})_l, \pi(\mathbf{A})_l})^{-1}) \mathbf{\Lambda}_{l,l}^{(3)} \tag{51}
\end{aligned}$$

639 where we use \mathbf{E} 's definition (35). Let $\pi^{-1}(\mathbf{A}) \in \mathbb{N}^d$ denote a permutation inverse to $\pi(\mathbf{A})$. By
640 rearranging terms in the sum, we have:

$$\sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{\pi(\mathbf{A})_l, \pi(\mathbf{A})_l})^{-1}) \mathbf{\Lambda}_{l,l}^{(3)} = \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{l,l})^{-1}) \mathbf{\Lambda}_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}.$$

641 Therefore, we have:

$$\begin{aligned}
\mathcal{G}(\mathbf{A}) & = \log \det(\mathbf{I}_d - 4\mathbf{A}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\
& \quad + \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{l,l})^{-1}) \mathbf{\Lambda}_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}. \tag{52}
\end{aligned}$$

642 Define a new function $\mathcal{G}(\hat{\mathbf{A}}, \mathbf{A})$, where $\hat{\mathbf{A}} \in \mathbb{D}_d$ satisfies $8\hat{\mathbf{A}} \prec \mathbf{I}_d$, as follows:

$$\mathcal{G}(\hat{\mathbf{A}}, \mathbf{A}) = \log \det(\mathbf{I}_d - 4\hat{\mathbf{A}}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\hat{\mathbf{A}}) + \sum_{l=1}^d (1 - 8\hat{\mathbf{A}}_{l,l})^{-1} \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}.$$

643 By the definition of $\mathcal{G}(\hat{\mathbf{A}}, \mathbf{A})$, we have:

$$\mathcal{G}(\mathbf{A}) = \mathcal{G}(\mathbf{A}, \mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 + \sum_{l=1}^d \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}.$$

644 Hence, it holds:

$$\mathcal{G}(\mathbf{A}) \geq \inf_{\hat{\mathbf{A}} \in \mathbb{D}_d, 8\hat{\mathbf{A}} \prec \mathbf{I}_d} \mathcal{G}(\hat{\mathbf{A}}, \mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 + \sum_{l=1}^d \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}. \quad (53)$$

645 Next, we show that there is a closed-form expression for the solution of $\inf_{\hat{\mathbf{A}} \in \mathbb{D}_d, 8\hat{\mathbf{A}} \prec \mathbf{I}_d} \mathcal{G}(\hat{\mathbf{A}}, \mathbf{A})$.

646 Since $\hat{\mathbf{A}} \in \mathbb{D}_d$, we have: $\log \det(\mathbf{I}_d - 4\hat{\mathbf{A}}) = \sum_{l=1}^d \log(1 - 4\hat{\mathbf{A}}_{l,l})$, $\log \det(\mathbf{I}_d - 8\hat{\mathbf{A}}) =$
 647 $\sum_{l=1}^d \log(1 - 8\hat{\mathbf{A}}_{l,l})$. We further have:

$$\mathcal{G}(\hat{\mathbf{A}}, \mathbf{A}) = \sum_{l=1}^d \left(\log(1 - 4\hat{\mathbf{A}}_{l,l}) - \frac{1}{2} \log(1 - 8\hat{\mathbf{A}}_{l,l}) + (1 - 8\hat{\mathbf{A}}_{l,l})^{-1} \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)} \right). \quad (54)$$

648 From (54), we see that minimization $\inf_{\hat{\mathbf{A}} \in \mathbb{D}_d, 8\hat{\mathbf{A}} \prec \mathbf{I}_d} \mathcal{G}(\hat{\mathbf{A}}, \mathbf{A})$ reduces to d independent minimization

649 problems with respect to $\hat{\mathbf{A}}_{l,l}$ such that $8\hat{\mathbf{A}}_{l,l} < 1$. l 'th problem, $1 \leq l \leq d$, is solved using Lemma

650 A.1 where we set $\phi = \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}$. Let $\mathbf{A}^{**} \in \mathbb{D}_d$ denote the corresponding solution. Then,

651 for all $1 \leq l \leq d$, we have:

$$\mathbf{A}_{l,l}^{**} = \frac{1}{16} \left(1 - 2\Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)} - \sqrt{\left(2\Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)} + 1 \right)^2 + 8\Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}} \right). \quad (55)$$

652 From (53) it follows that

$$\begin{aligned} \mathcal{G}(\mathbf{A}) &\geq \mathcal{G}(\mathbf{A}^{**}, \mathbf{A}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 + \sum_{l=1}^d \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)} \\ &= \log \det(\mathbf{I}_d - 4\mathbf{A}^{**}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}^{**}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\ &\quad + \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{l,l}^{**})^{-1}) \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)}. \end{aligned} \quad (56)$$

653 Denote $\mathbf{E}^{**} = -\frac{1}{2}\mathbf{I}_d - \frac{1}{2}(\mathbf{I}_d - 8\mathbf{A}^{**})^{-1}$. Then we have:

$$\begin{aligned} \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{l,l}^{**})^{-1}) \Lambda_{\pi^{-1}(\mathbf{A})_l, \pi^{-1}(\mathbf{A})_l}^{(3)} &= -2 \text{Trace} \left(\mathbf{E}^{**} \Pi(\mathbf{A})^{-1} \Lambda^{(3)} (\Pi(\mathbf{A})^{-1})^\top \right) \\ &\leq -2 \sum_{l=1}^d \mathbf{E}_{\pi(\mathbf{A}^{**})_l, \pi(\mathbf{A}^{**})_l}^{**} \Lambda_{l,l}^{(3)} \end{aligned} \quad (57)$$

654 where the second transition follows from Lemma A.4 and the fact that $\pi(\mathbf{A}^{**})$ sorts diagonal

655 entries of \mathbf{A}^{**} in a non-descending order, hence its sorts diagonal entries of \mathbf{E}^{**} in a non-ascending

656 order (recall the definition of $\pi(\mathbf{A}^{**})$ and \mathbf{E}^{**}). Denote $\mathbf{E}^* = \Pi(\mathbf{A}^{**}) \mathbf{E}^{**} \Pi(\mathbf{A}^{**})^\top$. Then

657 $\mathbf{E}_{\pi(\mathbf{A}^{**})_l, \pi(\mathbf{A}^{**})_l}^{**} = \mathbf{E}_{l,l}^*$ for all $1 \leq l \leq d$ and

$$\sum_{l=1}^d \mathbf{E}_{\pi(\mathbf{A}^{**})_l, \pi(\mathbf{A}^{**})_l}^{**} \Lambda_{l,l}^{(3)} = \sum_{l=1}^d \mathbf{E}_{l,l}^* \Lambda_{l,l}^{(3)}. \quad (58)$$

Further, we have:

$$\begin{aligned}\mathbf{E}^* &= \Pi(\mathbf{A}^{**}) \left(-\frac{1}{2}\mathbf{I}_d - \frac{1}{2}(\mathbf{I}_d - 8\mathbf{A}^{**})^{-1} \right) \Pi(\mathbf{A}^{**})^\top \\ &= -\frac{1}{2}\mathbf{I}_d - \frac{1}{2}(\mathbf{I}_d - 8\Pi(\mathbf{A}^{**})\mathbf{A}^{**}\Pi(\mathbf{A}^{**})^\top)^{-1} \\ &= -\frac{1}{2}\mathbf{I}_d - \frac{1}{2}(\mathbf{I}_d - 8\mathbf{A}^*)^{-1}\end{aligned}$$

where we denote $\mathbf{A}^* = \Pi(\mathbf{A}^{**})\mathbf{A}^{**}\Pi(\mathbf{A}^{**})^\top$, i.e. $\mathbf{A}_{l,l}^* = \mathbf{A}_{\pi(\mathbf{A}^{**})_l, \pi(\mathbf{A}^{**})_l}^{**}$ for all $1 \leq l \leq d$.
Given the definition of \mathbf{A}^{**} (55), for all $1 \leq l \leq d$ we have:

$$\mathbf{A}_{l,l}^* = \frac{1}{16} \left(1 - 2\Lambda_{l,l}^{(3)} - \sqrt{\left(2\Lambda_{l,l}^{(3)} + 1\right)^2 + 8\Lambda_{l,l}^{(3)}} \right). \quad (59)$$

That is, \mathbf{A}^* is independent of \mathbf{A} . Based on (59), we see that smaller values of $\Lambda_{l,l}^{(3)}$ result in bigger values of $\mathbf{A}_{l,l}^*$. Since $\Lambda_{1,1}^{(3)}, \dots, \Lambda_{d,d}^{(3)}$ are ordered in a non-ascending order, we deduce that $\mathbf{A}_{1,1}^*, \dots, \mathbf{A}_{d,d}^*$ are ordered in a non-descending order. By the definition of $\pi(\mathbf{A}^*)$, we then have $\mathbf{A}_{l,l}^* = \mathbf{A}_{\pi(\mathbf{A}^*)_l, \pi(\mathbf{A}^*)_l}^*$ for all $1 \leq l \leq d$. Therefore,

$$\sum_{l=1}^d \mathbf{E}_{l,l}^* \Lambda_{l,l}^{(3)} = -\frac{1}{2} \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{l,l}^*)^{-1}) \Lambda_{l,l}^{(3)} = -\frac{1}{2} \sum_{l=1}^d (1 + (1 - 8\mathbf{A}_{\pi(\mathbf{A}^*)_l, \pi(\mathbf{A}^*)_l}^*)^{-1}) \Lambda_{l,l}^{(3)}.$$

Combining this with (58), (57), we can continue the chain of inequalities (56):

$$\begin{aligned}\mathcal{G}(\mathbf{A}) &\geq \log \det(\mathbf{I}_d - 4\mathbf{A}^{**}) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}^{**}) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\ &\quad + \sum_{l=1}^d \left(1 + (1 - 8\mathbf{A}_{\pi(\mathbf{A}^*)_l, \pi(\mathbf{A}^*)_l}^*)^{-1} \right) \Lambda_{l,l}^{(3)} \\ &= \log \det(\mathbf{I}_d - 4\mathbf{A}^*) - \frac{1}{2} \log \det(\mathbf{I}_d - 8\mathbf{A}^*) - L^{-1} \sum_{i=1}^L \|\mathbf{x}^{(i)}\|^2 - L^{-1} \sum_{j=1}^L \|\mathbf{y}^{(j)}\|^2 \\ &\quad + \sum_{l=1}^d \left(1 + (1 - 8\mathbf{A}_{\pi(\mathbf{A}^*)_l, \pi(\mathbf{A}^*)_l}^*)^{-1} \right) \Lambda_{l,l}^{(3)} = \mathcal{G}(\mathbf{A}^*)\end{aligned} \quad (60)$$

where in the second transition we use the fact that

$$\det(\mathbf{I}_d - 4\mathbf{A}^*) = \det(\Pi(\mathbf{A}^{**})(\mathbf{I}_d - 4\mathbf{A}^{**})\Pi(\mathbf{A}^{**})^\top) = \det(\mathbf{I}_d - 4\mathbf{A}^{**})$$

and, similarly, $\det(\mathbf{I}_d - 8\mathbf{A}^*) = \det(\mathbf{I}_d - 8\mathbf{A}^{**})$. In the third transition, we use definition of $\mathcal{G}(\cdot)$ (51). Note that (60) holds for all $\mathbf{A} \in \mathbb{D}_d$ such that $8\mathbf{A} \prec \mathbf{I}_d$ and also $8\mathbf{A}^* \prec \mathbf{I}_d$ since $8\mathbf{A}^{**} \prec \mathbf{I}_d$. We conclude that, when $\mathbf{B}, \mathbf{C}, \mathbf{D}$ are chosen optimally with a given \mathbf{A} , the minimum of $\bar{\mathcal{L}}(\theta_{\text{SDE}}; \mathcal{X}, \mathcal{Y}, \mathcal{T}_{\text{SDE}})$ is reached when $\mathbf{A} = \mathbf{A}^*$. As we have already deduced, diagonal entries of $\mathbf{A} = \mathbf{A}^*$ are sorted in the non-descending order. Hence, using Lemma A.4's notation, diagonal entries of \mathbf{E} are already sorted in a non-ascending sorting order and $\pi = (1, \dots, d)$, $\Pi = \mathbf{I}_d$ satisfy requirements of the Lemma. Hence, with $\mathbf{A} = \mathbf{A}^*$, the optimal \mathbf{B} has a form $(\mathbf{I}_d - 4\mathbf{A})^{1/2} \mathbf{Q}^*$ where $\mathbf{Q}^* = \mathbf{I}_d(\mathbf{Q}^{(3)})^\top = (\mathbf{Q}^{(3)})^\top$. Optimal \mathbf{C} and \mathbf{D} are further determined by (33). (11) follows from (60) and the fact that, as discussed above, we can replace $\pi(\mathbf{A}^*)_l$ with l in (60), $1 \leq l \leq d$. \square

B Additional experimental details

B.1 Compute resources and implementation

We use NumPy [24] in Google Colaboratory the variance comparison and kernel classification experiment. For the Transformer setups, we use TPU cluster and JAX [7] library. All tested Transformer variants were trained and tested on a TPU pods containing 4 TPU v3 chips with JAX and on GPUs (V100).

682 B.2 Variance comparison

683 We repeat the setup of [30] closely: we draw 5 pairs of sets $\{\mathbf{x}^{(i)}\}_{1 \leq i \leq L}$, $\{\mathbf{y}^{(j)}\}_{1 \leq j \leq L}$, $L = 1024$.
 684 On each pair, we compute the relative variance for all pairs of points and for all indicated RF methods.
 685 Further, the shifted log-variance is optimized separately on each pair of sets for GERF, ADERF and
 686 SDERF.

687 We take $M = 1$ since M 's value is not important in this experiment: bigger M would just shift the
 688 curves below. The reported curves are means over all pairs of points and over all 5 sets.

689 B.3 Kernel classification

690 As in [30], we obtain training, validation and test splits by shuffling the raw dataset and taking 90%,
 691 5%, 5% objects respectively. The splits are fixed for all RF methods. We tune σ on a logarithmic grid
 692 of 10 values on $[10^{-2}, 10^2]$. For each σ and each RF type, we try 50 seeds for drawing RFs during
 693 validation and testing. Testing is performed for the best σ only. Figure 3 reports averages over 50
 694 seeds. We use orthogonal ω 's for all types of RFs as described in [30], since orthogonal random
 695 features work better in practice [15, 30].

696 B.4 DERFs for long-sequence Transformers

697 B.4.1 Speech modelling

698 Our Conformer-Transducer variant was characterized by: **20** conformer layers, model_dim = **512**,
 699 relative position embedding dimensionality rped = **512** and $h = 8$ heads. We used batch size
 700 bs = **2048** and trained with the adam optimizer on TPUs. For the regular Conformer-Transducer
 701 training, we run ablation studies over different number of random features: $m = 8, 32, 128$. In the
 702 NST setting, we run experiments with $m = 8$. We reported commonly used metric: normalized word
 703 error rate (NWER).

704 B.4.2 Natural language processing

705 We pretrained BERT model on two publicly available datasets (see: Table 3). Following the original
 706 BERT training, we mask 15% of tokens in these two datasets, and train to predict the mask. All
 707 methods were warm started from exactly the same pre-trained checkpoint after 1M iteration of BERT
 708 pretraining. We used the exact same hyperparameter-setup for all the baselines (FAVOR++[30],
 709 FAVOR+ [15], ELU [25], ReLU [15]) and FAVOR++. The hyperparameters for pretraining are shown
 710 in Table 2. We finetuned on GLUE task, warm-starting with the weights of the pretrained model. The
 711 setup is analogous to the one from the original BERT paper.

Table 2: Hyperparameters for the base models for pre-training for all methods

Parameter	Value
# of heads	12
# of hidden layers	12
Hidden layer size	768
# of tokens	512
Batch size	256
M	256
Pretrain Steps	$1M$
Loss	MLM
Activation layer	gelu
Dropout prob	0.1
Attention dropout prob	0.1
Optimizer	Adam
Learning rate	10^{-4}
Compute resources	8×8 TPUv3

Table 3: Dataset used for pre training.

Dataset	# tokens	Avg. doc len.
Books [52]	1.0B	37K
Wikipedia	3.1B	592