



# CHORUS

This is the accepted manuscript made available via CHORUS. The article has been published as:

## Surface nematic order in iron pnictides

Kok Wee Song and Alexei E. Koshelev

Phys. Rev. B **94**, 094509 — Published 9 September 2016

DOI: [10.1103/PhysRevB.94.094509](https://doi.org/10.1103/PhysRevB.94.094509)

# Surface Nematic Order in Iron Pnictides

Kok Wee Song and Alexei E. Koshelev

*Materials Science Division, Argonne National Laboratory, Illinois, 60439, USA*

(Dated: August 26, 2016)

Electronic nematicity plays important role in iron-based superconductors. These materials have layered structure and theoretical description of their magnetic and nematic transitions has been well established in two-dimensional approximation, i.e., when the layers can be treated independently. However, the interaction between iron layers mediated by electron tunneling may cause non-trivial three-dimensional behavior. Starting from the simplest model for orbital nematic in a single layer, we investigate the influence of interlayer tunneling on bulk nematic order and possible preemptive state where this order is only formed near the surface. We found that the interlayer tunneling suppresses the bulk nematicity which makes favorable formation of a surface nematic above the bulk transition temperature. The purely electronic tunneling Hamiltonian, however, favors alternating from layer-to-layer nematic order parameter in the bulk. The uniform bulk state typically observed experimentally may be stabilized by the coupling with the elastic lattice deformation. Depending on strength of this coupling, we found three regimes: (i) surface nematic and alternating bulk order, (ii) surface nematic and uniform bulk order, and (iii) uniform bulk order without the intermediate surface phase. The intermediate surface-nematic state may resolve the current controversy about the existence of the weak nematic transition in the compound  $\text{BaFe}_2\text{As}_{2-x}\text{P}_x$ .

## I. INTRODUCTION

Iron-based superconductors are multiple-band layered materials.<sup>1-3</sup> The correlations of itinerant electrons in different bands create several collective excitations: magnetic, orbital, superconducting. The complex interplay between these excitations can be tuned by doping or pressure leading to rich phase diagrams with antiferromagnetic, nematic, and superconducting phases.

Parent materials have stripe antiferromagnetic and structural tetragonal-to-orthorhombic phase transition. The corresponding orders can be either established simultaneously, via a single first-order transition, or via two sequential second-order phase transitions with the structural transition always preceding the antiferromagnetic one. In particular, the first scenario is realized in  $\text{Ba}_{1-x}\text{K}_x\text{Fe}_2\text{As}_2$  (122 structure),<sup>4</sup> while the second scenario is realized above threshold doping level in  $\text{Ba}(\text{Fe}_{1-x}\text{Co}_x)_2\text{As}_2$ ,<sup>5,6</sup> in  $\text{ReFeAsO}_{1-x}\text{F}_x$  compounds (1111 structure), where Re is the rare-earth element (La, Pr, Sm, Ce),<sup>7-11</sup> and in  $\text{Na}_{1-\delta}\text{FeAs}$ .<sup>12,13</sup>

When the four-fold crystal symmetry breaks, at least three types of order emerge simultaneously: (i) orthorhombic lattice deformation, (ii) spin Ising-nematic order lifting degeneracy between the stripe-antiferromagnetic fluctuations in two orthogonal directions,<sup>14-17</sup> and (iii) energy split between the  $d_{zx}$  and  $d_{zy}$  Fe orbitals leading to density difference between the two electron bands (ferro-orbital order).<sup>18-21</sup> Two latter orders are realizations of elec-

tronic nematicity<sup>22</sup> and, most likely, the orthorhombic deformation is its consequence. This interpretation is supported by measurement of unusual resistivity anisotropy very sensitive to elastic stress,<sup>23,24</sup> softening of the elastic shear modulus,<sup>25-27</sup> optical conductivity,<sup>28</sup> and asymmetric shifts of orbital energies observed by ARPES.<sup>29</sup> Whether the electronic nematicity has predominantly magnetic or orbital origin is the subject of ongoing debate.<sup>30,31</sup>

In addition to the strong and well-established “main” simultaneous AFM and structural transition in the parent and P-doped 122 materials, torque,<sup>32</sup> NMR,<sup>33</sup> and optical-pumping<sup>34</sup> experiments suggested existence of an intermediate nematic phase with broken C4 symmetry emerging at temperatures  $\sim 20\text{K}$  above the main transition. These observations are also consistent with finite orbital splitting persisting above the bulk transition found for unstressed  $\text{Ba}(\text{Fe}_{1-x}\text{Co}_x)_2\text{As}_2$  crystals by ARPES.<sup>29</sup> On the other hand, the recent high-resolution specific heat measurements<sup>35</sup> clearly excluded possibility of the bulk phase transition in this temperature range.

One possibility to resolve this controversy is to assume the existence of a preemptive state, at which the nematic order nucleates first only at the surface and decays inside the bulk. In this paper, we investigate the role of interlayer tunneling on bulk nematic order and possibility of such preemptive surface nematic. We use the simplest model for a single layer in which we take into account only electron pockets and the orbital order appears due to Pomeranchuk instability caused by interaction between the pockets. We

found that the interlayer tunneling suppresses the bulk transition. As a consequence, it is favorable for the nematic order to form near the surface first. We found that the purely electronic tunneling Hamiltonian favors bulk nematic order parameter which alternates from layer to layer. Such alternating state is not realized in iron pnictides. The uniform bulk state observed experimentally may be stabilized by the coupling with the elastic lattice deformation. The similar nematic-lattice coupling has also been studied in Ref. 36 for the single-layer model using Monte-Carlo simulations and it was demonstrated that this coupling lifts the nematic transition above the magnetic transition. In this paper, depending on strength of the nematic-lattice coupling, we found three regimes: (i) surface nematic and alternating bulk order, (ii) surface nematic and uniform bulk order, and (iii) uniform bulk order without intermediate surface phase. In the first two scenarios, the nematic order at the onset of ordering instability has strong spatial variation in the out-of-plane direction with maximum at the surface. In the later discussions, this type of instability will be referred to as surface instability.<sup>37</sup> The decay length typically is of the order of several layer spacings from the surface. In the vicinity of transition between the second and third regimes, the decay length rapidly increases, and eventually diverges at the transition to the third regime. In the third regime the order parameter nucleates mostly uniformly inside the sample with some suppression near the surface (bulk instability).

The paper is organized as follows. In section II, we discuss the model of the nematic order for finite-size system. In section III, we consider the system free energy, locate the nematic instability, and find the stable ground state at the onset of transition for both infinite and finite-size system. In section IV, we discuss effects of the lattice on the electronic nematic order. We summarize and conclude the paper in section V.

## II. MODEL

### A. Single layer

We start from a single-layer model Hamiltonian

$$H = \sum_{\alpha, \mathbf{k}} \varepsilon_{\mathbf{k}}^{\alpha} c_{\alpha, \mathbf{k} s}^{\dagger} c_{\alpha, \mathbf{k} s} - \frac{uS}{2} \sum_{\mathbf{q}} \rho_{\mathbf{q}} \rho_{-\mathbf{q}}, \quad (1)$$

where  $s$  is the spin (the summation is implicitly assumed),  $S$  is the total area of the layer,  $\mathbf{k} = (k_x, k_y)$

is the in-plane momentum,  $\alpha = X$  and  $Y$  represent the electron pockets at  $(\pi, 0)$  and  $(0, \pi)$  in the 1-Fe Brillouin zone respectively. The electrons energy dispersions near the  $X$  and  $Y$  pockets are  $\varepsilon_{\mathbf{k}}^X = \frac{k_x^2}{2m_x} + \frac{k_y^2}{2m_y}$  and  $\varepsilon_{\mathbf{k}}^Y = \frac{k_x^2}{2m_y} + \frac{k_y^2}{2m_x}$  with the band masses  $m_x$  and  $m_y$ , where  $\mathbf{k}$  is measured from  $(\pi, 0)$  in the  $X$  pocket and from  $(0, \pi)$  in the  $Y$  pocket.  $\rho_{\mathbf{q}}$  is a charge collective mode with

$$\rho_{\mathbf{q}} = \frac{1}{S} \sum_{\mathbf{k}} (c_{Y, \mathbf{k}+\mathbf{q}, s}^{\dagger} c_{Y, \mathbf{k} s} - c_{X, \mathbf{k}+\mathbf{q}, s}^{\dagger} c_{X, \mathbf{k} s}).$$

The Hamiltonian is invariant under the exchange between  $X$  and  $Y$  which preserves the 4-fold rotational symmetry. This Hamiltonian has been discussed in Ref. 38 and can be obtained from a more general itinerant model.<sup>39</sup> We do not include the hole bands in the middle of the Brillouin zone which do not play role in the consideration.

If the coupling constant  $u$  is large enough and positive, the model can give rise to Pomeranchuk instability at  $\mathbf{q} = (0, 0)$ , and the Fermi surface (FS) is distorted in the ground state. The nematic order parameter of the model (1) can be identified as

$$\Delta_{\mathbf{q}} = u \langle \rho_{\mathbf{q}} \rangle.$$

where  $\langle \dots \rangle = \text{Tr}(\dots e^{-\beta(H-\mu\mathcal{N})}) / \text{Tr} e^{-\beta(H-\mu\mathcal{N})}$  is a trace over all many-body quantum states with  $\beta = 1/T$ , chemical potential  $\mu$ , and total number operator  $\mathcal{N} = \sum_{\mathbf{k}\alpha} c_{\alpha, \mathbf{k} s}^{\dagger} c_{\alpha, \mathbf{k} s}$ . This order parameter measures the difference between the electron densities in  $X$  and  $Y$  pockets in the ground state. Since the main orbital component of the itinerant electrons at  $X$  and  $Y$  pockets is  $d_{zx}$  and  $d_{zy}$  orbital respectively,<sup>40</sup> this nematic order has a close connection with the orbital ordering.

We will use the standard mean-field approximation, see, e.g., Ref. 41, which assumes that the fluctuations  $\langle \delta\rho_{\mathbf{q}} \delta\rho_{-\mathbf{q}} \rangle$  with  $\delta\rho_{\mathbf{q}} = \rho_{\mathbf{q}} - \langle \rho_{\mathbf{q}} \rangle$  are small, and  $\delta\rho_{\mathbf{q}} \delta\rho_{-\mathbf{q}}$  can be neglected in Eq. (1). This yields the mean-field Hamiltonian

$$H \simeq \sum_{\alpha, \mathbf{k}} \varepsilon_{\mathbf{k}}^{\alpha} c_{\alpha, \mathbf{k} s}^{\dagger} c_{\alpha, \mathbf{k} s} - S \sum_{\mathbf{q}} \rho_{\mathbf{q}} \Delta_{-\mathbf{q}} + \sum_{\mathbf{q}} \frac{S |\Delta_{\mathbf{q}}|^2}{2u}.$$

Furthermore, one may assume that homogeneous order is the most energetically favorable state in an ideal crystal,  $\Delta_{\mathbf{q}} = \Delta \delta(\mathbf{q})$ . This immediately leads to

$$H \simeq \sum_{\alpha, \mathbf{k}} (\varepsilon_{\mathbf{k}}^{\alpha} + V^{\alpha}) c_{\alpha, \mathbf{k}, s}^{\dagger} c_{\alpha, \mathbf{k} s} + \frac{S \Delta^2}{2u}, \quad (2)$$

where  $V^X = \Delta$  and  $V^Y = -\Delta$ .

We remark that the particular form of the microscopic model in Eq. (1) is not crucial for our study. The approach in this paper can also be applied to other microscopic models which have the similar effective mean-field Hamiltonian. Furthermore, the general framework of nematic order induced by Pomeranchuk instabilities was first discussed in the  $d$ -wave nematic order,<sup>42,43</sup> and also was used in Refs. 44–46. Recently, Pomeranchuk instability in FeSC has also been investigated using renormalization group<sup>31</sup> and quantum Monte Carlo<sup>47</sup> techniques, see also subsequent discussion on comparing the Monte-Carlo and analytic results.<sup>48</sup> In the next section we consider tunneling terms for the layered systems.

## B. Interlayer tunneling

In this section, we discuss the effective three-dimensional model which takes into account interlayer electronic tunneling. This model will be used for analyzing nematic transition in layered materials. The building block of the multilayer model in this paper is the single-layer mean-field Hamiltonian from Eq. (1) with homogeneous nematic order within the plane,  $\mathbf{q} = (0, 0)$ . We will start with the simple model taking into account only nearest neighbor interlayer tunneling terms. This model does not mix the X- and Y- pockets and allows for several analytical results in the limit of weak interlayer hopping constant. Unfortunately, this simple model does not quite describe situation for real crystal structure in 122 materials, where interlayer hoppings via pnictogen atoms extend beyond nearest neighbors, break 4-fold rotational symmetry, and mix X- and Y- pockets.<sup>49</sup> We, therefore, will extend our interlayer model to include these effects.

### 1. Nearest-neighbor hopping

With the single-layer Hamiltonian (2), the  $N$ -layer system with nearest-neighbor interlayer hopping can be modeled by the following mean-field Hamiltonian

$$\mathcal{H}_N = \sum_{\ell=1}^N \left[ \sum_{\alpha, \mathbf{k}} (\varepsilon_{\mathbf{k}}^{\alpha} + V_{\ell}^{\alpha}) c_{\alpha \ell \mathbf{k} \mathbf{s}}^{\dagger} c_{\alpha \ell \mathbf{k} \mathbf{s}} + \frac{S \Delta_{\ell}^2}{2u} \right] - t_z \sum_{\ell=1}^{N-1} \sum_{\alpha, \mathbf{k}} c_{\alpha \ell \mathbf{k} \mathbf{s}}^{\dagger} c_{\alpha, \ell+1, \mathbf{k} \mathbf{s}} + h.c., \quad (3)$$

where each layer is labeled by  $\ell = 1, \dots, N$ , and  $V_{\ell}^X = \Delta_{\ell}$ ,  $V_{\ell}^Y = -\Delta_{\ell}$  are the nematic order parameters in the  $\ell$ -th layer. Here, we let the amplitude of the order parameters vary along the  $z$ -direction, since in a finite-size system translational symmetry is broken explicitly at the surfaces. Since the important electronic correlations are intralayer, we assumed that the direct overlapping between the Fe-orbitals in different layers are negligible and ignore the interlayer electron-electron interaction in the model.

### 2. Hoppings beyond the nearest neighbors: X-Y hybridization

The nearest-neighbor tunneling in Eq. (3) gives the simplest model for the interlayer coupling. However, for the crystalline structure of the 122 family of iron pnictides, such as  $\text{BaFe}_2\text{As}_2$ , the interlayer tunneling is a complicated process which involves the hopping between the Fe-orbitals and pnictogen orbitals (see Fig. 1a). As a consequence, the interlayer hoppings beyond the nearest neighbor are as important as the nearest-neighbor hopping.

This leads to two modifications in Eq. (3) (see Appendix A). First, due to the hoppings to the second neighbors, the hopping term becomes  $\mathbf{k}$ -dependent.

$$\mathcal{H}'_{\text{tun}} = - \sum_{\ell=1}^{N-1} \sum_{\mathbf{k}} s_{\mathbf{k}} (c_{\alpha \ell \mathbf{k} \mathbf{s}}^{\dagger} c_{\alpha, \ell+1, \mathbf{k} \mathbf{s}} + h.c.), \quad (4)$$

where  $s_{\mathbf{k}} = 2t_z(\cos k_x + \cos k_y)$ . Second, the hoppings to the third neighbors break the 4-fold rotation symmetry and this introduces the hybridization between the electrons in X- and Y- pockets.<sup>49</sup> The hybridization Hamiltonian is given by

$$\mathcal{H}_{\text{hyb}} = \sum_{\ell=1}^{N-1} (-1)^{\ell+1} \sum_{\mathbf{k}} \lambda_{\mathbf{k}} (c_{X \ell \mathbf{k} \mathbf{s}}^{\dagger} c_{Y, \ell+1, \mathbf{k} \mathbf{s}} + c_{Y \ell \mathbf{k} \mathbf{s}}^{\dagger} c_{X, \ell+1, \mathbf{k} \mathbf{s}} + h.c.), \quad (5)$$

where  $\lambda_{\mathbf{k}} = 2t'_z \sin k_x \sin k_y$ .

## III. NEMATIC PHASE TRANSITION

### A. Single-layer nematic transition

First, we consider the transition temperature for a single-layer system, Eq. (2). The free energy per

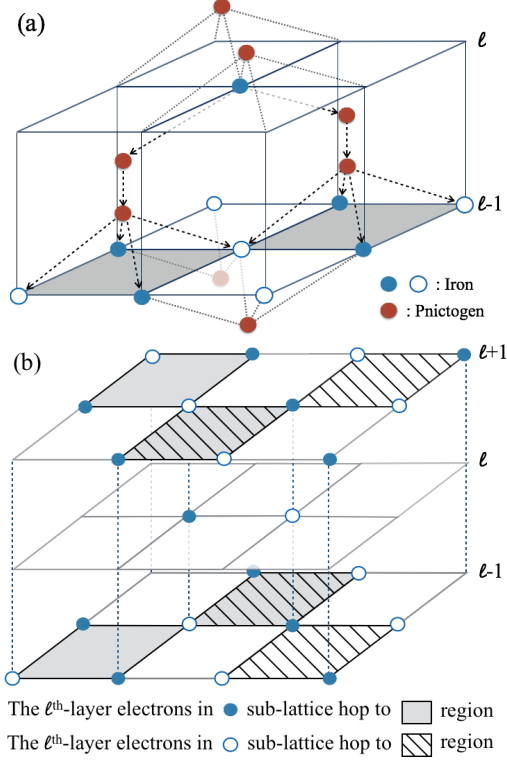


FIG. 1. Interlayer tunneling in 122 crystal: (a) The electron in the upper layer can tunnel to the shaded region in the lower layer via the pnictogens. These hopping processes should be treated at equal footing, since they all have equal tunneling probability. (b) The electron in the  $\ell$ -th layer odd (even) sub-lattice hopping to the shaded (tilted-lines) region in  $(\ell \pm 1)$ -th layer. The blue open circles and dots are the even and odd sublattice in the iron layer respectively, and the pnictogens in the lattice are not shown in the diagram.

unit area is

$$F_1[\Delta] = \frac{\Delta^2}{2u} - 2 \sum_{\alpha} \int_{\mathbf{k}} \ln \left[ 1 + e^{-\beta(\xi_{\mathbf{k}}^{\alpha} + V^{\alpha})} \right], \quad (6)$$

where  $\xi_{\mathbf{k}}^{\alpha} = \varepsilon_{\mathbf{k}}^{\alpha} - \mu$ , and  $\frac{1}{S} \sum_{\mathbf{k}} \rightarrow \int_{\mathbf{k}} = \int \frac{d^2 k}{(2\pi)^2}$ . The factor of two in the second term in Eq. (6) accounts for the spin degeneracy. To obtain the transition temperature, one can expand the free energy near the critical point and focus on the quadratic order term in  $F_1$ . Namely,

$$F_1[\Delta] \simeq F_1[0] + \frac{r_1}{2} \Delta^2$$

with the inverse nematic susceptibility

$$r_1 = \frac{1}{u} + 2 \sum_{\alpha} \int_{\mathbf{k}} n'_F(\xi_{\mathbf{k}}^{\alpha}), \quad (7)$$

where  $n_F(z) = [1 + \exp(\beta z)]^{-1}$  is the Fermi-Dirac distribution function. At the transition temperature, the inverse nematic susceptibility changes sign. Therefore, the nematic phase transition temperature  $T_0$  can be determined by setting this coefficient to zero and solve for  $\beta$ .

Furthermore, since the electronic modes far from the FS are suppressed by the Fermi-Dirac distribution factor, the upper limit of the momentum integral can be evaluated as  $\int_{\mathbf{k}} = \frac{\tilde{m}}{2\pi} \int_0^{\infty} d\varepsilon$  with  $\varepsilon = k_x^2/(2m_x) + k_y^2/(2m_y)$  and  $\tilde{m} = \sqrt{m_x m_y}$ . This yields the explicit result for the single-layer inverse susceptibility,

$$r_1 = \frac{1}{u} - \frac{\tilde{m}}{\pi} \left( 1 + \tanh \frac{\beta\mu}{2} \right). \quad (8)$$

In order to have a non-trivial solution, the coupling constant must satisfy the following condition

$$\frac{1}{4} < \frac{\tilde{m}u}{2\pi} < \frac{1}{2}, \quad (9)$$

since  $0 \leq \tanh \frac{\beta\mu}{2} \leq 1$ . For  $\frac{\tilde{m}u}{2\pi} \leq 1/4$ , no nematic order can be sustained in any temperature. For  $\frac{\tilde{m}u}{2\pi} \geq 1/2$ , the system is in the nematic phase for all temperatures with no phase transition. Therefore, the rest of the paper, we only consider the coupling constant in the region given by Eq. (9).

## B. Bulk nematic transition

In this section, we consider the bulk system in thermodynamic limit,  $N \rightarrow \infty$ . We start with the simplest model described by the Hamiltonian which takes into account only the nearest-neighbor hopping, Eq. (3). Then, we will generalize the model by including the hopping terms (4) and (5) beyond the nearest neighbor.

### 1. Nearest-neighbor hopping

The free energy of the system in the  $N \rightarrow \infty$  limit is convenient to calculate in the momentum space. The Fourier transformation of the field operators are

$$c_{\alpha\ell\mathbf{k}s} = \sum_{\mathbf{k}_z} \frac{e^{-ik_z\ell}}{\sqrt{N}} c_{\alpha k_z \mathbf{k}s}, \quad c_{\alpha k_z \mathbf{k}s} = \sum_{\ell=-\infty}^{\infty} \frac{e^{ik_z\ell}}{\sqrt{N}} c_{\alpha\ell\mathbf{k}s}.$$

The momentum space representation of the Hamiltonian (3) is

$$\mathcal{H} = \sum_{\mathbf{k}k'_z k_z} \psi_{k'_z}^{\dagger} (\hat{H}_b + \hat{V}_b) \psi_{k_z} + \frac{S}{2u} \frac{1}{N} \sum_q |\tilde{\Delta}_q|^2, \quad (10)$$

where  $k_z \in [-\pi, \pi]$  is the out-of-plane momentum,  $|\tilde{\Delta}_q|^2 = \tilde{\Delta}_q \tilde{\Delta}_{-q}$ , and  $\psi_{k_z}^\dagger = (c_{Xk_z\mathbf{k}s}^\dagger, c_{Yk_z\mathbf{k}s}^\dagger)$ . The  $2 \times 2$  matrices are

$$\hat{H}_b = \begin{pmatrix} \varepsilon_{\mathbf{k}}^X - 2t_z \cos k_z & 0 \\ 0 & \varepsilon_{\mathbf{k}}^Y - 2t_z \cos k_z \end{pmatrix} \delta_{k_z, k'_z}, \quad (11)$$

$$\hat{V}_b = \mathcal{V}_{k_z k'_z} = \frac{1}{N} \begin{pmatrix} \tilde{\Delta}_{k_z - k'_z} & 0 \\ 0 & -\tilde{\Delta}_{k_z - k'_z} \end{pmatrix}. \quad (12)$$

Furthermore, the order parameter in the momentum space is

$$\tilde{\Delta}_q = \sum_{\ell=-\infty}^{\infty} \Delta_\ell e^{-iq\ell}, \quad \Delta_\ell = \int_q \tilde{\Delta}_q e^{iq\ell},$$

where we replaced  $\frac{1}{N} \sum_q$  by  $\int_q = \int \frac{dq}{2\pi}$  in the limit  $N \rightarrow \infty$ . Note that  $\tilde{\Delta}_q$  is a periodic function of  $q$  with period  $2\pi$ ,  $\tilde{\Delta}_q = \tilde{\Delta}_{q \pm 2\pi}$ , and we choose the range of  $q$  to be  $[-\pi, \pi]$ . The free energy of this simple model can be immediately written down as

$$F_b[\tilde{\Delta}_q] = -\frac{1}{\beta S} \ln \text{Tr} e^{-\beta(\mathcal{H} - \mu\mathcal{N})} \\ = \int_q \frac{|\tilde{\Delta}_q|^2}{2u} - \frac{2}{\beta} \frac{1}{S} \text{tr} \ln [i\omega_n - \hat{H}_b - \hat{V}_b + \mu],$$

where  $\omega_n = 2\pi T(n + 1/2)$  are the fermionic Matsubara frequencies and the trace is the sum over all  $\omega_n$ ,  $\mathbf{k}$ ,  $k_z$ , and  $\alpha$ . To find the transition temperature, we expand the functional with respect to  $\tilde{\Delta}_q$  up to second order,

$$F_b[\tilde{\Delta}_q] \simeq F_b[0] + \int_q \frac{|\tilde{\Delta}_q|^2}{2u} \\ + \frac{1}{\beta} \sum_{\omega_n} \sum_{k_z, k'_z} \int_{\mathbf{k}} \text{tr} (\mathcal{G}_{k_z} \mathcal{V}_{k_z, k'_z} \mathcal{G}_{k'_z} \mathcal{V}_{k'_z, k_z}) \quad (13)$$

with  $\mathcal{G}_{k_z} = (i\omega_n - \hat{H}_b + \mu)^{-1}$ . Carrying out the trace explicitly in Eq. (13), we obtain

$$F_b[\tilde{\Delta}_q] = F_b[0] + \frac{1}{2} \int_q r_{b,q} |\tilde{\Delta}_q|^2, \quad (14)$$

where the inverse nematic susceptibility

$$r_{b,q} = \frac{1}{u} - \sum_{\alpha} \int_{\mathbf{k}} \int_{k_z} \frac{n_F(z_{k_z}^{\alpha}) - n_F(z_{k_z+q}^{\alpha})}{t_z [\cos k_z - \cos(k_z + q)]} \quad (15)$$

with  $\frac{1}{N} \sum_{k_z} \rightarrow \int_{k_z} = \int \frac{dk_z}{2\pi}$ ,  $z_{k_z}^{\alpha} = \xi_{\mathbf{k}}^{\alpha} - 2t_z \cos k_z$  and  $\xi_{\mathbf{k}}^{\alpha} = \varepsilon_{\mathbf{k}}^{\alpha} - \mu$ . In the small- $t_z$  limit, expansion over  $t_z$  gives

$$r_{b,q} \simeq r_1 + \left(1 + \frac{\cos q}{2}\right) r_t, \quad (16)$$

where  $r_1$  is given by Eq. (8) and  $r_t = \frac{8}{3} t_z^2 \int_{\mathbf{k}} n_F'''(\xi_{\mathbf{k}}^{\alpha})$ . Performing integration over the in-plane momentum, we obtain

$$r_t = \frac{2\tilde{m}\beta^2 t_z^2}{3\pi} \tanh\left(\frac{\beta\mu}{2}\right) \text{sech}^2\left(\frac{\beta\mu}{2}\right). \quad (17)$$

Similar to the single-layer case, the nematic phase transition occurs at the instability point:  $r_{b,q} = 0$  for some  $q$ . We can see that the second term in Eq. (16) is positive and monotonically decreasing function in  $q \in [0, \pi]$ . As a result, the inverse susceptibility  $r_{b,q}$  always has the minimum at  $q = \pi$  corresponding the leading instability for the system. This indicates that the purely electronic Hamiltonian considered in this section favors the bulk nematic order with alternating sign across different layers. This  $XY$  alternating order has broken symmetry with odd and even layers having enhanced density in  $X$ - and  $Y$ -pockets correspondingly. Such alternating broken symmetry pattern is consistent to the previous finding in Ref. 50. A similar behavior also occurs in the more general model with  $XY$ -hybridization that we discuss in the next section.

## 2. Interlayer-hopping with hybridization

The model which includes hoppings beyond the nearest neighbors is somewhat more complicated, since such hoppings make the even and odd layers not equivalent in 122 crystal structure<sup>49</sup> (see Fig. 1b). This doubles the unit cell in  $z$  direction and reduces the size of the Brillouin zone by half. The momentum space representation for the next-nearest hopping (4) and  $XY$ -hybridization (5) terms are

$$\mathcal{H}'_{\text{tun}} = -2 \sum_{\alpha} \sum_{\mathbf{k} k_z} s_{\mathbf{k}} \cos k_z c_{\alpha, k_z \mathbf{k}s}^\dagger c_{\alpha, k_z \mathbf{k}s}, \\ \mathcal{H}_{\text{hyb}} = -2i \sum_{\mathbf{k}} \sum'_{k_z} \lambda_{\mathbf{k}} \sin k_z (c_{X, k_z \mathbf{k}s}^\dagger c_{Y, k_z + \pi, \mathbf{k}s} \\ + c_{Y, k_z \mathbf{k}s}^\dagger c_{X, k_z + \pi, \mathbf{k}s} + h.c.),$$

where  $\sum'_{k_z}$  is the summation in the reduced space with  $k_z \in [-\pi, 0]$ . The mean-field Hamiltonian for the bulk becomes

$$\mathcal{H}'_b = \sum'_{k_z k'_z} \sum_{\mathbf{k}} \Psi_{k_z}^\dagger (\hat{H}'_b + \hat{V}'_b) \Psi_{k'_z} + \frac{S}{2u} \frac{1}{N} \sum_q |\tilde{\Delta}_q|^2, \quad (18)$$

where  $\Psi_{k_z}^\dagger = (c_{Xk_z\mathbf{k}s}^\dagger, c_{X, k_z + \pi, \mathbf{k}s}^\dagger, c_{Yk_z\mathbf{k}s}^\dagger, c_{Y, k_z + \pi, \mathbf{k}s}^\dagger)$ , and the  $4 \times 4$  matrices are

$$\hat{H}'_b = \begin{pmatrix} \varepsilon_{\mathbf{k}}^X - 2t_{\mathbf{k}} \sigma^z \cos k_z & 2\lambda_{\mathbf{k}} \sigma^y \sin k_z \\ 2\lambda_{\mathbf{k}} \sigma^y \sin k_z & \varepsilon_{\mathbf{k}}^Y - 2t_{\mathbf{k}} \sigma^z \cos k_z \end{pmatrix} \delta_{k_z, k'_z},$$

where  $t_{\mathbf{k}} = t_z(1 + 2 \cos k_x + 2 \cos k_y)$ ,  $\sigma^{x,y,z}$  are the Pauli matrices in the  $(k_z, k_z + \pi)$  space, and

$$\hat{V}'_b = \gamma'_{k_z k'_z} = \frac{1}{N} \begin{pmatrix} V_{k_z, k'_z} & 0 \\ 0 & -V_{k_z, k'_z} \end{pmatrix} \quad (19)$$

with the block matrix

$$V_{k_z k'_z} = \begin{pmatrix} \tilde{\Delta}_{k_z - k'_z} & \tilde{\Delta}_{k_z - k'_z - \pi} \\ \tilde{\Delta}_{k_z - k'_z + \pi} & \tilde{\Delta}_{k_z - k'_z} \end{pmatrix}.$$

Diagonalizing  $\hat{H}'_b$  yields the two energy-dispersion branches in the three-dimensional space

$$\varepsilon_{\mathbf{k}, k_z}^{\pm} = \varepsilon_{\mathbf{k}} \pm 2\sqrt{(\frac{\delta_{\mathbf{k}}}{2} - t_{\mathbf{k}} \cos k_z)^2 + \lambda_{\mathbf{k}}^2 \sin^2 k_z}, \quad (20)$$

$$r'_{b,q} = \frac{1}{u} - \sum_{\gamma=\pm 1} \int_{\mathbf{k}k_z} \left[ \frac{(\delta_{\mathbf{k}} - 2t_{\mathbf{k}} \cos k_z + \gamma \eta'_{\mathbf{k}k_z})(\delta_{\mathbf{k}} - 2t_{\mathbf{k}} \cos(k_z + q) + \gamma \eta'_{\mathbf{k}, k_z}) - 4\lambda_{\mathbf{k}}^2 \sin k_z \sin(k_z + q)}{4\gamma \eta'_{\mathbf{k}k_z} [\delta_{\mathbf{k}} t_{\mathbf{k}} - (t_{\mathbf{k}}^2 - \lambda_{\mathbf{k}}^2)(\cos(k_z + q) + \cos k_z)] (\cos(k_z + q) - \cos k_z) \coth \frac{\beta z_{\mathbf{k}}^{\gamma}}{2}} \right. \\ \left. - \left( \eta'_{\mathbf{k}k_z} \rightarrow \eta'_{\mathbf{k}, k_z + q} \text{ and } z_{\mathbf{k}k_z}^{\gamma} \rightarrow z_{\mathbf{k}, k_z + q}^{\gamma} \right) \right] \quad (22)$$

with  $z_{\mathbf{k}k_z}^{\pm} = \varepsilon_{\mathbf{k}}^{\pm} - \mu$ , and  $\eta'_{\mathbf{k}k_z} = 2[(\frac{1}{2}\delta_{\mathbf{k}} + t_{\mathbf{k}} \cos k_z)^2 + \lambda_{\mathbf{k}}^2 \sin^2 k_z]^{1/2}$ . Some special-case results for  $r'_{b,q}$  (circular FS with  $\delta_{\mathbf{k}} = 0$ , and  $q = 0, \pi$ ) can be found in Appendix B.

To find the transition temperature, we evaluated  $r'_{b,q}$  numerically. In the calculation, we let  $\mu = \varepsilon_0$  be the Fermi energy and introduce the reduced parameters as follows:  $\bar{\beta} = \beta \varepsilon_0$ ,  $\bar{t}_z = t_z / \varepsilon_0$ , and  $\bar{u} = um / (2\pi)$ . Furthermore, we note that the inverse susceptibility in Eq. (22) is just a linear function of the inverse coupling constant ( $1/\bar{u}$ ). Although the transition temperature depends on the coupling constants explicitly, changing  $\bar{u}$  does not change the qualitative behavior. Therefore, we use only one representative value  $\bar{u} = 0.35$  throughout the paper. Figure 2 shows the  $q$  dependence of  $r'_{b,q}$  for different temperatures. We found that, the  $q = \pi$  mode again has the smallest inverse susceptibility near the transition temperature meaning that the XY alternating order also persists in this general model. However, this does not correspond to experiment, in iron pnictides the nematic order is uniform in  $z$  direction. This may mean that the  $q = 0$  state is stabilized by external factors, such as elastic energy due to the lattice distortions induced by the nematic order. We address such stabilization below, in Section IV.

where  $\varepsilon_{\mathbf{k}} = (\varepsilon_{\mathbf{k}}^X + \varepsilon_{\mathbf{k}}^Y)/2$  and  $\delta_{\mathbf{k}} = (\varepsilon_{\mathbf{k}}^X - \varepsilon_{\mathbf{k}}^Y)/2$ .

The calculation for the free-energy functional of  $\mathcal{H}'_b$  is completely parallel to the previous section (see Appendix B for detail). Expanding the free energy up to the second order, we obtain

$$F'_b[\tilde{\Delta}_q] = -\frac{1}{\beta S} \ln \text{Tr} e^{-\beta(\mathcal{H}'_b - \mu \mathcal{N})} \\ = F'_b[0] + \frac{1}{2} \int_q r'_{b,q} |\tilde{\Delta}_q|^2, \quad (21)$$

where the inverse susceptibility  $r'_{b,q}$  in this general case is

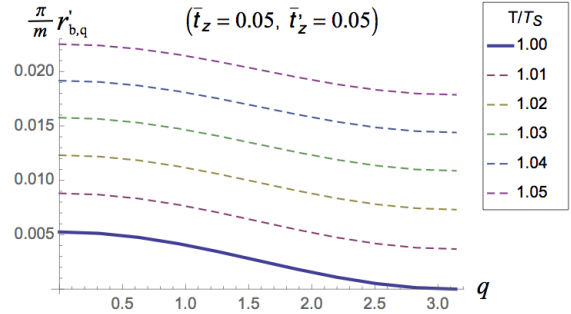


FIG. 2. The plots show the bulk inverse susceptibility  $r'_{b,q}$  for different temperatures. The temperature of the thick line is the nematic transition temperature ( $T_S$ ). The modes with  $q = \pi$  has the lowest value indicating that the system favors oscillating order parameter,  $\Delta_{\ell} = (-1)^{\ell} \Delta$ . The parameters in this plot are  $m_x = m_y$ ,  $\bar{u} = 0.35$ ,  $\bar{t}_z = \bar{t}'_z = 0.05$ .

### C. Finite-size system and surface nematic order

For the finite-size system, following the outline used for the bulk system in Sec. III B, we begin the discussion with only nearest-neighbor hopping. This model can be solved analytically allowing us to gain some insight about the system properties. Many of

these basic properties can also be found in the more realistic interlayer-hopping model. The model with  $XY$  pocket hybridization has to be solved numerically, and we also discuss the method in this section.

### 1. Nearest-neighbor hopping

Considering only the nearest-neighbor hopping, the free energy is

$$\begin{aligned} F_N[\Delta_\ell] &= -\frac{1}{\beta S} \ln \text{Tr} e^{-\beta(\mathcal{H}_N - \mu N)} \\ &= \sum_\ell \frac{\Delta_\ell^2}{2u} - \frac{2}{\beta} \frac{1}{S} \sum_\alpha \text{tr} \ln[(\mathcal{G}_{0,k}^\alpha)^{-1} + \mathcal{V}^\alpha]. \end{aligned} \quad (23)$$

Here we have introduced the notation  $k = (\omega_n, \mathbf{k})$ , the matrix  $(\mathcal{G}_{0,k}^\alpha)^{-1}$  is defined as

$$(\mathcal{G}_{0,k}^\alpha)^{-1} = \begin{pmatrix} (G_{0,k}^\alpha)^{-1} & -t_z & \cdots & 0 \\ -t_z & (G_{0,k}^\alpha)^{-1} & \ddots & \vdots \\ \vdots & \ddots & \ddots & -t_z \\ 0 & \cdots & -t_z & (G_{0,k}^\alpha)^{-1} \end{pmatrix},$$

where  $(G_{0,k}^\alpha)^{-1} = i\omega_n - \xi_{\mathbf{k}}^\alpha$  is the in-plane one-particle Green's function, and  $[\mathcal{V}^\alpha]_{\ell\ell'} = V_\ell^\alpha \delta_{\ell\ell'}$  is the diagonal matrix of order parameters.

To find the critical point, we have to expand the free energy, Eq. (23), with respect to  $\Delta_\ell$  up to the second order,

$$F_N[\Delta_\ell] \simeq F_N[0] + \frac{1}{2} \sum_{\ell\ell'} r_{\ell\ell'} \Delta_\ell \Delta_{\ell'}.$$

The inverse susceptibility is

$$r_{\ell\ell'} = \frac{\delta_{\ell\ell'}}{u} + \frac{2}{\beta} \sum_{\omega_n, \alpha} \int_{\mathbf{k}} [\mathcal{G}_0^\alpha]_{\ell\ell'} [\mathcal{G}_0^\alpha]_{\ell'\ell}, \quad (24)$$

where  $[\mathcal{G}_0^\alpha]_{\ell\ell'}$  is the tight-binding Green's function in the layer-index basis (see Appendix C),

$$[\mathcal{G}_0^\alpha]_{\ell\ell'} = \sum_p \frac{2/(N+1) \sin \ell \vartheta_p \sin \ell' \vartheta_p}{i\omega_n - \varepsilon_{\mathbf{k}}^\alpha + \mu + 2t_z \cos \vartheta_p} \quad (25)$$

with  $\vartheta_p = \frac{p\pi}{N+1}$  and  $p = 1 \dots N$ .

Carrying out the Matsubara frequency summation in Eq. (24) explicitly, we obtain

$$\begin{aligned} r_{\ell\ell'} &= \frac{\delta_{\ell\ell'}}{u} - \sum_{p,\alpha} \int_{\mathbf{k}} \left[ \frac{\beta}{2} (S_{\ell\ell'}^p)^2 \text{sech}^2 \left( \frac{\beta z_p^\alpha}{2} \right) \right. \\ &\quad \left. - \frac{1}{2} \sum_{p' \neq p} S_{\ell\ell'}^p S_{\ell'\ell}^{p'} \frac{\tanh \frac{\beta z_p^\alpha}{2} - \tanh \frac{\beta z_{p'}^\alpha}{2}}{2t_z (\cos \vartheta_p - \cos \vartheta_{p'})} \right]. \end{aligned} \quad (26)$$

where  $S_{\ell\ell'}^p = 2 \sin \ell \vartheta_p \sin \ell' \vartheta_p / (N+1)$ , and  $z_p^\alpha = \varepsilon_{\mathbf{k}}^\alpha - \mu - 2t_z \cos \vartheta_p$ . The quadratic matrix  $r_{\ell\ell'}$  in Eq. (26) contains the essential information for analyzing the nematic phase transition at the critical point.

Expanding Eq. (26) with respect to  $t_z$  and performing integration with respect to  $\mathbf{k}$  (see Appendix D), we obtain

$$\begin{aligned} r_{\ell\ell'} &\simeq r_1 \delta_{\ell\ell'} \\ &\quad + r_t \left[ \delta_{\ell\ell'} + \frac{\delta_{\ell, \ell'+1} + \delta_{\ell, \ell'-1}}{4} - \frac{\delta_{\ell, 1} + \delta_{\ell, N}}{2} \delta_{\ell\ell'} \right]. \end{aligned} \quad (27)$$

The first line is just the single-layer term derived before, Eq. (8), and the second line is the correction from the interlayer tunneling. The first two terms in the second line describe bulk interactions and correspond to previous derivation in the momentum space, Eqs. (16) and (17). As expected, the lowest-order expansion with respect to  $t_z$  gives only interaction between  $\Delta_\ell$  in neighboring layers. The coefficients  $r_{\ell\ell'}$  with  $|\ell - \ell'| > 1$  correspond to the higher-order corrections in  $t_z$ . This implies that the nematic order parameters do not have long-range interactions in the out-of-plane direction. The third term in the interlayer contribution represents the surface correction. Its negative sign implies that the surface favors the nematic order. Thus, in this model, surface transition should be expected. These properties, derived analytically from the nearest-neighbor model, also preserve in the more realistic model that includes the hoppings beyond the nearest neighbor, which we consider in the next section.

### 2. Inter-layer hopping with hybridization

For completeness, we also consider hopping processes beyond the nearest neighbor. Taking into account tunneling terms described by Eqs. (4) and (5), the free energy becomes

$$F'_N[\Delta_\ell] = -\frac{1}{\beta S} \ln \text{Tr} e^{-\beta(\mathcal{H}_N + \mathcal{H}'_{\text{tun}} + \mathcal{H}_{\text{hyb}} - \mu N)}.$$

As the  $N$ -layers Hamiltonian is just a one-body field operator, the free energy can be immediately written down as

$$F'_N = \sum_\ell \frac{\Delta_\ell^2}{2u} - \frac{2}{\beta} \sum_p \int_{\mathbf{k}} \ln \left[ 1 + e^{-\beta f_{p,\mathbf{k}}[\Delta_\ell]} \right], \quad (28)$$

where  $f_{p,\mathbf{k}}$  is the quasiparticle energy, which are the  $p$ -th eigenvalue ( $p = 1 \dots 2N$ ) of the following Hamiltonian matrix

$$\hat{H} = \begin{pmatrix} \mathbf{E}_{\mathbf{k}}^X & \mathbf{P}_{\mathbf{k}} \\ \mathbf{P}_{\mathbf{k}} & \mathbf{E}_{\mathbf{k}}^Y \end{pmatrix} + \begin{pmatrix} \nu^X & 0 \\ 0 & \nu^Y \end{pmatrix}. \quad (29)$$



The block matrices are

$$\mathbf{E}_{\mathbf{k}}^{\alpha} = \begin{pmatrix} \varepsilon_{\mathbf{k}}^{\alpha} - \mu & -t_{\mathbf{k}} & \cdots & 0 \\ -t_{\mathbf{k}} & \varepsilon_{\mathbf{k}}^{\alpha} - \mu & \ddots & \vdots \\ \vdots & \ddots & \ddots & -t_{\mathbf{k}} \\ 0 & \cdots & -t_{\mathbf{k}} & \varepsilon_{\mathbf{k}}^{\alpha} - \mu \end{pmatrix}, \quad (30a)$$

with  $t_{\mathbf{k}} = t_z(1 + 2 \cos k_x + 2 \cos k_y)$ ,

$$\mathbf{P}_{\mathbf{k}} = \lambda_{\mathbf{k}} \begin{pmatrix} 0 & (-1)^1 & \cdots & 0 \\ (-1)^1 & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & (-1)^{N-1} \\ 0 & \cdots & (-1)^{N-1} & 0 \end{pmatrix}, \quad (30b)$$

and  $[\mathcal{V}^{\alpha}] = V_{\ell}^{\alpha} \delta_{\ell\ell'}$ .

Due to the off-diagonal block matrix  $\mathbf{P}_{\mathbf{k}}$  in  $\hat{H}$ , there is no simple analytical expression for the free energy in this case. To evaluate the free energy, we use Eq. (28) and solve the eigenvalues problem numerically. The eigenvalues ( $f_{p,\mathbf{k}}$ ) of  $\hat{H}$  can be calculated by treating the second term in Eq. (29) as a perturbation, since  $\Delta_{\ell}$  are small at the critical point. To expand the eigenvalues at  $\Delta_{\ell} \simeq 0$ , we first solve for the eigenvalues and eigenvectors of the first term in Eq. (29), which we notate as

$$f_{p,\mathbf{k}}^{(0)}, \text{ and } \mathbf{x}_{p,\mathbf{k}}^T = (x_{p,\mathbf{k}}^1, \dots, x_{p,\mathbf{k}}^{2N})$$

respectively. The eigenvectors are normalized as  $\mathbf{x}_p^T \mathbf{x}_p = 1$ .

We apply perturbation expansion to  $f_{p,\mathbf{k}}[\Delta_{\ell}]$ . If  $N$  is even or odd with  $\lambda_{\mathbf{k}} \neq 0$  and  $m_x \neq m_y$ , all  $f_{p,\mathbf{k}}^{(0)}$  are distinct and non-degenerate. Therefore, the approximate eigenvalues are

$$\begin{aligned} f_{p,\mathbf{k}} &\simeq f_{p,\mathbf{k}}^{(0)} + \mathbf{x}_p^T \hat{V} \mathbf{x}_p + \sum_{p' \neq p} \frac{(\mathbf{x}_p^T \hat{V} \mathbf{x}_{p'})^2}{f_{p,\mathbf{k}}^{(0)} - f_{p',\mathbf{k}}^{(0)}}, \\ &= f_{p,\mathbf{k}}^{(0)} + \sum_{\ell} a_{\ell}^p \Delta_{\ell} + \sum_{\ell\ell'} b_{\ell\ell'}^p \Delta_{\ell} \Delta_{\ell'}, \end{aligned} \quad (31)$$

where  $\hat{V}$  is the second term in Eq. (29), and

$$a_{\ell}^p = v_{pp}^{\ell}, \quad b_{\ell\ell'}^p = \sum_{p' \neq p} \frac{v_{pp'}^{\ell} v_{pp'}^{\ell'}}{f_{p,\mathbf{k}}^{(0)} - f_{p',\mathbf{k}}^{(0)}} \quad (32)$$

with  $v_{pp'}^{\ell} = x_{p,\mathbf{k}}^{\ell} x_{p',\mathbf{k}}^{\ell} - x_{p,\mathbf{k}}^{N+\ell} x_{p',\mathbf{k}}^{N+\ell}$ . Therefore, using Eq. (31) near  $\Delta_{\ell} \simeq 0$  and expanding the free energy, we obtain the inverse susceptibilities with hybridization,

$$\begin{aligned} r'_{\ell\ell'} &= \frac{\delta_{\ell\ell'}}{u} - \sum_p \int_{\mathbf{k}} \left[ \frac{\beta}{2} a_{\ell}^p a_{\ell'}^p \operatorname{sech}^2 \left( \frac{\beta z_p}{2} \right) \right. \\ &\quad \left. + 2b_{\ell\ell'}^p \tanh \left( \frac{\beta z_p}{2} \right) \right], \end{aligned} \quad (33)$$

where  $z_p = f_{p,\mathbf{k}}^{(0)} - \mu$ . To facilitate the integration over the in-plane momentum  $\mathbf{k}$ , we approximate the energy dispersion near the FS as

$$\varepsilon_{\mathbf{k}}^{X,Y} \simeq \frac{\mathbf{k}^2}{2m} \pm \delta_2 \cos 2\theta, \quad (34)$$

where  $m = 2m_x m_y / (m_x + m_y)$  and  $\delta_2 = \varepsilon_0(1 - m_x/m_y)/2$ . The upper '+' (lower '-') sign is for the X (Y) pocket electrons. Furthermore, the momentum integration can be done by using  $\int_{\mathbf{k}} = \frac{m}{2\pi} \int_0^{\mu+\epsilon_c} d\varepsilon \int_0^{2\pi} \frac{d\theta}{2\pi}$ , where  $\epsilon_c$  is some cutoff energy of the model with the scale of bandwidth energy.

The case when  $N$  is odd and  $m_x = m_y$  requires special consideration in numerical calculations, see Appendix E. In this case the eigenspace of the first term in Eq. (29) breaks into  $N$  2-fold degenerate subspaces meaning that the formula in Eq. (32) has to be modified.

### 3. Calculation of transition temperature

We will use the same notations for the reduced parameters as in Sec. III B, i. e.,  $\beta = \beta\varepsilon_0$ ,  $\bar{t}_z = t_z/\varepsilon_0$ ,  $\bar{u} = um/(2\pi)$ , and, also,  $\bar{\delta}_2 = \delta_2/\varepsilon_0$ . To obtain the transition temperature,  $T_S$ , we search for the  $\bar{\beta} = \varepsilon_0/T_S$  such that the lowest eigenvalue of  $r'_{\ell\ell'}$ , Eq. (33), approaches zero meaning that at  $T_S$  the equation  $\sum_{\ell'=1}^N r'_{\ell,\ell'} \Delta_{\ell'} = 0$  has a nontrivial solution  $\Delta_{\ell'} \neq 0$ . We examine evolution of the transition temperature with increasing number of layers  $N$ .

For the nearest-neighbor model at  $t'_z=0$  and small  $t_z$ , when the inverse susceptibility is given by Eq. (27), the problem has simple analytical solution for  $N \gg 1$ . Near the surface  $\ell=0$  we obtain the following system

$$(r_1 + r_t) \Delta_{\ell} + \frac{r_t}{4} (\Delta_{\ell-1} + \Delta_{\ell+1}) - \frac{r_t}{2} \delta_{\ell,1} \Delta_1 = 0,$$

for  $\ell \geq 1$  and  $\Delta_0 = 0$ . Looking for solution in the form  $\Delta_{\ell} \propto (-1)^{\ell} \exp(-\varkappa\ell)$ , we obtain

$$r_1 + r_t - \frac{r_t}{2} \cosh \varkappa = 0, \text{ for } \ell > 1, \quad (35a)$$

$$r_1 + \frac{r_t}{2} - \frac{r_t}{4} \exp(-\varkappa) = 0, \text{ for } \ell = 1. \quad (35b)$$

These two equations yield

$$\exp \varkappa = 2 \quad (36)$$

and equation

$$r_1 + \frac{3r_t}{8} = 0, \quad (37)$$

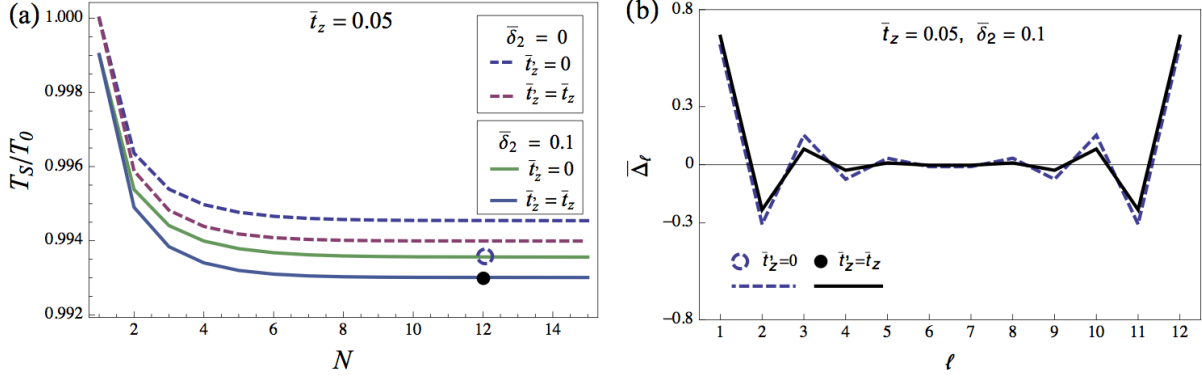


FIG. 3. (a) The plot of the nematic transition temperature  $T_S$  versus the number of layers in the system.  $T_0$  is the transition temperature of the single layer at  $\mu = \varepsilon_0$  with  $\bar{u} = 0.35$ . This plot indicates that the interlayer hoppings (both  $t$  and  $t'_z$ ) always lower the system transition temperature. Note that due to the  $\mathbf{k}$ -dependent in the hopping terms ( $t_{\mathbf{k}}$  and  $\lambda_{\mathbf{k}}$ ), the FS ellipticity ( $\delta_2$ ) also influences the  $T_S$ . (b) The spatial configuration of the lowest eigenmode near  $T_S$  for the system with  $N = 12$ . The dashed (solid) line correspond to the point in (a) that are marked by the open circle (filled circle). The order parameter has the maximum at the surface and decays in the bulk. We have defined  $\bar{\Delta}_\ell = \Delta_\ell / \sqrt{\sum_\ell \Delta_\ell^2}$ .

which determines the surface instability temperature for the nearest-neighbor model. This result has to be compared with the bulk-transition equation,

$$r_1 + \frac{r_t}{2} = 0,$$

which can be obtained by setting  $\varkappa = 0$  in Eq. (35a).

In general case with arbitrary  $t_z$  and  $t'_z$  we solved equations for  $T_S$  numerically. Fig. 3a shows representative dependences  $T_S(N)$  obtained for  $\bar{t}_z = 0.05$  and different  $\bar{t}'_z$  and  $\bar{\delta}_2$ . The transition temperature decreases as more layers are added to the system. For very large  $N$ ,  $T_S$  eventually approaches a finite limiting value. Furthermore,  $T_S$  always decreases as the hopping energies  $t_z$  and  $t'_z$  increase. This implies that the hopping between layers suppresses the nematic order.

Near the transition points, the sign change in the eigenvalue of  $r'_{\ell\ell'}$  also indicates the divergence in nematic susceptibility for the corresponding eigenmode, which signals an instability of this eigenmode. Examining the lowest eigenmode with zero eigenvalue at  $T_S$  allows us to deduce the most energetically favorable nematic spatial configuration. Figure 3b shows coordinate dependence of the unstable eigenmode for different parameters. We see that this eigenmode decays away from the surface so that for sufficiently large  $N$ , nematic order becomes vanishingly small at the center. This result implies that instability for finite-size systems at large  $N$  corresponds to formation of surface nematic and limiting values of  $T_S$  at large  $N$  in Fig. 3a correspond to sur-

face instability. In Fig. 4 we compare  $t_z$ -dependences of the bulk and surface transition temperatures. The split between these transitions increases with  $t_z$ . In

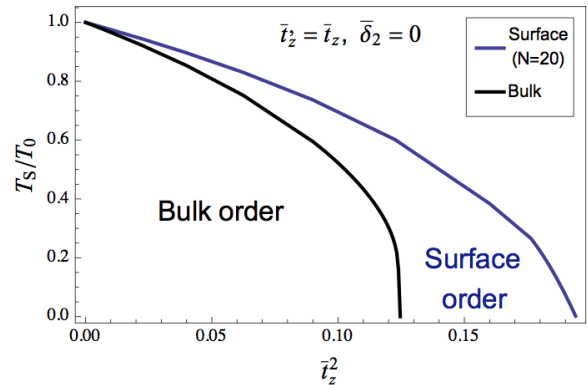


FIG. 4. This plot shows the  $t_z$  dependence of the bulk and surface transition temperatures. Both transitions are suppressed by the interlayer tunneling and nematic order can be completely destroyed if  $t_z$  is too large. The surface-nematic range rapidly increases with increasing  $t_z$ . Within some range of  $t_z$  only the surface nematic state exists.

addition, if  $t_z$  is too strong, the nematic phase transition disappears. In this large  $t_z$  case, the lowest eigenvalue of  $r'_{\ell\ell'}$  never becomes negative and the symmetry unbroken phase,  $\Delta_\ell = 0$ , always remains the true global minimum of  $F_N[\Delta_\ell]$ . Within some range of  $t_z$  only the surface nematic exists without bulk transition.

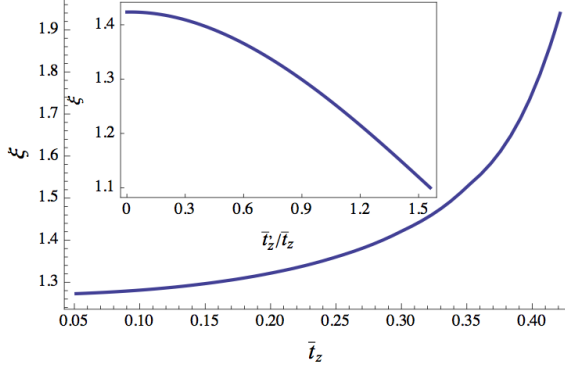


FIG. 5. The decay length of surface nematic obtained using system with  $N = 12$  for  $\bar{t}'_z = 0.05$ . Note that, the nematic order vanishes for  $\bar{t}'_z > 0.44$ . The inset shows  $\bar{t}'_z$  dependence of  $\xi$  for fixed  $\bar{t}'_z = 0.05$ .

To see how spatial configuration of the nematic order parameter depends on different hopping effects, we fit the nematic order parameter from the surface ( $1 \leq \ell < N/2$ ) to the exponential function,

$$|\bar{\Delta}_\ell| = \Delta_0 e^{-\ell/\xi} \quad (38)$$

where  $\xi$  is the decay length in units of the interlayer spacing, and  $\Delta_0$  is some constant. Figure 5 illustrates dependence of the decay length at  $T_S$  on the interlayer hopping parameters. We can see that the decay length is typically very small, only 1-2 interlayer spacings. It increases with increasing  $\bar{t}'_z$ , but decreases with increasing  $\bar{t}'_z$  (see the inset). Note that the value of  $\xi$  at  $\bar{t}'_z = 0$  reproduces the value  $1/\varkappa = 1/\ln(2) \approx 1.44$  analytically derived above, Eq. (36).

#### IV. THE EFFECTS FROM LATTICE DISTORTION

We found that the interlayer interactions mediated by the electronic tunneling favor  $XY$ -alternating nematic order. However, such order is not favorable for the lattice elastic energy. Indeed, if one layer is stressed in the  $X$ -direction and its neighboring layer is stressed in the  $Y$ -direction, this distortion increases the interlayer ion-ion distances in lattice. Thus, this costs higher elastic energy than stressing all layers in the same direction. Therefore, in the electron-lattice coupled system,  $XY$ -alternating order may not yield the lowest free energy.

For quantitative treatment of this problem, we consider the free energy with the following simple

extension

$$\mathcal{F}_N[\Delta_\ell, u_\ell] = F'_N[\Delta_\ell] + F_{el}[\Delta_\ell, u_\ell], \quad (39)$$

where the elastic part is modeled by

$$F_{el}[\Delta_\ell, u_\ell] = -g \sum_\ell u_\ell \Delta_\ell + \frac{1}{2} \sum_{\ell\ell'} C_{\ell,\ell'} u_\ell u_{\ell'}. \quad (40)$$

Here  $u_\ell = \frac{a_\ell - b_\ell}{a_\ell + b_\ell}$  is the  $\ell$ -th layer lattice distortion, and  $a_\ell$  and  $b_\ell$  are the in-plane lattice constant in  $x$ - and  $y$ -direction respectively.  $C_{\ell,\ell'}$  is the shear modulus constants, and  $g$  is the coupling between the nematic order parameter and lattice distortion. We will neglect temperature dependences of these parameters. For simplicity, we only consider the elastic matrix up to nearest neighbor and use notations  $C_{\ell,\ell} = C_s$  and  $C_{\ell,\ell\pm 1} = -C'_s$ . We assume  $C_{\ell,\ell\pm 1}$  to be negative so that it is favorable for layers in the lattice to be stressed in the same direction. The bulk shear modulus is given by  $C_{66} = (C_s - 2C'_s)/c_z$  with  $c_z$  being the  $c$ -axis lattice parameter meaning that the elastic constants must satisfy the condition  $C_s/C'_s > 2$  in order to have a stable lattice.

#### A. Bulk nematic transition

For the  $N \rightarrow \infty$  bulk limit, the free energy of the elastic part in the momentum space is

$$F_{el} = -g \int_q \tilde{\Delta}_q \tilde{u}_{-q} + \frac{1}{2} \int_q (C_s - 2C'_s \cos q) \tilde{u}_q \tilde{u}_{-q}$$

where  $\tilde{u}_q = \sum_{\ell=-\infty}^{\infty} e^{iq\ell} u_\ell$ . Minimizing  $F_{el}$  with respect to  $u_{-q}$ , we obtain

$$\tilde{u}_q = \frac{g \tilde{\Delta}_q}{C_s - 2C'_s \cos q}. \quad (41)$$

Therefore, the optimal elastic free energy is

$$F_{el} = -\frac{1}{2} \int_q \frac{g^2 |\tilde{\Delta}_q|^2}{C_s - 2C'_s \cos q}. \quad (42)$$

This new negative term modifies the inverse susceptibility as

$$\bar{r}_{b,q} = r'_{b,q} - \frac{g^2}{2(C_s - 2C'_s \cos q)}. \quad (43)$$

The elastic correction always reduces the transition temperature. Moreover, one can check that the elastic correction reduces the value of  $\bar{r}_{b,0}$  more than  $\bar{r}_{b,\pi}$ . This implies that with the elastic correction

$\bar{r}_{b,\pi}$  may not be the minimum inverse susceptibility any more. For sufficiently strong coupling  $g > g_{c1}$ ,  $\bar{r}_{b,0}$  drops below  $\bar{r}_{b,\pi}$  at the transition temperature and the system starts to favor a uniform order. The condition for the critical coupling strength is determined by

$$\bar{r}_{b,0}|_{T=T_S} = \bar{r}_{b,\pi}|_{T=T_S} = 0, \quad (44)$$

giving

$$\frac{2g_{c1}^2/C'_s}{(C_s/C'_s)^2 - 4} + r'_{b,0}|_{T=T_S} - r'_{b,\pi}|_{T=T_S} = 0. \quad (45)$$

In particular, for the model with only nearest-neighbor hopping in the small- $t_z$  limit, the inverse susceptibility has simple analytical form, Eq. (16). In this case Eq. (45) simply becomes

$$\frac{2g_{c1}^2/C'_s}{(C_s/C'_s)^2 - 4} + r_t|_{T=T_S} = 0, \quad (46)$$

where  $r_t$  is defined in Eq. (17).

For numerical analysis we introduce the reduced parameters  $\bar{g} = g\sqrt{\pi/mC'_s}$ , and  $C_s/C'_s$ . Figure 6 illustrates dependences of the difference  $\bar{r}_{b,0} - \bar{r}_{b,\pi}$  at the transition point on the reduced coupling strength  $\bar{g}^2$  for different ratios  $C_s/C'_s$ . The zero-crossing of these plots determines the critical coupling strength  $g_{c1}$ . We can see that it rapidly increases with  $C_s/C'_s$ . For more realistic model, which takes into account  $XY$ -pocket hybridization, the critical value of coupling  $g_{c1}$  has to be found numerically from Eq. (45) using full expression for the electronic inverse susceptibility  $r'_{b,q}$ , Eq. (22).

### B. Finite-size system and surface-nematic transition

For the finite-size system, to find the nematic order ground state, we minimize the free energy with respect to  $u_\ell$ .

$$\frac{\delta \mathcal{F}}{\delta u_\ell} = \frac{\delta F_{el}}{\delta u_\ell} = 0.$$

This yields a system of linear equations

$$-g\Delta_\ell + C_s u_\ell - C'_s(u_{\ell+1} + u_{\ell-1}) = 0 \quad (47)$$

with  $u_0 = u_{N+1} = 0$ . Solving the equation by inverting the Toeplitz tridiagonal matrix,<sup>51</sup> we obtain the required  $u_\ell$  which minimizes  $F_{el}$ ,

$$u_\ell = \frac{g}{C'_s} \sum_{\ell'} M_{\ell\ell'} \Delta_{\ell'} \quad (48)$$

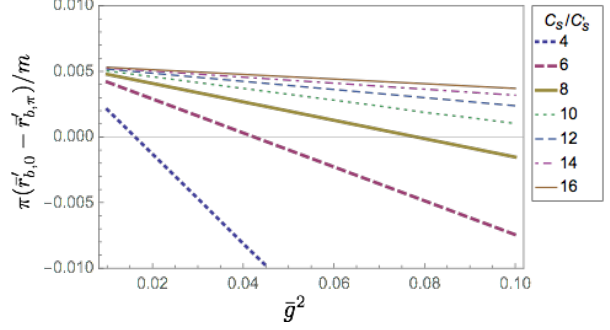


FIG. 6. The  $\bar{g}^2$  dependence of the inverse susceptibility  $r_q^b$  at the temperature with  $r_\pi^b = 0$ . At some  $\bar{g} = \bar{g}_{c1}$ , the difference  $r_{q=0}^b - r_{q=\pi}^b$  becomes negative at the transition temperature. This indicates that the uniform bulk nematic order becomes more favorable than the  $XY$ -alternating order. The plots are made using  $\bar{u} = 0.35$ ,  $\bar{t}_z = 0.05$ , and  $\bar{t}'_z = 0$ .

where the matrix  $M_{\ell\ell'}$  is

$$M_{\ell\ell'} = \frac{\cosh(\kappa\varphi_{\ell\ell'}^-) - \cosh(\kappa\varphi_{\ell\ell'}^+)}{2 \sinh \kappa \sinh[(N+1)\kappa]}$$

with  $\varphi_{\ell\ell'}^\pm = N+1 - |\ell \pm \ell'|$  and  $\cosh \kappa = C_s/(2C'_s)$ . Substituting these lattice distortions into  $F_{el}$ , we immediately obtain the optimal free energy as

$$F_{el} = - \sum_{\ell\ell'} \frac{g^2}{2C'_s} M_{\ell\ell'} \Delta_\ell \Delta_{\ell'}. \quad (49)$$

The optimal spatial configuration for the nematic order can be determined using Eqs. (39) and (49). We remark that the coupling to lattice leads to higher transition temperature. This indicates that lattice distortion actually promotes the formation of nematic order. Evolution of spatial dependence of the order parameter with increasing coupling strength  $\bar{g}$  is shown in Fig. 7a for  $N = 20$ . Other parameters are  $\bar{u} = 0.35$ ,  $\bar{t}_z = 0.05$ ,  $\bar{t}'_z = 0$ , and  $C_s/C'_s = 10$ . We can see that coupling to the lattice distortion can change the spatial configuration of the nematic order drastically. If the electron-lattice coupling  $\bar{g}$  is large enough, the order parameter no longer oscillates across different layers. On the other hand, the coupling to the lattice suppresses surface instability. When the coupling strength  $\bar{g}$  exceeds certain critical value,  $\bar{g}_{c2}$ , the intermediate surface nematic disappears and only bulk transition remains, see, e. g., the plot for  $\bar{g}^2 = 0.3$  in Fig. 7a. Finding this critical value requires careful study of finite-size effects for very large  $N$ . Fig. 7b shows the size dependence of the ratio of nematic order parameters at the center,  $\ell = N/2$ , and at the surface,

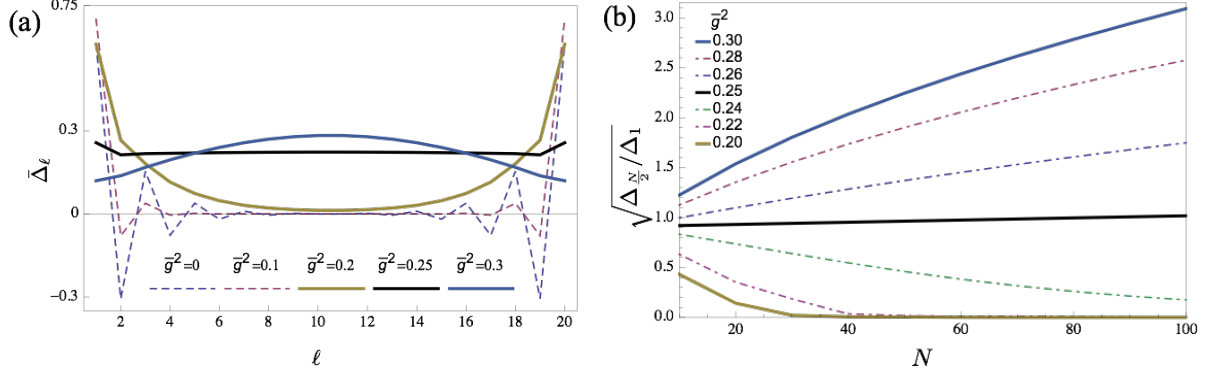


FIG. 7. (a) Spatial configuration of the nematic order with the effects of lattice distortion for  $N = 20$ ,  $\bar{t}_z = 0.05$ ,  $\bar{t}'_z = 0$ ,  $C_s/C'_s = 10$ , and different coupling strengths  $\bar{g}^2$ . Depending on the coupling strength  $\bar{g}$ , the nematic order ceases to oscillate across different layers at  $\bar{g}^2 \simeq 0.2$ . For  $\bar{g}^2 > 0.25$ , the instability corresponds to bulk transition. (b)  $N$  dependence of the order parameter at the center for three values of  $\bar{g}$  close to  $\bar{g}_{c2}$ . For  $\bar{g}^2 \lesssim 0.25$  the nematic order at the center approaches zero with increasing  $N$  corresponding to the surface order. The instability is bulk for  $\bar{g}^2 \gtrsim 0.25$ .

$\ell = 1$ , for seven values of  $\bar{g}$  close to  $\bar{g}_{c2}$ . We can see that for  $\bar{g} < 0.25$  this ratio decays to zero with increasing  $N$  (surface order) while for  $\bar{g} \geq 0.25$  it grows (bulk order).<sup>52</sup> This means that  $0.24 < \bar{g}_{c2}^2 < 0.25$ .

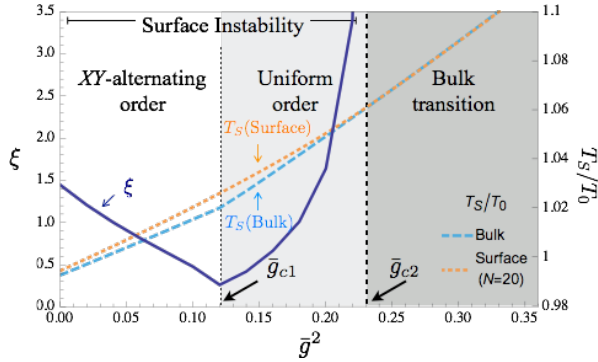


FIG. 8. The decay length (solid line) and the transition temperature (dotted line) of the system with  $N = 20$ ,  $\bar{t}_z = 0.05$ ,  $\bar{t}'_z = \bar{\delta}_2 = 0$ , and  $C_s/C'_s = 10$ . The dashed line is the transition temperature of the bulk system ( $N = \infty$ ). The kink is approximately located at  $\bar{g}_{c1}^2 \simeq 0.12$ , which is the transition point from XY-alternating order to uniform order. The decay length diverges at  $\bar{g}_{c2}^2 \simeq 0.25$  and only bulk transition remains after this point.

To further characterize the influence of lattice distortion on surface nematic, we computed its decay length by fitting the order parameters with Eq. (38). Figure 8 shows the dependence of the decay length,  $\xi$ , on the reduced coupling constant. We can see that this length is a nonmonotonic function of  $\bar{g}^2$ ; it decreases with  $\bar{g}^2$  for  $\bar{g} < \bar{g}_{c1}$  and increases with  $\bar{g}^2$  for

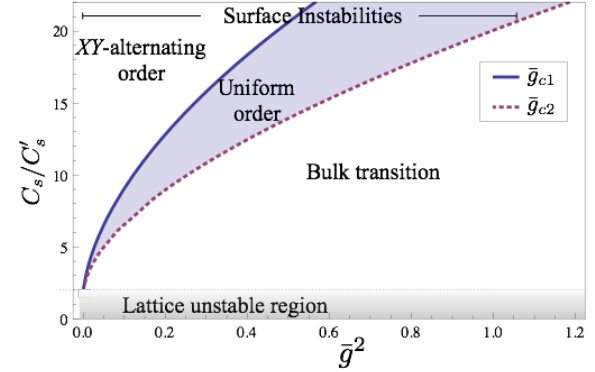


FIG. 9. This diagram shows three regions with different transition behaviors in the parameter plane  $\bar{g}^2 - C_s/C'_s$ . The shaded region corresponds to existence of surface instability for uniform bulk ground state.

$\bar{g} > \bar{g}_{c1}$ . A very small minimum value,  $\xi \approx 0.3$  is realized at the transition point from the alternating to uniform order,  $\bar{g} = \bar{g}_{c1}$ , where the dependence has a kink.<sup>53</sup> The decreasing  $\xi(\bar{g})$  in the alternating-order region can be easily understood.  $\xi^2$  is proportional to the stiffness  $r'' = d^2 r_{b,q}/dq^2|_{q=q_0}$ , where  $q_0 = 0, \pi$  is the ground-state wave vector. The lattice contribution to  $r_{b,q}$  has minimum at  $q = 0$  and maximum at  $q = \pi$ . Therefore, in the alternating state ( $q_0 = \pi$ ) the electronic part of  $r''$  is positive and the lattice contribution is negative meaning that the increase of  $\bar{g}$  reduces the stiffness  $r''$  and, correspondingly, decreases  $\xi$ . In the uniform state ( $q_0 = 0$ ) the growth of the lattice contribution enhances stiffness at  $q_0 = 0$

and increases  $\xi$ . We can see that the decay length diverges rapidly as  $\bar{g}$  approaching  $\bar{g}_{c2}$  from below, where the surface-nematic state disappears.

Figure 8 also shows the coupling-constant dependences of the surface and bulk transition temperatures. We can see that the bulk  $T_S$  has a kink at  $g = \bar{g}_{c1}$  while the surface transition is smooth at this point. The kink in bulk  $T_S$  is a natural consequence of qualitative change of the ground-state configuration from alternating to uniform order. At  $\bar{g} = \bar{g}_{c2}$  the surface transition smoothly merges with the bulk one. For  $\bar{g} > \bar{g}_{c2}$  the transition temperature evaluated for finite-size system with  $N = 20$  is indistinguishable from the bulk  $T_S$ .

Figure 8 summarizes the three possible phase-transition scenarios with the lattice distortion effects: (i) XY-alternating order with surface instability for  $\bar{g} < \bar{g}_{c1}$ , (ii) uniform order with surface instability for  $\bar{g}_{c1} < \bar{g} < \bar{g}_{c2}$ , and (iii) uniform bulk order without surface instability for  $\bar{g}_{c2} < \bar{g}$ . We explored in detail these different transition scenarios for  $\bar{t}_z = 0.05$  and  $\bar{t}'_z = 0$  and the results are summarized in the phase diagram presented in Fig. 9. The most interesting scenario, uniform bulk order with surface instability, is realized within the intermediate range of  $\bar{g}$  highlighted by the shaded area in the phase diagram. We can see that this range rapidly grows with increasing the ratio  $C_s/C'_s$  and can be rather wide.

## V. DISCUSSION AND CONCLUSIONS

To summarize our work, we have considered the bulk and surface nematic phase transitions in layered materials. The consideration is based on the simple single-layer two-band mean-field Hamiltonian and electronic tunneling between the layers. Evaluating the nematic free energy near the critical point, we have demonstrated that nematic order forms near surface before the main bulk transition. We also found that purely electronic tunneling Hamiltonian favors XY-alternating order in which sign of the order parameter changes from layer to layer.

Furthermore, lattice elasticity plays an important role in determining the ordering pattern in the multilayer system. In particular, coupling to the lattice may stabilize the uniform bulk nematic order. Depending on the coupling strength between the electrons and lattice, we found three different scenarios: (i) XY-alternating order with surface transition, (ii) Uniform order with surface transition, and (iii) Uniform bulk transition without surface instability. Scenario (ii) spans a considerable region

in the parameter space and may plausible for iron-pnictides. The surface nematic may be realization of the intermediate nematic state reported for P-doped 122 materials.<sup>32,34</sup> As several powerful experimental techniques, such as STM and ARPES, are inherently surface probes, the formation of preemptive surface nematic may strongly influence interpretation of experimental data.

Typically nematic instability in iron pnictides occurs either simultaneously or in the close proximity with the antiferromagnetic transition meaning that the spin fluctuations may strongly influence the formation of the nematic order. A proper microscopic treatment of these fluctuations is complicated and requires consideration of the hole bands in the zone center, which we leave for the future work. On the phenomenological level, the nematic order couples to the spin fluctuations linearly in the Landau-theory free-energy expansion,<sup>30</sup>  $\delta F \propto \Delta_\ell(M_{X,\ell}^2 - M_{Y,\ell}^2)$ , where  $M_{\alpha,\ell}$  are the fluctuating stripe-antiferromagnetic magnetizations in two perpendicular directions. Integrating out the spin fluctuations, this linear-coupling term generates the negative correction to the inverse nematic susceptibility which decays with  $|\ell - \ell'|$ . This means that the spin fluctuations promote the nematic ordering in each layer and favor the uniform order with respect to the XY-alternating order. We also expect that these fluctuations should suppress the surface instability. Therefore, effects of the spin fluctuations are qualitatively similar to ones of coupling to the lattice distortions.

In our study, the e-e interactions between different layers are ignored by assuming that overlapping between the Fe-orbital wave functions is negligible, since FeSC are layered materials. If the interlayer e-e interactions play an important role in driving the nematic order, long-range correlation can be built up between different layers. In this case, bulk nematic order can be more energetically favorable. We also note that the approach in this paper is only valid for the study of second-order phase transition near the critical temperature. To accurately describe the case of first-order phase transition or far away from the critical temperature, higher-order terms in the free energy have to be taken into account. These higher-order terms may change the ground state configuration drastically.

## ACKNOWLEDGMENTS

We would like to thank Ian Fisher for useful discussion. This work was supported by the Center

for Emergent Superconductivity, an Energy Frontier Research Center funded by the US DOE, Office of Science, under Award No. DEAC0298CH1088.

### Appendix A: Interlayer tight-binding model and XY-FS pocket hybridization

In iron pnictides with 122 composition, such as BaFe<sub>2</sub>As<sub>2</sub>, the off-diagonal out-of-plane hopping is important, because it modifies electronic spectrum qualitatively. The tight-binding Hamiltonian that includes the off-diagonal hopping in term of orbital basis is

$$\mathcal{H}_{\text{tun}} = - \left[ \sum_{\ell}^{\text{odd}} \left( \sum_{\mathbf{n}=\mathbf{n}_1} \sum_{\delta=\delta_1} + \sum_{\mathbf{n}=\mathbf{n}_2} \sum_{\delta=\delta_2} \right) + \sum_{\ell}^{\text{even}} \left( \sum_{\mathbf{n}=\mathbf{n}_1} \sum_{\delta=\delta_2} + \sum_{\mathbf{n}=\mathbf{n}_2} \sum_{\delta=\delta_1} \right) \right] \times \sum_{\bar{o}\bar{o}'} h^{\bar{o}\bar{o}'} d_{\bar{o},\ell,\mathbf{n}s}^{\dagger} d_{\bar{o}',\ell+1,\mathbf{n}+\delta,s} + h.c., \quad (\text{A1})$$

where  $d_{\bar{o},\ell,\mathbf{n}}^{\dagger}$  ( $d_{\bar{o},\ell,\mathbf{n}}$ ) is the orbital creation (annihilation) field operator with orbital index  $\bar{o} = 1, 2, 3$  standing for  $d_{xz}$ ,  $d_{yz}$ , and  $d_{xy}$  respectively. Furthermore,  $h^{\bar{o}\bar{o}'}$  is the tight-binding constant,  $\ell$  is the layer index, and  $\mathbf{n} = (n_x, n_y)$  is the lattice site index in the Fe-layer, and  $\mathbf{n}_1$  ( $\mathbf{n}_2$ ) is the lattice site with  $n_x + n_y = \text{odd}$  ( $n_x + n_y = \text{even}$ ). The following next-nearest hoppings have almost equal strength  $\delta_1 = (0, 0), (\pm 1, 0), (0, \pm 1), (-1, 1), (1, -1)$  and  $\delta_2 = (0, 0), (\pm 1, 0), (0, \pm 1), (-1, -1), (1, 1)$ .

We express the orbital field operators in the momentum space as

$$d_{\bar{o},\ell,\mathbf{n}s} = \frac{1}{\sqrt{A}} \sum_{\mathbf{k}} e^{-i\mathbf{k}\cdot\mathbf{n}} d_{\bar{o},\ell,\mathbf{k}s}, \quad (\text{A2})$$

where  $A$  is the total number of unit-cells in the Fe-layer. Substituting the above equation into  $\mathcal{H}_{\text{tun}}$ ,

$$\mathcal{H}_{\text{hyb}} = - \frac{1}{A} \sum_{\ell}^{\text{odd}} \sum_{\bar{o}\bar{o}'} h^{\bar{o}\bar{o}'} \sum_{\mathbf{k}\mathbf{k}'} \sum_{\mathbf{n}} 2(\cos k'_x \cos k'_y - e^{i\mathbf{Q}\cdot\mathbf{n}} \sin k'_x \sin k'_y) e^{i(\mathbf{k}-\mathbf{k}')\cdot\mathbf{n}} d_{\bar{o},\ell,\mathbf{k}s}^{\dagger} d_{\bar{o}',\ell+1,\mathbf{k}'s} + h.c. \\ - \frac{1}{A} \sum_{\ell}^{\text{even}} \sum_{\bar{o}\bar{o}'} h^{\bar{o}\bar{o}'} \sum_{\mathbf{k}\mathbf{k}'} \sum_{\mathbf{n}} 2(\cos k'_x \cos k'_y + e^{i\mathbf{Q}\cdot\mathbf{n}} \sin k'_x \sin k'_y) e^{i(\mathbf{k}-\mathbf{k}')\cdot\mathbf{n}} d_{\bar{o},\ell,\mathbf{k}s}^{\dagger} d_{\bar{o}',\ell+1,\mathbf{k}'s} + h.c. \quad (\text{A6})$$

Note that in the calculation we have used  $(-1)^{n_x+n_y} = e^{i\mathbf{Q}\cdot\mathbf{n}}$ . Carrying out the  $\mathbf{n}$  summation and combining

we break the tight-binding Hamiltonian as follows:  $\mathcal{H}_{\text{tun}} = \mathcal{H}_{\text{tun}}^0 + \mathcal{H}'_{\text{tun}} + \mathcal{H}_{\text{hyb}}$ .

For the direct hopping:  $\delta_1 = \delta_2 = (0, 0)$ ,

$$\mathcal{H}_{\text{tun}}^0 = \sum_{\ell}^{N-1} \sum_{\bar{o}\bar{o}',\mathbf{k}} h^{\bar{o}\bar{o}'} d_{\bar{o},\ell,\mathbf{k}s}^{\dagger} d_{\bar{o}',\ell+1,\mathbf{k}s} + h.c. \quad (\text{A3})$$

Note that, in this case, we have combined the even and odd layers and sub-lattice in the summation of  $\ell$  and  $\mathbf{n}$ . The orbital field operator can be expressed in term of band electron field operator as follows.

$$d_{\bar{o},\ell,\mathbf{k}s} = \sum_{\alpha} [\gamma_{\bar{\alpha}\mathbf{k}}^{\bar{o}}]^{-1} c_{\ell,\bar{\alpha},\mathbf{k}s}. \quad (\text{A4})$$

where  $\gamma_{\bar{\alpha}\mathbf{k}}^{\bar{o}}$  is the rotation matrix which diagonalized the single layer tight-binding Hamiltonian in the  $\mathbf{k}$ -space, and  $\bar{\alpha}$  is the corresponding band index. Substituting the Eq. (A4) into Eq. (A1), we obtain

$$\mathcal{H}_{\text{tun}}^0 = - \sum_{\ell}^{N-1} \sum_{\mathbf{k};\bar{\alpha}\bar{\alpha}'} \lambda_{1,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'} c_{\ell,\bar{\alpha},\mathbf{k}s}^{\dagger} c_{\ell+1,\bar{\alpha}',\mathbf{k}s} + h.c.,$$

where  $\lambda_{1,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'} = \sum_{\bar{o}\bar{o}'} h^{\bar{o}\bar{o}'} [(\gamma_{\bar{\alpha}\mathbf{k}}^{\bar{o}})^*]^{-1} [\gamma_{\bar{\alpha}\mathbf{k}}^{\bar{o}'}]^{-1}$ . If we ignore the  $\mathbf{k}$  dependence in  $\lambda_1$  and restrict the momentum near the FS then  $\lambda_{1,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'} \simeq t_z$  which is the direct hopping term in the Hamiltonian (3).

For the next-nearest neighbor hoppings:  $\delta_1 = (\pm 1, 0), (0, \pm 1)$ ,

$$\mathcal{H}'_{\text{tun}} = - \sum_{\ell}^{N-1} \sum_{\mathbf{k};\bar{\alpha}\bar{\alpha}'} 2\lambda_{1,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'} (\cos k_x + \cos k_y) \times (c_{\ell,\bar{\alpha},\mathbf{k}s}^{\dagger} c_{\ell+1,\bar{\alpha}',\mathbf{k}s} + h.c.). \quad (\text{A5})$$

This term modifies the interlayer nearest-neighbor hopping constant.

Turning to the next-next nearest-neighbor hopping:  $\delta_1 = (1, -1), (-1, 1)$  and  $\delta_2 = (1, 1), (-1, -1)$ , we derive

the  $\ell = \text{even}$  and odd terms, we obtain

$$\mathcal{H}_{\text{hyb}} = - \sum_{\ell} \sum_{\bar{o}\bar{o}'}^{N-1} h^{\bar{o}\bar{o}'} \sum_{\mathbf{k}} 2(\cos k_x \cos k_y d_{\bar{o},\ell,\mathbf{k}s}^{\dagger} d_{\bar{o}',\ell+1,\mathbf{k}s} + (-1)^{\ell} \sin k_x \sin k_y d_{\bar{o},\ell,\mathbf{k}s}^{\dagger} d_{\bar{o}',\ell+1,\mathbf{k}+\mathbf{Q},s}) + h.c. \quad (\text{A7})$$

The momentum  $\mathbf{k}$  is measured from  $(0, \pi)$  and  $\mathbf{Q} = (\pi, \pi)$ . The last term generates the hybridization between  $X$ - and  $Y$ - pockets. To see this, we write  $\mathcal{H}_{\text{hyb}}''$  in the band basis, and keeping only the last term in Eq. (A7),

$$\mathcal{H}_{\text{hyb}} = \sum_{\ell} \sum_{\mathbf{k}, \bar{\alpha}\bar{\alpha}'}^{N-1} (-1)^{\ell+1} 2\lambda_{2,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'} \sin k_x \sin k_y \times c_{\bar{\alpha},\ell,\mathbf{k}s}^{\dagger} c_{\bar{\alpha}',\ell+1,\mathbf{k}+\mathbf{Q},s}, \quad (\text{A8})$$

where  $\lambda_{2,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'} = \sum_{\bar{o}\bar{o}'} h^{\bar{o}\bar{o}'} [(\gamma_{\bar{\alpha}\mathbf{k}}^{\bar{o}})^*]^{-1} [\gamma_{\bar{\alpha}\mathbf{k}+\mathbf{Q}}^{\bar{o}'}]^{-1}$ . If we restrict the momentum to be near the FS, and regroup the band index into  $X$  and  $Y$  according to their momentum, this immediately lead to the  $X$ - and  $Y$ - pockets hybridization. As in the direct hopping term, for simplicity, ignoring the  $\mathbf{k}$ -dependence in  $\lambda_{2,\mathbf{k}}^{\bar{\alpha}\bar{\alpha}'}$  and set it to  $t'_z$ , we therefore obtain

$$\mathcal{H}_{\text{hyb}} \simeq \sum_{\ell=1}^{N-1} \sum_{\mathbf{k}} \lambda_{\mathbf{k}} (-1)^{\ell-1} (c_{X,\ell,\mathbf{k}s}^{\dagger} c_{Y,\ell+1,\mathbf{k},s} + c_{Y,\ell,\mathbf{k}s}^{\dagger} c_{X,\ell+1,\mathbf{k},s}) + h.c. \quad (\text{A9})$$

$$\mathcal{G}_{k_z} = \begin{pmatrix} [(i\omega_n - \varepsilon_{\mathbf{k}}^Y + \mu)\mathbb{I} + 2t_{\mathbf{k}}\sigma^z \cos k_z] \Omega_{k_z}^{-1} & -2\lambda_{\mathbf{k}}\sigma^y \sin k_z \Omega_{k_z+\pi}^{-1} \\ -2\lambda_{\mathbf{k}}\sigma^y \sin k_z \Omega_{k_z}^{-1} & [(i\omega_n - \varepsilon_{\mathbf{k}}^X + \mu)\mathbb{I} + 2t_{\mathbf{k}}\sigma^z \cos k_z] \Omega_{k_z+\pi}^{-1} \end{pmatrix}$$

with

$$\Omega_{k_z} = \begin{pmatrix} (G_{k_z}^X G_{k_z+\pi}^Y)^{-1} - 4\lambda_{\mathbf{k}}^2 \sin^2 k_z & 0 \\ 0 & (G_{k_z+\pi}^X G_{k_z}^Y)^{-1} - 4\lambda_{\mathbf{k}}^2 \sin^2 k_z \end{pmatrix}$$

and  $G_{k_z}^{\alpha} = i\omega_n - \varepsilon_{\mathbf{k}}^{\alpha} + 2t_{\mathbf{k}} \cos k_z + \mu$ . Writing out the trace explicitly and using the periodic condition in  $q \rightarrow q + 2\pi$ , this yields

$$F_b^{(2)} = -\frac{1}{2} \int_q r'_{b,q} \tilde{\Delta}_{-q} \tilde{\Delta}_q$$

with the inverse nematic susceptibility

$$r'_{b,q} = \frac{1}{u} - 2 \int_{\mathbf{k}k_z} \text{Res} \left[ \frac{[(z - z_{k_z+\pi}^Y)(z - z_{k_z+q+\pi}^Y) - (2\lambda_{\mathbf{k}})^2 \sin k_z \sin(k_z + q)] \frac{1}{2} \tanh \frac{\beta z}{2}}{[(z - z_{k_z}^X)(z - z_{k_z+\pi}^Y) - 4\lambda_{\mathbf{k}}^2 \sin^2 k_z][(z - z_{k_z+q}^X)(z - z_{k_z+q+\pi}^Y) - 4\lambda_{\mathbf{k}}^2 \sin^2(k_z + q)]} \right] + (X \leftrightarrow Y)$$

where  $z_{k_z}^{\alpha} = \varepsilon_{\mathbf{k}}^{\alpha} - \mu - 2t_{\mathbf{k}} \cos k_z$ . Note that the summation of  $k_z$  is running over  $k_z \in [-\pi, \pi]$ . The denominator has poles at  $z = z_{k_z}^{\pm} = \varepsilon_{\mathbf{k}k_z}^{\pm} - \mu$  and  $z = z_{k_z+q}^{\pm} = \varepsilon_{\mathbf{k},k_z+q}^{\pm} - \mu$ , where  $\varepsilon_{\mathbf{k}k_z}^{\pm}$  is defined in Eq. (20). Factorizing

with  $\lambda_{\mathbf{k}} = 2t'_z \sin k_x \sin k_y$ .

## Appendix B: Derivation of the bulk free energy

In this section, we derive the free energy in  $N \rightarrow \infty$  limit. By the definition of free energy,

$$F'_b[\tilde{\Delta}_q] = \int_q \frac{\tilde{\Delta}_q \tilde{\Delta}_{-q}}{2u} - \frac{2}{\beta S} \text{tr} \ln [i\omega_n - \hat{H}'_b - \hat{V}'_b + \mu]. \quad (\text{B1})$$

To obtain the inverse nematic susceptibility, we expand the free energy up to second order in  $\tilde{\Delta}$ ,

$$F_b^{(2)} = \int_q \frac{|\tilde{\Delta}_q|^2}{2u} + \frac{1}{\beta} \sum_{k_z, k'_z} \int_{\mathbf{k}} \text{tr} (\mathcal{G}'_{k_z} \mathcal{V}'_{k_z-k'_z} \mathcal{G}'_{k'_z} \mathcal{V}'_{k'_z-k_z}),$$

where  $\mathcal{V}'_{k_z-k'_z}$  is given by Eq. (19), and



the denominator with these poles, and using the symmetry by exchanging  $X \leftrightarrow Y$ , we obtain

$$r'_{b,q} = \frac{1}{u} - 2 \int_{\mathbf{k}k_z} \text{Res} \left[ \frac{[(z - z_{k_z+\pi}^Y)(z - z_{k_z+q+\pi}^Y) - (2\lambda_{\mathbf{k}})^2 \sin(k_z + \pi) \sin(k_z + q + \pi)] \tanh \frac{\beta z}{2}}{(z - z_{k_z}^+)(z - z_{k_z}^-)(z - z_{k_z+q}^+)(z - z_{k_z+q}^-)} \right] \quad (\text{B2})$$

Applying the residue theorem straightforwardly, we finally obtain Eq. (22).

For the rest of this section, we evaluate the inverse susceptibility for some special cases. For circular, FS  $\delta_{\mathbf{k}} = 0$ , the susceptibility reduced to

$$r'_{b,q} = \frac{1}{u} - \int_{\mathbf{k},k_z} \frac{1}{4(t_{\mathbf{k}}^2 - \lambda_{\mathbf{k}}^2)} \left[ \frac{[\eta_{\mathbf{k}k_z}^2 + 4t_{\mathbf{k}}^2 \cos k_z \cos(k_z + q) - 4\lambda_{\mathbf{k}}^2 \sin k_z \sin(k_z + q)] n_-(k_z)}{\eta_{\mathbf{k}k_z} (\cos^2(k_z + q) - \cos^2 k_z)} \right. \\ \left. - \frac{[\eta_{\mathbf{k},k_z+q}^2 + 4t_{\mathbf{k}}^2 \cos k_z \cos(k_z + q) - 4\lambda_{\mathbf{k}}^2 \sin k_z \sin(k_z + q)] n_-(k_z + q)}{\eta_{\mathbf{k},k_z+q} (\cos^2(k_z + q) - \cos^2 k_z)} - \frac{2t_{\mathbf{k}}[n_+(k_z) - n_+(k_z + q)]}{\cos(k_z + q) - \cos k_z} \right] \quad (\text{B3})$$

where  $n_{\pm}(k_z) = \tanh(\beta z_{k_z}^{\pm}/2) \pm \tanh(\beta z_{k_z}^{\mp}/2)$ ,  $z_{k_z}^{\pm} = \mathbf{k}^2/(2m) - \mu \pm \eta_{\mathbf{k}k_z}$ , and  $\eta_{\mathbf{k}k_z} = 2(t_{\mathbf{k}}^2 \cos^2 k_z + \lambda_{\mathbf{k}}^2 \sin^2 k_z)^{1/2}$ . Furthermore, using equation (B2) with  $\delta_{\mathbf{k}} = 0$ , we have

$$r'_{b,q=0} = \frac{1}{u} - 2 \int_{\mathbf{k}k_z} \left[ \frac{2\lambda_{\mathbf{k}}^2 \sin^2 k_z \left( \tanh \frac{\beta z_{k_z}^+}{2} - \tanh \frac{\beta z_{k_z}^-}{2} \right)}{\eta_{\mathbf{k}k_z}^3} + \frac{\beta t_{\mathbf{k}}^2 \cos^2 k_z}{\eta_{\mathbf{k},k_z}^2} \left( \text{sech}^2 \frac{\beta z_{k_z}^+}{2} + \text{sech}^2 \frac{\beta z_{k_z}^-}{2} \right) \right. \\ \left. - \frac{t_{\mathbf{k}} \cos k_z \beta}{\eta_{\mathbf{k},k_z}} \frac{\beta}{2} \left( \text{sech}^2 \frac{\beta z_{k_z}^+}{2} - \text{sech}^2 \frac{\beta z_{k_z}^-}{2} \right) \right] \quad (\text{B4})$$

Similarly, for  $q = \pi$  with  $\delta_{\mathbf{k}} = 0$ ,

$$r'_{b,q=\pi} = \frac{1}{u} - 2 \int_{\mathbf{k}k_z} \left[ \frac{2t_{\mathbf{k}}^2 \cos^2 k_z}{\eta_{\mathbf{k}k_z}^3} \left( \tanh \frac{\beta z_{k_z}^+}{2} - \tanh \frac{\beta z_{k_z}^-}{2} \right) + \frac{\beta \lambda_{\mathbf{k}} \sin^2 k_z}{\eta_{\mathbf{k}k_z}^2} \left( \text{sech}^2 \frac{\beta z_{k_z}^+}{2} + \text{sech}^2 \frac{\beta z_{k_z}^-}{2} \right) \right] \quad (\text{B5})$$

### Appendix C: Tight-binding Green's function $\mathcal{G}_0^\alpha$

To find  $\mathcal{G}_0^\alpha$  in section III C 1, we first note that the eigenmodes of the Hamiltonian operator without hybridization is

$$c_{\alpha,p,\mathbf{k}\sigma} = \sqrt{\frac{2}{N+1}} \sum_{\ell} \sin \ell \vartheta_p c_{\alpha,\ell,\mathbf{k}\sigma}, \quad (\text{C1})$$

where  $\vartheta_p = \frac{p\pi}{N+1}$  with  $p = 1 \dots N$ . This can be checked by substituting the above equation into the Hamiltonian

$$\mathcal{H}_N^0 = \sum_{\ell=1}^N \sum_{\alpha,\mathbf{k}} \varepsilon_{\mathbf{k}}^\alpha c_{\alpha,\ell,\mathbf{k}\sigma}^\dagger c_{\alpha,\ell,\mathbf{k}\sigma} \\ - t_z \sum_{\ell=1}^{N-1} \sum_{\alpha,\mathbf{k}} c_{\alpha,\ell,\mathbf{k}\sigma}^\dagger c_{\alpha,\ell+1,\mathbf{k}\sigma} + h.c. \quad (\text{C2})$$

By using the orthogonal relation  $\sum_{\ell} \sin \ell \vartheta_p \sin \ell \vartheta_{p'} = \frac{2}{N+1} \delta_{pp'}$ , we obtain

$$\mathcal{H}_N^0 = \sum_{p=1}^N \sum_{\alpha,\mathbf{k}} (\varepsilon_{\mathbf{k}}^\alpha - 2t_z \cos \vartheta_p) c_{\alpha,p,\mathbf{k}\sigma}^\dagger c_{\alpha,p,\mathbf{k}\sigma}, \quad (\text{C3})$$

which is diagonal in ' $p$ '-basis.

Therefore, expanding the Green's function in this basis, this immediately lead to

$$[\mathcal{G}_0^\alpha]_{\ell\ell'} = \sum_p \frac{2/(N+1) \sin \ell \vartheta_p \sin \ell' \vartheta_p}{i\omega_n - \varepsilon_{\mathbf{k}}^\alpha + 2t_z \cos \vartheta_p + \mu}. \quad (\text{C4})$$

### Appendix D: Small- $t_z$ expansion of the inverse susceptibility $r_{\ell\ell'}$ for finite-size system with $t'_z = 0$

To make the expansion with respect to  $t_z$ , we start from the following presentation for the second term

in Eq. (24)

$$\frac{2}{\beta} \sum_{\omega_n} [\mathcal{G}_0^\alpha]_{\ell\ell'} [\mathcal{G}_0^\alpha]_{\ell'\ell} = - \sum_{pp'} \text{Res} \left[ \frac{S_{\ell\ell'}^p S_{\ell'\ell}^{p'} \tanh \frac{\beta z}{2}}{(z - z_p^\alpha)(z - z_{p'}^\alpha)} \right], \quad (\text{D1})$$

where the analytic-continuation technique has been used in the frequency summation, and  $z_p^\alpha = \xi_{\alpha, \mathbf{k}} - 2t_z \cos \vartheta_p$  with  $\alpha = X, Y$ ,  $\xi_{X, \mathbf{k}} = k_x^2/(2m_x) + k_y^2/(2m_y) - \mu$ ,  $\xi_{Y, \mathbf{k}} = k_x^2/(2m_x) + k_y^2/(2m_y) - \mu$ , and  $S_{\ell\ell'}^p = 2 \sin(\ell\vartheta_p) \sin(\ell'\vartheta_p)/(N+1)$ . Using the orthogonality relation  $\sum_{p=1}^N \sin(\ell\vartheta_p) \sin(\ell'\vartheta_p) = \delta_{\ell\ell'}(N+1)/2$ , the expansion of Eq. (D1) is

$$\text{Res} \left[ \frac{\tanh \frac{\beta z}{2}}{(z - \xi_{\mathbf{k}})^2} \left( \delta_{\ell\ell'} + \frac{2t_z \delta_{\ell\ell'} (\delta_{\ell, \ell'+1} + \delta_{\ell, \ell'-1})}{(z - \xi_{\mathbf{k}})} \right) + t_z^2 \frac{(\delta_{\ell, \ell'+1} + \delta_{\ell, \ell'-1})^2}{(z - \xi_{\mathbf{k}})^2} + 4t_z^2 \frac{(1 - \frac{\delta_{\ell, 1} + \delta_{\ell, N}}{2}) \delta_{\ell\ell'}}{(z - \xi_{\mathbf{k}})^2} \right],$$

where we used

$$\sum_p^N S_{\ell\ell'}^p \cos \vartheta_p = \frac{\delta_{\ell, \ell'+1} + \delta_{\ell, \ell'-1}}{2},$$

$$\sum_p^N S_{\ell\ell}^p \cos^2 \vartheta_p = 1 - \frac{\delta_{\ell, 1} + \delta_{\ell, N}}{2}.$$

Applying the residue theorem, we immediately obtain the approximation of  $r_{\ell\ell'}$  for small  $t_z$ ,

$$r_{\ell\ell'} \simeq \left[ \frac{1}{u} - \beta \int_{\mathbf{k}} \text{sech}^2 \frac{\beta \xi_{\mathbf{k}}}{2} \right] \delta_{\ell\ell'} - \frac{\beta^3 t_z^2}{3} \int_{\mathbf{k}} \left[ \text{sech}^2 \frac{\beta \xi_{\mathbf{k}}}{2} (3 \tanh^2 \frac{\beta \xi_{\mathbf{k}}}{2} - 1) \right] \times \left[ \delta_{\ell\ell'} - \frac{\delta_{\ell, 1} + \delta_{\ell, N}}{2} \delta_{\ell\ell'} + \frac{\delta_{\ell, \ell'+1} + \delta_{\ell, \ell'-1}}{4} \right] \quad (\text{D2})$$

Integration over the in-plane momentum  $\mathbf{k}$  using  $\int_{\mathbf{k}} \rightarrow (\tilde{m}/2\pi) \int d\xi_{\mathbf{k}}$  gives the result (27) of the main text.

### Appendix E: Perturbation calculation for the case when $N$ is odd and $m_x = m_y$

For  $N$  is odd with isotropic FS ( $m_x = m_y$ ), the matrix in the first term of Eq. (29) is degenerate. In this case, Eq. (31) cannot be used directly due to zero denominators in some terms. Before computing the eigenvalue of  $\hat{H}$  perturbatively, we first rotate the eigenspace of  $\hat{H}$  matrix in  $F'_N[\Delta_\ell]$  by

$$R = \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbb{I} & -\mathbb{I} \\ \mathbb{I} & \mathbb{I} \end{pmatrix}$$

where  $\mathbb{I}$  is a  $N \times N$  identity matrix. Namely,

$$F'_N[\Delta_\ell] = \sum_\ell \frac{\Delta_\ell^2}{2u} + \frac{2}{S\beta} \text{tr}[R^{-1} \ln(i\omega_n - \hat{H})R]$$

$$= \sum_\ell \frac{\Delta_\ell^2}{2u} + \frac{2}{S\beta} \text{tr}[\ln(i\omega_n - R^{-1} \hat{H} R)].$$

Therefore, we obtain

$$R^{-1} \hat{H} R = \begin{pmatrix} \mathbf{E}_{\mathbf{k}} + \mathbf{P}_{\mathbf{k}} & 0 \\ 0 & \mathbf{E}_{\mathbf{k}} - \mathbf{P}_{\mathbf{k}} \end{pmatrix} + \begin{pmatrix} 0 & -\mathcal{V} \\ -\mathcal{V} & 0 \end{pmatrix},$$

where  $\mathbf{E}_{\mathbf{k}} = \mathbf{E}_{\mathbf{k}}^X = \mathbf{E}_{\mathbf{k}}^Y$  is given by equation (30a) (the  $X$  and  $Y$  pockets are identical),  $\mathbf{P}_{\mathbf{k}}$  is defined by Eq. (30b), and  $[\mathcal{V}]_{\ell\ell'} = \Delta_\ell \delta_{\ell\ell'}$ . The first term in  $R^{-1} \hat{H} R$  becomes block-diagonalized and more convenient for perturbation calculation. The second term in  $R^{-1} \hat{H} R$  is treated as perturbation.

Now, we let  $\tilde{\mathbf{x}}_{p, \mathbf{k}}^T = [\tilde{x}_{p, \mathbf{k}}^1, \dots, \tilde{x}_{p, \mathbf{k}}^N]$  and  $\tilde{\mathbf{y}}_{p, \mathbf{k}}^T = [\tilde{y}_{p, \mathbf{k}}^1, \dots, \tilde{y}_{p, \mathbf{k}}^N]$  with  $p = 1 \dots N$ , and they satisfies

$$(\mathbf{E}_{\mathbf{k}} + \mathbf{P}_{\mathbf{k}}) \tilde{\mathbf{x}}_{p, \mathbf{k}} = f_{p, \mathbf{k}}^{(0)} \tilde{\mathbf{x}}_{p, \mathbf{k}},$$

$$(\mathbf{E}_{\mathbf{k}} - \mathbf{P}_{\mathbf{k}}) \tilde{\mathbf{y}}_{p, \mathbf{k}} = f_{p, \mathbf{k}}^{(0)} \tilde{\mathbf{y}}_{p, \mathbf{k}},$$

where  $f_{p, \mathbf{k}}^{(0)}$  is the unperturbed eigenvalue.

Furthermore, in order to make connection with the standard notation in quantum mechanics perturbation theory, we introduce the following ‘bra’ and ‘ket’ notation.

$$|\tilde{\mathbf{x}}_p\rangle = \begin{pmatrix} \tilde{\mathbf{x}}_{p, \mathbf{k}} \\ \mathbf{0} \end{pmatrix}, \quad |\tilde{\mathbf{y}}_p\rangle = \begin{pmatrix} \mathbf{0} \\ \tilde{\mathbf{y}}_{p, \mathbf{k}} \end{pmatrix}, \quad (\text{E1})$$

where  $\mathbf{0}$  is a  $1 \times N$  zero matrix. Thus,  $|\tilde{\mathbf{x}}_p\rangle$  and  $|\tilde{\mathbf{y}}_p\rangle$  span the  $p$ -th 2-fold degenerate subspace. Also, we set the perturbation operator as

$$\hat{V} = \begin{pmatrix} 0 & -\mathcal{V} \\ -\mathcal{V} & 0 \end{pmatrix}. \quad (\text{E2})$$

One can immediately see that, any linear combination of  $|\tilde{\mathbf{x}}_p\rangle$  and  $|\tilde{\mathbf{y}}_p\rangle$  are still the eigenvector of the first term in  $R^{-1} \hat{H} R$ . Therefore, the choices of eigenvector are not unique. Exploiting this fact, we can choose a basis such that the numerators with overlapping degenerate eigenvectors in Eq. (31) vanish. Hence, the zero denominator terms are dropped out in the calculation. The procedure to obtain such basis is as follows.

First, we calculate the first-order correction for the eigenvalue in the  $p$ -th degenerate subspace. In

this subspace, the operator  $\hat{V}$  can be represented as the following matrix form,

$$\begin{pmatrix} \langle \tilde{\mathbf{x}}_p | \hat{V} | \tilde{\mathbf{x}}_p \rangle & \langle \tilde{\mathbf{x}}_p | \hat{V} | \tilde{\mathbf{y}}_p \rangle \\ \langle \tilde{\mathbf{y}}_p | \hat{V} | \tilde{\mathbf{x}}_p \rangle & \langle \tilde{\mathbf{y}}_p | \hat{V} | \tilde{\mathbf{y}}_p \rangle \end{pmatrix} = \begin{pmatrix} 0 & \tilde{\mathbf{x}}_{p,\mathbf{k}}^T \mathcal{V} \tilde{\mathbf{y}}_{p,\mathbf{k}} \\ \tilde{\mathbf{y}}_{p,\mathbf{k}}^T \mathcal{V} \tilde{\mathbf{x}}_{p,\mathbf{k}} & 0 \end{pmatrix}.$$

Solving the eigenvalues of the above  $2 \times 2$  matrix yields the first-order correction. Writing out the nematic order parameters explicitly, this matrix becomes

$$\sum_{\ell} \begin{pmatrix} 0 & \tilde{x}_{p,\mathbf{k}}^{\ell} \tilde{y}_{p,\mathbf{k}}^{\ell} \Delta_{\ell} \\ \tilde{x}_{p,\mathbf{k}}^{\ell} \tilde{y}_{p,\mathbf{k}}^{\ell} \Delta_{\ell} & 0 \end{pmatrix}, \quad (\text{E3})$$

and has the following eigenvalues and eigenvectors

$$\pm \sum_{\ell} \tilde{x}_{p,\mathbf{k}}^{\ell} \tilde{y}_{p,\mathbf{k}}^{\ell} \Delta_{\ell}, \text{ and } \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ \pm 1 \end{pmatrix}. \quad (\text{E4})$$

These results yield the desirable ‘rotated’ eigenvectors

$$|p, \pm\rangle = \frac{1}{\sqrt{2}} (|\tilde{\mathbf{x}}_p\rangle \pm |\tilde{\mathbf{y}}_p\rangle) \quad (\text{E5})$$

with the first-order corrected eigenvalue

$$\tilde{f}_{\mathbf{k}}^{(p,\pm)} \simeq \tilde{f}_{p,\mathbf{k}}^{(0)} \pm \sum_{\ell} \tilde{a}_{\ell}^{(p,\pm)} \Delta_{\ell} + \mathcal{O}(\Delta^2), \quad (\text{E6})$$

where  $\tilde{a}_{\ell}^{(p,\pm)} = \pm \tilde{x}_{p,\mathbf{k}}^{\ell} \tilde{y}_{p,\mathbf{k}}^{\ell}$ . The first order correction has lifted the degeneracy and single out the particular choice of linear combination:  $|p, \pm\rangle$ . Also, note that, only this choice can be smoothly approached from the perturbed eigenvectors when the perturbations are turning off.

Further, the second-order corrected eigenvalues for eigenvectors  $|p, \pm\rangle$ , Eq. (E5), are evaluated as

$$\sum_{p' \neq p} \frac{|\langle p, \pm | \hat{V} | p', + \rangle|^2 + |\langle p, \pm | \hat{V} | p', - \rangle|^2}{f_{p,\mathbf{k}}^{(0)} - f_{p',\mathbf{k}}^{(0)}}. \quad (\text{E7})$$

Note that, in the summation, not only the terms  $\langle p, \pm | \hat{V} | p, \pm \rangle$  are excluded, but also  $\langle p, - | \hat{V} | p, + \rangle$  and  $\langle p, + | \hat{V} | p, - \rangle$  (also having zero denominator), since they vanish in the rotated new basis.

Therefore, the approximation of the  $2N$  eigenvalues up to second order is

$$\tilde{f}_{p,\pm,\mathbf{k}} \simeq \tilde{f}_{p,\mathbf{k}}^{(0)} + \sum_{\ell} \tilde{a}_{\ell}^{p,\pm} \Delta_{\ell} + \sum_{\ell\ell'} \tilde{b}_{\ell\ell'}^{p,\pm} \Delta_{\ell} \Delta_{\ell'} \quad (\text{E8})$$

with

$$\tilde{b}_{\ell\ell'}^{p,\pm} = \sum_{p' \neq p} \frac{\tilde{y}_{p,\mathbf{k}}^{\ell} \tilde{x}_{p',\mathbf{k}}^{\ell} \tilde{y}_{p,\mathbf{k}}^{\ell'} \tilde{x}_{p',\mathbf{k}}^{\ell'} + \tilde{x}_{p,\mathbf{k}}^{\ell} \tilde{y}_{p',\mathbf{k}}^{\ell} \tilde{x}_{p,\mathbf{k}}^{\ell'} \tilde{y}_{p',\mathbf{k}}^{\ell'}}{2(f_{p,\mathbf{k}}^{(0)} - f_{p',\mathbf{k}}^{(0)})}.$$

Using (E8), we expand the free energy near  $\Delta_{\ell} \simeq 0$  and obtain

$$\begin{aligned} r'_{\ell\ell'} &= \frac{\delta_{\ell\ell'}}{u} - \sum_p \int_{\mathbf{k}} \left[ \frac{\beta}{2} (a_{\ell}^{p,+} a_{\ell'}^{p,+} + a_{\ell}^{p,-} a_{\ell'}^{p,-}) \right. \\ &\quad \left. \times \text{sech}^2\left(\frac{\beta z_p}{2}\right) + 2(b_{\ell\ell'}^{p,+} + b_{\ell\ell'}^{p,-}) \tanh\left(\frac{\beta z_p}{2}\right) \right], \end{aligned} \quad (\text{E9})$$

where  $z_p = f_{p,\mathbf{k}}^{(0)} - \mu$ .

<sup>1</sup> J. Paglione and R. L. Greene, Nat. Phys. **6**, 645 (2010).  
<sup>2</sup> G. R. Stewart, Rev. Mod. Phys. **83**, 1589 (2011).  
<sup>3</sup> H. Hosono and K. Kuroki, Physica C **514**, 399 (2015).  
<sup>4</sup> S. Avci, O. Chmaissem, D. Y. Chung, S. Rosenkranz, E. A. Goremychkin, J. P. Castellan, I. S. Todorov, J. A. Schlueter, H. Claus, A. Daoud-Aladine, D. D. Khalyavin, M. G. Kanatzidis, and R. Osborn, Phys. Rev. B **85**, 184507 (2012).  
<sup>5</sup> M. G. Kim, R. M. Fernandes, A. Kreyssig, J. W. Kim, A. Thaler, S. L. Bud'ko, P. C. Canfield, R. J. McQueeney, J. Schmalian, and A. I. Goldman, Phys. Rev. B **83**, 134522 (2011).  
<sup>6</sup> C. R. Rotundu and R. J. Birgeneau, Phys. Rev. B **84**, 092501 (2011).  
<sup>7</sup> J. Zhao, Q. Huang, C. De La Cruz, S. Li, J. Lynn,

Y. Chen, M. Green, G. Chen, G. Li, Z. Li, *et al.*, Nat. Mater. **7**, 953 (2008).

<sup>8</sup> Q. Huang, J. Zhao, J. W. Lynn, G. F. Chen, J. L. Luo, N. L. Wang, and P. Dai, Phys. Rev. B **78**, 054529 (2008).

<sup>9</sup> H. Luetkens, H.-H. Klauss, M. Kraken, F. Litterst, T. Dellmann, R. Klingeler, C. Hess, R. Khasanov, A. Amato, C. Baines, *et al.*, Nat. Mater. **8**, 305 (2009).

<sup>10</sup> C. R. Rotundu, D. T. Keane, B. Freelon, S. D. Wilson, A. Kim, P. N. Valdivia, E. Bourret-Courchesne, and R. J. Birgeneau, Phys. Rev. B **80**, 144517 (2009).

<sup>11</sup> A. Jesche, C. Krellner, M. de Souza, M. Lang, and C. Geibel, Phys. Rev. B **81**, 134525 (2010).

<sup>12</sup> G. F. Chen, W. Z. Hu, J. L. Luo, and N. L. Wang, Phys. Rev. Lett. **102**, 227004 (2009).

<sup>13</sup> S. Li, C. de la Cruz, Q. Huang, G. F. Chen, T.-L.

- Xia, J. L. Luo, N. L. Wang, and P. Dai, Phys. Rev. B **80**, 020504 (2009).
- <sup>14</sup> C. Fang, H. Yao, W.-F. Tsai, J. P. Hu, and S. A. Kivelson, Phys. Rev. B **77**, 224509 (2008).
- <sup>15</sup> C. Xu, M. Müller, and S. Sachdev, Phys. Rev. B **78**, 020501 (2008).
- <sup>16</sup> R. M. Fernandes, A. V. Chubukov, J. Knolle, I. Eremin, and J. Schmalian, Phys. Rev. B **85**, 024534 (2012).
- <sup>17</sup> L. Fanfarillo, A. Cortijo, and B. Valenzuela, Phys. Rev. B **91**, 214515 (2015).
- <sup>18</sup> F. Krüger, S. Kumar, J. Zaanen, and J. van den Brink, Phys. Rev. B **79**, 054504 (2009).
- <sup>19</sup> W. Lv, J. Wu, and P. Phillips, Phys. Rev. B **80**, 224506 (2009).
- <sup>20</sup> C.-C. Lee, W.-G. Yin, and W. Ku, Phys. Rev. Lett. **103**, 267001 (2009).
- <sup>21</sup> C.-C. Chen, J. Maciejko, A. P. Sorini, B. Moritz, R. R. P. Singh, and T. P. Devereaux, Phys. Rev. B **82**, 100504 (2010).
- <sup>22</sup> E. Fradkin, S. A. Kivelson, M. J. Lawler, J. P. Eisenstein, and A. P. Mackenzie, Annu. Rev. Cond. Mat. Phys. **1**, 153 (2010).
- <sup>23</sup> J.-H. Chu, J. G. Analytis, K. De Greve, P. L. McMahon, Z. Islam, Y. Yamamoto, and I. R. Fisher, Science **329**, 824 (2010); J.-H. Chu, H.-H. Kuo, J. G. Analytis, and I. R. Fisher, Science **337**, 710 (2012).
- <sup>24</sup> M. A. Tanatar, E. C. Blomberg, A. Kreyssig, M. G. Kim, N. Ni, A. Thaler, S. L. Bud'ko, P. C. Canfield, A. I. Goldman, I. I. Mazin, and R. Prozorov, Phys. Rev. B **81**, 184508 (2010).
- <sup>25</sup> R. M. Fernandes, L. H. VanBebber, S. Bhattacharya, P. Chandra, V. Keppens, D. Mandrus, M. A. McGuire, B. C. Sales, A. S. Sefat, and J. Schmalian, Phys. Rev. Lett. **105**, 157003 (2010).
- <sup>26</sup> M. Yoshizawa, D. Kimura, T. Chiba, S. Simayi, Y. Nakanishi, K. Kihou, C.-H. Lee, A. Iyo, H. Eisaki, M. Nakajima, and S. ichi Uchida, Journal of the Physical Society of Japan **81**, 024604 (2012).
- <sup>27</sup> A. E. Böhrner, P. Burger, F. Hardy, T. Wolf, P. Schweiss, R. Fromknecht, M. Reinecker, W. Schranz, and C. Meingast, Phys. Rev. Lett. **112**, 047001 (2014).
- <sup>28</sup> A. Dusza, A. Lucarelli, F. Pfuner, J.-H. Chu, I. R. Fisher, and L. Degiorgi, Europhys. Lett. **93**, 37002 (2011); M. Nakajima, S. Ishida, Y. Tomioka, K. Kihou, C. H. Lee, A. Iyo, T. Ito, T. Kakeshita, H. Eisaki, and S. Uchida, Phys. Rev. Lett. **109**, 217003 (2012).
- <sup>29</sup> M. Yi, D. Lu, J.-H. Chu, J. G. Analytis, A. P. Sorini, A. F. Kemper, B. Moritz, S.-K. Mo, R. G. Moore, M. Hashimoto, W.-S. Lee, Z. Hussain, T. P. Devereaux, I. R. Fisher, and Z.-X. Shen, Proc. Natl. Acad. Sci. U.S.A. **108**, 6878 (2011).
- <sup>30</sup> R. M. Fernandes, A. V. Chubukov, and J. Schmalian, Nat. Phys. **10**, 97 (2014).
- <sup>31</sup> A. V. Chubukov, M. Khodas, and R. M. Fernandes, arXiv:1602.05503 (2016).
- <sup>32</sup> S. Kasahara, H. J. Shi, K. Hashimoto, S. Tonegawa, Y. Mizukami, T. Shibauchi, K. Sugimoto, T. Fukuda, T. Terashima, A. H. Nevidomskyy, and Y. Matsuda, Nature **486**, 382 (2012).
- <sup>33</sup> T. Iye, M.-H. Julien, H. Mayaffre, M. Horvatić, C. Berthier, K. Ishida, H. Ikeda, S. Kasahara, T. Shibauchi, and Y. Matsuda, J. Phys. Soc. Jpn. **84**, 043705 (2015).
- <sup>34</sup> E. Thewalt, J. P. Hinton, I. M. Hayes, T. Helm, D. H. Lee, J. G. Analytis, and J. Orenstein, arXiv:1507.03981 (2015).
- <sup>35</sup> X. Luo, V. Stanev, B. Shen, L. Fang, X. S. Ling, R. Osborn, S. Rosenkranz, T. M. Benseman, R. Divan, W.-K. Kwok, and U. Welp, Phys. Rev. B **91**, 094512 (2015).
- <sup>36</sup> S. Liang, A. Moreo, and E. Dagotto, Phys. Rev. Lett. **111**, 047004 (2013).
- <sup>37</sup> We note that in the case of continuous transition we consider, the surface instability actually smears the bulk transition in finite-size samples. With decreasing temperature the nematic order smoothly extends over larger distances away from the surface. In macroscopic samples, however, the bulk transition becomes a very sharp crossover occurring at the bulk transition temperature, the larger sample the sharper crossover. In contrast, in the case of first-order bulk transition, not considered in this paper, there should be two distinct phase transitions.
- <sup>38</sup> H. Yamase and R. Zeyher, Phys. Rev. B **88**, 180502 (2013).
- <sup>39</sup> A. V. Chubukov, D. V. Efremov, and I. Eremin, Phys. Rev. B **78**, 134512 (2008).
- <sup>40</sup> S. Graser, T. A. Maier, P. J. Hirschfeld, and D. J. Scalapino, New J. Phys. **11**, 025016 (2009).
- <sup>41</sup> E. Fradkin, *Field Theories of Condensed Matter Systems* (Perseus Books, 1991).
- <sup>42</sup> H. Yamase and H. Kohno, J. Phys. Soc. Jpn. **69**, 332 (2000); J. Phys. Soc. Jpn. **69**, 2151 (2000).
- <sup>43</sup> C. J. Halboth and W. Metzner, Phys. Rev. Lett. **85**, 5162 (2000).
- <sup>44</sup> H. Yamase, V. Oganesyan, and W. Metzner, Phys. Rev. B **72**, 035114 (2005).
- <sup>45</sup> H.-Y. Kee, E. H. Kim, and C.-H. Chung, Phys. Rev. B **68**, 245109 (2003).
- <sup>46</sup> I. Khavkine, C.-H. Chung, V. Oganesyan, and H.-Y. Kee, Phys. Rev. B **70**, 155110 (2004).
- <sup>47</sup> P. T. Dumitrescu, M. Serbyn, R. T. Scalettar, and A. Vishwanath, arXiv:1512.08523 (2015).
- <sup>48</sup> A. V. Chubukov and R.-Q. Xing, Phys. Rev. B **93**, 165141 (2016).
- <sup>49</sup> M. Khodas and A. V. Chubukov, Phys. Rev. B **86**, 144519 (2012).
- <sup>50</sup> H. Yamase, Phys. Rev. Lett. **102**, 116404 (2009).
- <sup>51</sup> G. Y. Hu and R. F. O'Connell, Journ. Phys. A: Mathematical and General **29**, 1511 (1996).
- <sup>52</sup> One can demonstrate that if the order parameter is suppressed at the surface then at the transition point the ratio  $\Delta_{N/2}/\Delta_1$  grows linearly with the system size  $N$ . We can expect that below the transition temperature this ratio becomes size independent at large  $N$ .
- <sup>53</sup> Note, however, that in the vicinity of  $\bar{g}_{c1}$  the decay of

the order parameter is not well described by a simple exponent, Eq. (38).