



CHORUS

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

Recovery of massless Dirac fermions at charge neutrality in strongly interacting twisted bilayer graphene with disorder

Alex Thomson and Jason Alicea

Phys. Rev. B **103**, 125138 — Published 17 March 2021

DOI: [10.1103/PhysRevB.103.125138](https://doi.org/10.1103/PhysRevB.103.125138)

Recovery of massless Dirac fermions at charge neutrality in strongly interacting twisted bilayer graphene with disorder

Alex Thomson^{1,2,3} and Jason Alicea^{1,2,3}

¹*Institute for Quantum Information and Matter, California Institute of Technology, Pasadena, California 91125, USA*

²*Walter Burke Institute for Theoretical Physics, California Institute of Technology, Pasadena, California 91125, USA*

³*Department of Physics, California Institute of Technology, Pasadena, California 91125, USA*

(Dated: October 23, 2020)

Stacking two graphene layers twisted by the ‘magic angle’ $\theta \approx 1.1^\circ$ generates flat energy bands, which in turn catalyzes various strongly correlated phenomena depending on filling and sample details. At charge neutrality, transport measurements reveal superficially mundane semimetallicity (as expected when correlations are weak) in some samples yet robust insulation in others. We propose that the interplay between interactions and disorder admits either behavior, *even when the system is strongly correlated and locally gapped*. Specifically, we argue that strong interactions supplemented by weak, smooth disorder stabilize a network of gapped quantum valley Hall domains with spatially varying Chern numbers determined by the disorder landscape — even when an entirely different order is favored in the clean limit. Within this scenario, sufficiently small samples that realize a single domain display insulating transport characteristics. Conversely, multi-domain samples exhibit re-emergent massless Dirac fermions formed by gapless domain-wall modes, yielding semimetallic behavior except on the ultra-long scales at which localization becomes visible. We discuss experimental tests of this proposal via local probes and transport. Our results highlight the crucial role that randomness can play in ground-state selection of twisted heterostructures, an observation that we expect to have further ramifications at other fillings.

I. INTRODUCTION

The discovery of superconductivity and correlated insulators in magic-angle twisted bilayer graphene (mTBG) [1, 2] opened a fascinating new chapter in the field of strongly interacting quantum matter. The ‘magic’ stems from the fact that upon twisting the two graphene layers by an angle $\theta \approx 1.1^\circ$ from one another, the bands immediately above and below the charge neutrality point become exceptionally flat [3, 4] —bringing interactions center stage. Accounting for spin and valley degrees of freedom, each of these two flat bands is essentially fourfold degenerate. Correlated physics, including superconductivity, thus naturally arises when the number of charge carriers per moiré unit cell is between $\nu = -4$ (four holes) and $\nu = +4$ (four electrons).

The observed phenomenology of mTBG depends sensitively on sample details. Cao *et al.* [1, 2] originally observed correlated insulating states at $\nu = \pm 2$ along with superconducting domes upon doping away from the $\nu = -2$ insulator. Near the charge neutrality point at $\nu = 0$, the conductance exhibited a V-shaped suppression indicative of semimetallicity. Non-interacting band theory calculations [4] predict massless Dirac fermions at charge neutrality provided the system preserves $C_2\mathcal{T}$ symmetry, with C_2 a two-fold rotation and \mathcal{T} time reversal [5–8]; the latter observation thus at first sight suggests weak correlations at $\nu = 0$. The magic-angle device examined by Yankowitz *et al.* [9] additionally exhibited superconductivity adjacent to the $\nu = +2$ insulator and a resistive correlated state at $\nu = +3$. Near charge neutrality, transport again appeared consistent with the semimetallic behavior expected from band theory.

A second class of mTBG systems arises upon aligning the hexagonal boron nitride (hBN) substrate with one of the

graphene sheets [10, 11]. The alignment appears to underlie strikingly different correlated physics: an absence of superconductivity, removal of the $\nu = -2$ insulator, weak resistive peaks at $\nu = +2$ instead of robust insulation, and a quantum anomalous Hall state at $\nu = +3$. Furthermore, at charge neutrality the system becomes strongly insulating instead of semimetallic. The behavior at charge neutrality is, however, yet again consistent with band theory. Indeed, alignment-induced breaking of C_2 symmetry renders the Dirac fermions massive, yielding a band gap at $\nu = 0$. Explicit C_2 breaking has also been proposed as a catalyst for the observed quantum anomalous Hall state [12, 13].

Still different phenomenology emerges in the ultra-homogeneous samples studied by Lu *et al.* [14]. These samples featured resistive peaks evincing either well-developed or incipient insulators at *all* integer fillings $\nu = 0, \pm 1, \pm 2, \pm 3$, as well as additional superconducting domes beyond those reported previously. Notably, the strongest insulating state within the flat-band manifold occurred at charge neutrality, naively suggesting alignment with the hBN substrate as in Refs. 10 and 11. Several factors challenge this interpretation, however. First, Lu *et al.* make no attempt to align the hBN, and it is unlikely to occur at random. Second, hBN-aligned samples and those of Lu *et al.* realize a largely disparate set of phenomena, suggesting against a common microscopic origin. Finally, the gap reported by Lu *et al.* [14] dwarfs by roughly an order of magnitude that measured in hBN-aligned mTBG [11], making its formation by explicit symmetry-breaking seem unlikely in comparison. Thus the insulating behavior observed at all of the fillings indicated above—including $\nu = 0$ —seems most naturally rooted in strong correlations.

A conservative interpretation of the available charge-

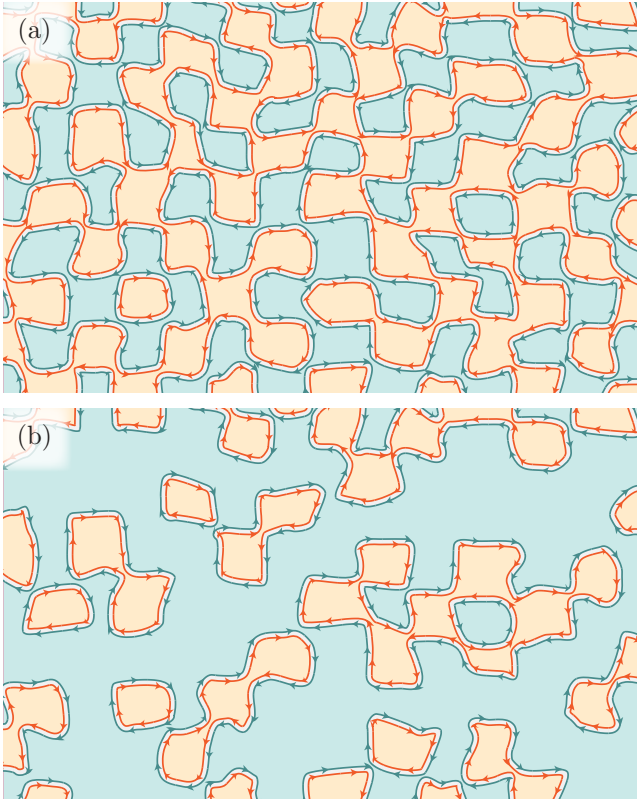


FIG. 1. Random tiling of quantum valley Hall states in a system where C_2 is (a) preserved on average and (b) explicitly broken (*e.g.*, by hBN alignment). For simplicity we show the domain structure only for a single valley. Blue regions carry Chern number $C = +1$ for both spins, whereas orange regions carry Chern number -1 . The arrows represent the two chiral edge modes (per spin) that traverse the domain boundaries. The domain structure corresponding to the valley-sector not depicted here is simply obtained by exchanging the colours and reversing edge modes in (a) and (b).

neutrality transport data is that greater inhomogeneity in the samples from Refs. 1, 2, and 9 merely obliterates the strong correlations operative at $\nu = 0$ in the Lu *et al.* samples. Such a viewpoint is supported by the fact that significant “twist-angle disorder” has been observed by multiple groups [1, 2, 9, 15–19]; moreover, deviations from the magic angle locally enhance the flat-band dispersion [20], potentially diminishing correlation effects. Scanning tunneling microscopy (STM) measurements from Refs. 15–18, however, do not simply fit this picture. All of these STM studies observed *local* correlation effects at charge neutrality, manifested by a pronounced splitting of the flat-band van Hove peaks upon approaching $\nu = 0$ and, in Ref. 18, evidence of a hard gap at charge neutrality¹. Much subtler signatures of correlated states were also seen at other integer fillings

(typically most prominently at $\nu = +2$). From a local perspective, it therefore appears that correlations in the STM samples are actually *strongest* at $\nu = 0$.

In this paper we propose a unifying explanation for the diverse phenomenology observed to date in mTBG at charge neutrality. Our scenario posits that strong correlations are ubiquitous—even in samples that observe semimetallic behavior expected from band theory—with disorder playing a secondary but still crucial role. We specifically assume that in a perfectly clean infinite system, interactions favor, or very nearly favor, correlated states that spontaneously break $C_2\mathcal{T}$ symmetry in a way that yields Chern number $C = \pm 1$ for a given spin/valley sector. This assumption is bolstered by existing numerical simulations [14, 16, 21, 22] and justified further below. Among the many possible insulators, only two preserve translation symmetry, spin rotation symmetry, and time reversal: the pair of ‘quantum valley Hall’ states [23–31] with $C = +1$ for both spins in one valley and $C = -1$ for both spins in the other valley, or vice versa. Note that C_2 transforms the quantum valley Hall states into one another; hence they are exactly degenerate provided C_2 is not explicitly broken.

Imagine now turning on smooth, non-magnetic disorder (arising from twist-angle inhomogeneity, strain, etc.) that explicitly violates the infinite system’s C_2 symmetry but preserves it in an average sense. Within the manifold specified above, quantum valley Hall states are unique in that their order parameter directly couples to the disorder potential—allowing the system to efficiently gain energy by locally forming one of those two phases. We further assume that the energy gain outweighs any energy cost (should one exist) for forming quantum valley Hall order in the clean limit. Under these circumstances the infinite system exhibits a random tiling of the two quantum valley Hall states, details of which are determined by the interplay between interactions and the disorder landscape; see Fig. 1(a) for an illustration. Similar domain structures have been discussed in several other contexts, *e.g.*, in systems with valley Hall nematic order [32, 33] or as a source of non-Abelian ‘PH-Pfaffian’ topological order [34–36].

Crucially, the infinite system is locally gapped within the quantum valley Hall domains but is not entirely electrically inert. Each domain wall binds four ‘right-moving’ and four ‘left-moving’ charge-carrying modes, reflecting the fact that the Chern numbers for the spin/valley sectors change by ± 2 upon passing between adjacent domains. Smoothness of the disorder potential suppresses scattering among these modes and thus justifies treating the spin/valley sectors as decoupled (to a first approximation). In this limit the system realizes four copies of a Chalker-Coddington network model [37] describing an integer-quantum-Hall plateau transition at

¹ The issue of hBN alignment is subtle given the different nature of these experiments. Nevertheless, alignment is expected to enhance the van

Hove peak splitting *independent of filling*, whereas the splitting observed by STM is significantly larger at charge neutrality compared to when the bands are fully filled.

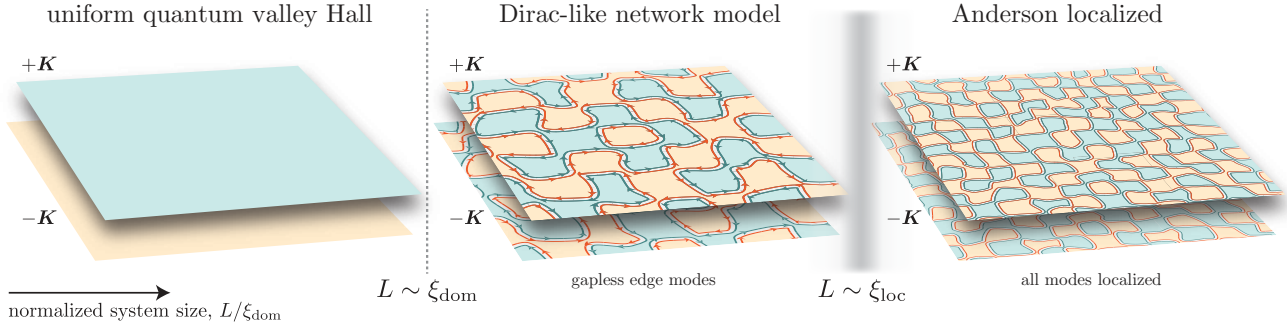


FIG. 2. Phase diagram versus system size L for a given disorder landscape. Left panel: With L below the typical domain size ξ_{dom} , a single domain is realized. Here each valley exhibits Chern number $C = +1$ (depicted in blue) or $C = -1$ (depicted in orange) throughout the entire sample, yielding insulating transport as observed by Lu *et al.* [14]. Central panel: When $L \gtrsim \xi_{\text{dom}}$ multiple domains may be present within a single sample. The domain structure results in a percolating network of gapless edge modes that underlies the Dirac-like conductance seen by Cao *et al.* [1, 2] and Yankowitz *et al.* [9]. Right panel: When L exceeds the localization length ξ_{loc} the sample localizes and ceases to conduct.

which the Hall conductivity changes by $\delta\sigma_{xy} = \pm 2e^2/h$ [38, 39]. This set of plateau transitions can be described by eight massless Dirac fermions (two per sector) with disorder acting within each cone [40, 41]. Hence essentially the same low-energy physics expected from band theory emerges from a strongly correlated framework! Residual scattering among the domain-wall modes generates inter-cone disorder that produces localization, but the localization lengths can be arbitrarily long.

We emphasize that strong correlations form the bedrock of the scenario outlined above. Without interactions and in the absence of explicit net $C_2\mathcal{T}$ -breaking, the fate of the system depends sensitively on the nonuniversal details of the disorder. The quantum valley Hall state is but one among many potential phases, both gapped and gapless, that disorder could locally favour. Moreover, even if local quantum valley Hall order happened to develop, the gap would be set by disorder and could be exceedingly small. In contrast, by inducing the spontaneous breaking of $C_2\mathcal{T}$, interactions select a small subset of energetically competitive states in our scenario. Disorder then plays a subordinate role by favoring one of the two quantum valley Hall orders in that set, thereby generating the domain structure. The local gap protecting the insulating domains is determined primarily by interactions rather than disorder.

Let us now revisit experiments in light of our proposed picture. Locally probing the quantum valley Hall domains should reveal signatures of a correlation-driven gapped spectrum (possibly dressed with disorder-induced subgap states), consistent with STM experiments [15–18]. The outcome of global transport measurements depends on the ratio of sample size L to the typical domain size ξ_{dom} that would occur in an infinite system. For homogeneous systems such that $L/\xi_{\text{dom}} \lesssim 1$, transport probes essentially a single domain, yielding insulating behavior as observed by Lu *et al.* [14]. (Strong intervalley scattering induced by the sample boundary is expected to suppress edge conduction.) Con-

versely, for more-disordered samples with $L/\xi_{\text{dom}} \gg 1$, transport probes many domains; here the massless Dirac fermions emerging from the gapless domain walls underpin semimetallic conduction as measured in Refs. 1, 2, and 9. See Fig. 2 for a summary. We can also make contact with alignment-induced insulation observed in Refs. 10 and 11. Turning on hBN alignment supplements the disorder landscape with a *uniform* C_2 -breaking potential that shrinks the area occupied by one of the quantum valley Hall states and expands the area of the other, as shown in Fig. 1(b). Domain walls then no longer percolate, thereby gapping the re-emergent massless Dirac fermions and producing insulating transport when $L/\xi_{\text{dom}} \gg 1$.

The arguments outlined above are justified through a Landau-Ginzburg theory describing the quantum valley Hall order parameter. With the inclusion of disorder, we arrive at a classical $2d$ random-field Ising model, which allows us to estimate the scaling of the typical domain size as a function of system parameters. Through a simple extension of this formulation, we can further study what occurs when a different phase that does *not* couple directly to disorder is energetically favoured over the quantum valley Hall state in the clean limit. As expected, when the (clean) ground state energy splitting between the two states is sufficiently small—in a sense that we quantify with our Ising formulation—quantum valley Hall order prevails throughout the majority of the sample.

Our scenario for ubiquitous strong correlations at charge neutrality is not only compatible with existing charge-neutrality data, but further leads to falsifiable predictions both for STM and transport as described in Sec. VI. We also propose that two elements of this work may have broader applications in the study of mTBG. First, disorder can play a key role in discriminating among nearly degenerate correlated states. And second, disorder need not obliterate correlations, but can mask them as seen by global transport experiments.

The rest of the paper is organized as follows. We begin by reviewing the low-energy theory and establishing our conventions for twisted bilayer graphene in Sec. II. Next, Sec. III describes the fate of non-interacting mTBG Dirac fermions at charge neutrality in the presence of disorder. We then discuss the clean interacting theory in Sec. IV. The interaction form is first outlined [Sec. IV A], and then used to argue that the quantum valley Hall state is energetically competitive at charge neutrality [Secs. IV B and IV C]. Our main thesis is presented in detail in Sec. V, where each of the three regions illustrated in Fig. 2 is described in turn. We conclude in Sec. VI by summarizing and highlighting future directions. Supplemental details appear in numerous appendices.

II. REVIEW OF LOW-ENERGY THEORY

In this section, we set the stage by reviewing the low-energy physics of mTBG at charge neutrality in the absence of interactions and disorder.

A. Continuum model

Consider two monolayer-graphene sheets stacked such that they are twisted relative to one another by an angle θ , as shown in Fig. 3(a) (for an arbitrary angle θ). The twist dramatically reduces the system's translational symmetry. While true translational symmetry requires special commensurate angles, when the twist angle is small, an effective moiré translational symmetry emerges; the resulting triangular superlattice of orange AA regions, each surrounded by a hexagon of alternating AB and BA regions, is clearly visible in the cartoon of Fig. 3(a). In this case, the band structure at charge neutrality descends from the band structure of the individual graphene layers in a relatively straightforward manner when described in momentum space [3, 4]. Figure 3(b) shows the Brillouin zones (BZs) of the top and bottom graphene monolayers after applying a rotation by an angle $+\theta/2$ and $-\theta/2$, respectively. The reciprocal lattice vectors of the resulting moiré pattern are given by $\mathbf{G}_\ell = \mathcal{R}_{\theta/2}[\mathbf{G}_\ell] - \mathcal{R}_{-\theta/2}[\mathbf{G}_\ell]$ where $\mathbf{G}_{1,2}$ denote the reciprocal lattice vectors of the unrotated graphene sheets and $\mathcal{R}_\phi[\mathbf{v}]$ rotates a vector \mathbf{v} by an angle ϕ . The length of the moiré reciprocal lattice vectors, $|\mathbf{G}_\ell|$, is therefore suppressed relative to the graphene reciprocal lattice vectors by a factor of $2 \sin(\theta/2) \sim \theta$, making it very small by assumption. Equivalently, the moiré lattice constant is enlarged by $\sim 1/\theta$ relative to the graphene lattice constant. We let $\pm\mathbf{K}_t = \mathcal{R}_{\theta/2}[\pm\mathbf{K}]$ and $\pm\mathbf{K}_b = \mathcal{R}_{-\theta/2}[\pm\mathbf{K}]$ denote the $\pm\mathbf{K}$ points of the top and bottom layers, respectively.

As tunnelling between the two layers turns on, states at momentum $\mathbf{K}_t + \mathbf{k}$ on the top layer mix with those at momentum $\mathbf{K}_b + \mathbf{k} + \mathbf{q}_\ell + \mathbf{G}$ on the bottom layer, where \mathbf{G} is a moiré reciprocal lattice vector and $\mathbf{q}_\ell = \mathcal{R}_{2\pi(\ell-1)/3}[\mathbf{K}_t - \mathbf{K}_b]$, $\ell = 1, 2, 3$ [see Fig. 3(b)]. Since the moiré BZ is much

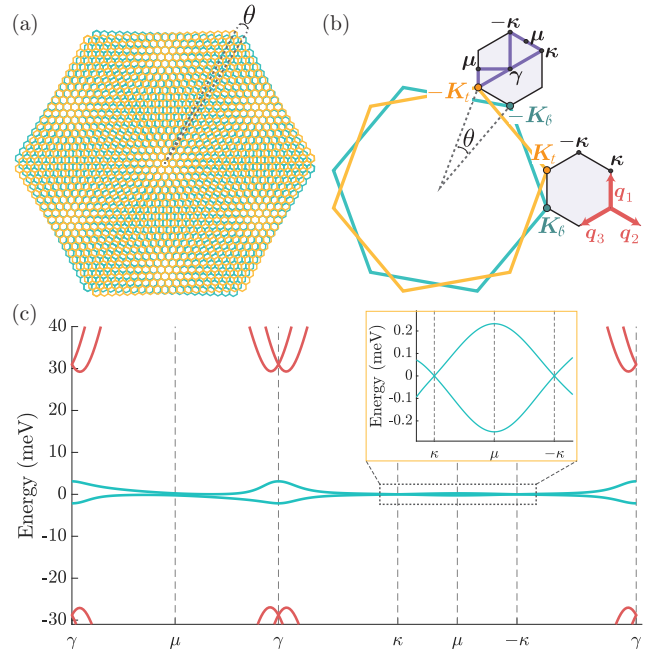


FIG. 3. (a) Cartoon representation of twisted bilayer graphene. The top and bottom graphene sheets are respectively represented by the orange and turquoise honeycomb lattices. The light, orange-tinted AA regions form a triangular superlattice, each of which is surrounded by a darker hexagonal rim whose vertices correspond to alternating AB and BA stacking regions. (b) Representation of the microscopic and moiré Brillouin zones. The large orange and turquoise hexagons represent the microscopic Brillouin zones of the underlying graphene layers. The moiré Brillouin zones, shown in purple, are defined by the distance between the the \mathbf{K} -points of the top and bottom layers. The size of the twist angles in both (a) and (b) has been exaggerated for clarity. (c) Flat bands corresponding to one \mathbf{K} -valley as calculated using the continuum model along the momentum line cut shown in purple in (b). Here, we took $\theta = 1.05^\circ$. Lattice relaxation is mimicked by decreasing the AA tunnelling amplitude w_0 relative to the AB tunnelling amplitude w_1 . In particular, we have $w_0 = 85$ meV and $w_1 = 110$ meV [6, 42, 43]. The inset zooms in on the Dirac cones at κ and $-\kappa$.

smaller than the BZ of monolayer graphene, mixing states proximate to \mathbf{K} with those proximate to $-\mathbf{K}$ is an extremely high-order tunnelling process; the two valleys of the original graphene monolayers thus effectively decouple. This decoupling is particularly convenient as it allows us to express the full band Hamiltonian H_{cont} as a sum of terms for the \mathbf{K} and $-\mathbf{K}$ valleys: $H_{\text{cont}} = H_+ + H_-$. For convenience, we explicitly reproduce H_+ in Appendix A. Our ability to decompose the Hamiltonian into \mathbf{K} -valley sectors is equivalent to the emergence of a U(1) “valley” symmetry, which we denote $U(1)_v$.

In addition to moiré translations and $U(1)_v$, the continuum model preserves the $SU(2)_s$ spin rotation symmetry (neglecting spin-orbit coupling), time reversal \mathcal{T} , C_2 rotations by π , C_3 rotations by $2\pi/3$, and a mirror symmetry \mathcal{M}_y that takes $(x, y) \rightarrow (x, -y)$ and interchanges the two layers. The latter

three should be regarded as emergent symmetries similar to $U(1)_v$. In our conventions the time-reversal operator \mathcal{T} does *not* flip the electronic spins and accordingly obeys $\mathcal{T}^2 = +1$. Both \mathcal{T} and \mathcal{C}_2 interchange the two valleys and hence are not symmetries of the individual single-valley Hamiltonians H_{\pm} . Rather than keep track of these two symmetries separately, it is therefore convenient to consider \mathcal{T} along with the composite operation $\mathcal{C}_2\mathcal{T}$ —which commutes with $U(1)_v$.

The momenta \mathbf{K}_t and \mathbf{K}_b map to the corners of the moiré BZ. In what follows, we denote these momenta by $\pm\boldsymbol{\kappa}$ to distinguish them from the $\pm\mathbf{K}$ -valleys of the microscopic graphene layers. Provided the $\mathcal{C}_2\mathcal{T}$ and \mathcal{C}_3 symmetries are present, the massless Dirac cones at $\mathbf{K}_{t,b}$ for the microscopic graphene layers evolve into massless Dirac cones at $\pm\boldsymbol{\kappa}$ even once tunnelling is turned on. This crucial property follows from the fact that the Berry phase enclosed within any loop is quantized to 0 or $\pi \pmod{2\pi}$ when $\mathcal{C}_2\mathcal{T}$ is preserved. Since a Dirac point necessarily exhibits Berry phase π , the Dirac cones at $+\boldsymbol{\kappa}$ and $-\boldsymbol{\kappa}$ are locally protected against a mass [5, 6]. Breaking \mathcal{C}_3 can shift the location of the Dirac cones, but cannot gap them. Importantly, since both cones in H_+ (H_-) descend from the Dirac cones at $\mathbf{K}_{t,b}$ ($-\mathbf{K}_{t,b}$) in a continuous fashion, they possess the same chirality [7, 8]—thereby obstructing the development of a two-band, single \mathbf{K} -valley tight-binding model in which all symmetries are realized in a local fashion [6, 44, 45].

B. Flat bands

The previous subsection highlighted generic features of small-angle twisted-bilayer graphene. At the *magic* angle, the velocity of the massless Dirac fermions becomes very small, and the bands immediately above and below the charge neutrality point separate from the remaining bands by a finite energy (provided lattice relaxation is incorporated [6, 42, 43]); see Fig. 3(c). The resulting energetically isolated “flat bands” are each (essentially) four-fold degenerate, reflecting spin and valley degrees of freedom. We now describe the flat-band Hamiltonian by first focusing on the $+\mathbf{K}$ valley and subsequently incorporating the $-\mathbf{K}$ valley.

Let $c_{\alpha j}(\mathbf{k})$ denote momentum-space annihilation operators associated with the flat bands at valley $+\mathbf{K}$; here $\alpha = \uparrow, \downarrow$ is a spin index and $j = 1, 2$ is a band index. Reference 46 showed that these operators can be defined such that they transform under $\mathcal{C}_2\mathcal{T}$ via

$$\mathcal{C}_2\mathcal{T} : \quad c(\mathbf{k}) \rightarrow \eta^x c(\mathbf{k}), \quad i \rightarrow -i \quad (1)$$

with Pauli matrices $\eta^{x,y,z}$ that act on the band indices. (Here and below we often suppress indices for notational simplicity.) It follows that the $\mathcal{C}_2\mathcal{T}$ -invariant flat-band Hamiltonian

takes the form

$$H_0 = \int_{\mathbf{k} \in BZ} c^\dagger(\mathbf{k}) \left[h_0(\mathbf{k}) + h_x(\mathbf{k})\eta^x + h_y(\mathbf{k})\eta^y \right] c(\mathbf{k}). \quad (2)$$

Next we project onto the massless Dirac fermions at $\pm\boldsymbol{\kappa}$ in the moiré BZ by defining Dirac spinors $\psi_{1\alpha j}(\mathbf{q}) \sim c_{\alpha j}(+\boldsymbol{\kappa} + \mathbf{q})$ and $\psi_{2\alpha j}(\mathbf{q}) \sim c_{\alpha j}(-\boldsymbol{\kappa} + \mathbf{q})$ and retaining only small \mathbf{q} modes. The fact that the massless Dirac cones exhibit the same chirality at $\pm\boldsymbol{\kappa}$ implies that $h_x(\pm\boldsymbol{\kappa} + \mathbf{q}) \sim +q_x$ and $h_y(\pm\boldsymbol{\kappa} + \mathbf{q}) \sim +q_y$. Upon shifting the energy such that $h_0(\pm\boldsymbol{\kappa}) = 0$ and reverting to real space, the low-energy Hamiltonian becomes

$$H_D = - \int_{\mathbf{r}} v_F \psi^\dagger (i\partial_x \eta^x + i\partial_y \eta^y) \psi. \quad (3)$$

The Fermi velocity v_F has been assumed isotropic and identical for both $\pm\boldsymbol{\kappa}$ Dirac cones, which is guaranteed when all symmetries outlined earlier are present.

An insulating phase at charge neutrality may only be obtained by either breaking the $\mathcal{C}_2\mathcal{T}$ symmetry or by closing the gap separating the flat bands and the dispersing bands [5, 6]. We focus entirely on the former scenario, which is straightforward to represent using the Dirac theory. Let $\tau^{x,y,z}$ and $\sigma^{x,y,z}$ denote Pauli matrices that respectively act on $\boldsymbol{\kappa}$ -valley indices and spin indices. Mass terms then take the form $\psi^\dagger \eta^z M \psi$ with $M = \{\mathbb{1}, \sigma^i, \tau^i, \tau^i \sigma^j\}$. The Chern number for a given spin/valley sector depends on the relative sign of the masses gapping the $\boldsymbol{\kappa}$ and $-\boldsymbol{\kappa}$ Dirac cones. When both cones have the same-sign mass, the sector acquires Chern number $C = \pm 1$, whereas opposite-sign masses yield $C = 0$. Consequently, mass terms with $M = \mathbb{1}, \sigma^i$ yield insulating bands with non-zero Chern number, while masses with $M = \tau^i, \tau^i \sigma^j$ yield trivial insulating bands².

We now restore the $-\mathbf{K}$ valley. In terms of the low-energy Dirac Hamiltonian, the chirality of the massless Dirac fermions in the $-\mathbf{K}$ valley is opposite that of the $+\mathbf{K}$ valley. Defining the spinor $\Psi = (\psi_+, \psi_-)^T$, where ψ_{\pm} describe Dirac fermions in valley $\pm\mathbf{K}$, the full Dirac Hamiltonian may be written

$$H_{D,\text{tot}} = - \int_{\mathbf{r}} v_F \Psi^\dagger (i\partial_x \mu^z \eta^x + i\partial_y \eta^y) \Psi, \quad (4)$$

where we introduced Pauli matrices $\mu^{x,y,z}$ that act on \mathbf{K} -valley indices. The presence of μ^z in the first term above implements the opposite-chirality requirement. Our discussion of the mass terms and the associated Chern numbers ex-

² Given the Hamiltonian H_D , the above conclusions regarding Chern number hold true regardless of the details of the high-energy theory from which it was derived. It is worth noting that for a single graphene sheet, expanding in small \mathbf{q} the functions analogous to $h_{x,y}(\pm\boldsymbol{\kappa} + \mathbf{q})$ would not yield the low-energy Hamiltonian of Eq. (3). Instead, the two Dirac cones would possess opposite chirality.

tends straightforwardly to the Ψ fermions. For details see Appendix B. A notable consequence of the opposing chiralities of the $\pm\mathbf{K}$ valleys is that a mass term $\Psi^\dagger \eta^z \Psi = \psi_+^\dagger \eta^z \psi_+ + \psi_-^\dagger \eta^z \psi_-$ generates an insulator with $C = +1$ for the $+\mathbf{K}$ valley and $C = -1$ for the $-\mathbf{K}$ valley (or vice versa depending on the overall sign of the mass term). These insulators correspond to the quantum valley Hall states that play a prominent role in this paper.

In the following sections, we use the operator ψ when restricting our discussion to a single \mathbf{K} -valley. We suppress the “ \pm ” indices in such cases but assume for concreteness that the $+\mathbf{K}$ valley is being considered (as in the beginning of this subsection). We reserve use of Ψ for occasions when both valleys are discussed simultaneously.

III. FREE FERMIONS WITH DISORDER

Next, we discuss the physics of non-interacting twisted bilayer graphene with disorder at charge neutrality.

A. Sources of disorder in twisted bilayer graphene

It is useful to review the specific types of disorder that are believed to be most relevant to experiments, though we attempt to keep the majority of our discussion as general as possible. Charge disorder appears to be quite low: Refs. 2, 9, and 19 estimate charge-carrier inhomogeneity in the range $\delta n \sim 1 - 2 \times 10^{10} \text{ cm}^{-2}$. Yankowitz *et al.* [9] further consider the observation of fractional quantum Hall states at magnetic fields as low as 4 T as additional proof of the high purity of their sample.

Twist-angle disorder is perhaps the most prevalent type of inhomogeneity in mTBG systems. Due to strain, different regions of a given sample may correspond to different twist angles, as directly imaged in STM [15–17]. From topography, the AA regions of the moiré structure are very clear, allowing one to locally establish the moiré lattice constants and thus the twist angle. Twist-angle variations were more recently characterized by Uri *et al.* [19] using a superconducting quantum interference device on a tip. Under an applied magnetic field, these authors measured the electron density of the sample as a function of the tip location, which in turn allowed them to map out the twist angle throughout the entire sample. Such measurements indicated local twist angles varying within a range $\delta\theta \sim 0.1^\circ$. Both samples they studied developed correlated insulating states, but only the sample with a continuous magic-angle region percolating across the sample displayed clear signs of superconductivity.

While transport measurements cannot access such local information, by comparing two-terminal conductance measurements between different pairs of contacts, Cao *et al.* [1, 2] and Yankowitz *et al.* [9] nevertheless note that some regions require different electron densities to achieve the band insulator at full-filling, again implying that unit cells differ

between regions. Similar measurements by Lu *et al.* [14] returned a much more uniform signal across the sample. Disorder signatures are also observable from within the superconducting states. Both Cao *et al.* and Yankowitz *et al.* observe phase-coherent Fraunhofer interference, indicating the coexistence of superconducting and normal regions. Conversely, the interference patterns measured by Lu *et al.* are comparatively weak, which they take as further indication of the high degree of sample homogeneity.

The hBN substrate may serve as yet another source of disorder. When uniformly aligned with one of the graphene monolayers, $C_2\mathcal{T}$ symmetry is explicitly broken and a gap at charge neutrality is opened [12, 13]. While the explicit gapping naturally explains the $\nu = 0$ insulator and anomalous Hall effect observed by Sharpe *et al.* [10] and Serlin *et al.* [11] at $\nu = +3$, hBN-alignment is believed to be an otherwise small effect in the majority of samples studied. Nevertheless, it is possible that a *local* alignment of the substrate, differing between regions, could weakly break the $C_2\mathcal{T}$ symmetry—just as for twist-angle disorder—even though it may be present on average.

B. Theoretical modelling of disorder

Motivated by the preceding discussion, we now incorporate weak, smooth disorder that preserves time-reversal and spin-rotation symmetries. We model such disorder by coupling spatially varying (but static) fields to fermion bilinears of the non-interacting Dirac theory reviewed in Sec. II B. The most relevant forms of disorder couple to bilinears that do not contain derivatives, and so we focus our study on this subset³. Time-reversal invariance and spin symmetry further reduce the number of bilinears capable of coupling to disorder; we enumerate all such symmetry-preserving terms in Appendix B. Collectively denoting the set of symmetry-allowed operators by $\{\Psi^\dagger T^i \Psi\}$, the most general disorder Hamiltonian takes the form

$$H_{\text{dis}} = \int_{\mathbf{r}} \sum_i R_i(\mathbf{r}) \Psi^\dagger(\mathbf{r}) T^i \Psi(\mathbf{r}). \quad (5)$$

We assume Gaussian-distributed $R_i(\mathbf{r})$ with zero mean and variance

$$\overline{R_i(\mathbf{r}) R_j(\mathbf{r}') } = \delta_{ij} g_i^2 K_i((\mathbf{r} - \mathbf{r}')/\xi_i). \quad (6)$$

Here, g_i is the disorder strength with units of energy, ξ_i is the disorder correlation length, and K_i is a dimensionless function that characterizes the spatial correlations of the disorder and obeys $K_i(0) = 1$. We frequently specialize to the case

³ In particular, we neglect disorder-induced variation in the Fermi velocity. This omission is supported by the numerics of Ref. 20, which show that the velocity remains largely unaffected by the presence of twist-angle disorder.

where the spatial correlations are Gaussian, *i.e.*,

$$K_i(\mathbf{r}/\xi_i) = e^{-\mathbf{r}^2/(2\xi_i^2)}. \quad (7)$$

Weakness of the disorder implies that g_i are small relative to the other scales of the theory, enabling a perturbative treatment. Smoothness of disorder is imposed by requiring that $\xi_i \gtrsim a_M$, with a_M the moiré lattice constant. We assume that the correlation lengths corresponding to different forms of disorder do not differ substantially and simply set $\xi_i = \xi_{\text{dis}}$ for all i .

The smoothness condition is physically very natural given that the existence of the moiré superlattice and the resulting band structure is predicated on the absence of fluctuations on the scale of the graphene lattice constant a . In momentum space, smoothness implies the suppression of inhomogeneities mediating momentum exchanges of order $\sim |\mathbf{K}|$, *i.e.*, disorder processes that couple to bilinears of the form $\Psi^\dagger \mu^{x,y} M \Psi$. In fact, we demonstrate in Appendix C that given Gaussian-correlated disorder [Eq. (7)], the disorder strengths corresponding to inter- \mathbf{K} -valley scattering are exponentially suppressed relative the intra- \mathbf{K} -scattering disorder strengths: if g is the magnitude of a typical intra- \mathbf{K} -valley disorder field, then

$$g_{\mathbf{K}\mathbf{K}'} \sim g e^{-\mathbf{K}^2 \xi_{\text{dis}}^2/4} = g e^{-4\pi^2 \xi_{\text{dis}}^2/a^2} \quad (8)$$

is the typical amplitude of an inter- \mathbf{K} -valley scattering event. Neglecting such exponentially suppressed events for now, we focus on a single \mathbf{K} -valley and couple disorder to the ψ fermions described by H_D .

Since time-reversal interchanges \mathbf{K} -valleys, it is *not* a symmetry of the single- \mathbf{K} -valley theory, implying that the system is described by the Wigner-Dyson class A [47, 48]. Disorder can thus couple to all spin-rotation-invariant bilinears and takes the form

$$\begin{aligned} H_{\text{dis}}^{\text{smooth}} = & \int_{\mathbf{r}} \psi^\dagger(\mathbf{r}) \left\{ \mathcal{M}_0(\mathbf{r}) \eta^z + \mathcal{M}_\ell(\mathbf{r}) \eta^z \tau^\ell \right. \\ & + \sum_{i=x,y} \left[\mathcal{A}_{i,0}(\mathbf{r}) \eta^i + \sum_{\ell=x,y,z} \mathcal{A}_{i,\ell}(\mathbf{r}) \eta^i \tau^\ell \right] \\ & \left. + \mathcal{V}_0(\mathbf{r}) + \sum_{\ell=x,y,z} \mathcal{V}_\ell(\mathbf{r}) \tau^\ell \right\} \psi(\mathbf{r}), \quad (9) \end{aligned}$$

where \mathcal{M} , \mathcal{A} , and \mathcal{V} respectively represent various forms of mass, vector potential, and scalar potential disorder.

It is also useful to consider the limit where disorder is sufficiently smooth relative to the moiré lattice scale that inter- κ -valley scattering may also be neglected. We can then further restrict our attention to one of the Dirac cones in the moiré BZ—say $+\kappa$. Denoting the spinor describing the Dirac cone at $+\kappa$ by $\chi(\mathbf{r})$, the disorder Hamiltonian be-

comes simply

$$\begin{aligned} H_{\text{dis}}^{\text{ultra-smooth}} = & \int_{\mathbf{r}} \chi^\dagger(\mathbf{r}) \left[m(\mathbf{r}) \eta^z \right. \\ & \left. + \sum_{i=x,y} a_i(\mathbf{r}) \eta^i + v(\mathbf{r}) \right] \chi(\mathbf{r}), \quad (10) \end{aligned}$$

where the random mass m , vector potential $a_{x,y}$ and scalar potential v satisfy Eq. (6) with $R_i = m, a_{x,y}, v$. Since each moiré unit cell encompasses $\sim 10\,000$ carbon atoms, distilling the disorder Hamiltonian down to Eq. (10) is significantly more suspect than merely omitting inter- \mathbf{K} -valley scattering terms. Moreover, though it might naively appear that inter- κ -scattering should be suppressed in a manner analogous to Eq. (8), we only expect such an effect to be manifest for extremely large correlation lengths ξ_{dis} relative to a_M as discussed at the end of Appendix C. We nevertheless argue in Sec. V that interactions greatly enhance the validity of Eq. (10) over a broader parameter regime.

C. Free Dirac fermions coupled to disorder

While we are interested in the situation where the disorder strength is subleading relative to interactions, it is instructive to review the expected fate of the free Dirac theory at charge neutrality in several limits. Consider first the single-Dirac-cone theory with disorder described by $H_{\text{dis}}^{\text{ultra-smooth}}$. Having restricted to this minimal theory, it is convenient to abandon smooth disorder and instead take white-noise correlations such that $K_i(\mathbf{r}/\xi_{\text{dis}}) = \xi_{\text{dis}}^2 \delta^2(\mathbf{r})$. Physically, this simplification implies that we are probing the system at long enough scales relative to ξ_{dis} that all correlations in R_i are washed away. The disorder correlation length is then encoded in the dimensionless (up to factors of \hbar and v_F) disorder strength parametrized by $g_i^2 \xi_{\text{dis}}^2$.

Ludwig *et al.* [49] analyzed the effect of each of the three remaining disorder fields— $m(\mathbf{r})$, $a_{x,y}(\mathbf{r})$, and $v(\mathbf{r})$. In the absence of all other types of disorder, the random mass, vector potential, and scalar potential fields were individually found to be marginally irrelevant, exactly marginal, and marginally relevant in turn. Ludwig *et al.* further postulated that when all three disorder types are simultaneously present, the system flows to the integer quantum Hall (IQH) plateau transition fixed point.

The correspondence between Landau-level physics and disordered Dirac theories may be understood from the perspective of a Chalker-Coddington network model [37]. This model can be employed to efficiently study the transition between a trivial insulator with Landau-level filling $\tilde{\nu} = 0$ an IQH state with $\tilde{\nu} = 1$ (the tilde distinguishes Landau-level filling from the mTBG filling). The system is assumed to locally prefer either $\tilde{\nu} = 0$ or $\tilde{\nu} = 1$ —thus forming domains of trivial and IQH states whose detailed structure depends on the total filling and the disorder potential. As the total filling varies, either the trivial state percolates, with small “lakes” of

$\tilde{\nu} = 1$, or vice versa. At some critical value, the system transitions between these two limits and becomes gapless. The network model exploits the fact that each boundary between $\tilde{\nu} = 0$ and $\tilde{\nu} = 1$ regions binds a chiral edge mode, and maps the problem onto one of directed links scattering at different nodes.

The key observation for our purposes is that the network model may be directly mapped onto a massless Dirac cone coupled to random fields $m(\mathbf{r})$, $a_{x,y}(\mathbf{r})$, and $v(\mathbf{r})$ [40, 41]. The correspondence between a lone disordered Dirac cone and the IQH plateau transition has been studied more recently in the context of *monolayer* graphene [50], where reducing the problem to that of a single Dirac cone only requires that disorder correlations are smooth on the scale of the microscopic lattice. When the effective time-reversal symmetry of the single Dirac cone is broken by strain [51], the appropriate nonlinear σ -model was shown to possess a topological term with $\theta = \pi$; consequently, the system exhibits universal conductivity, just as predicted for the IQH plateau transition by Pruisken [52, 53].

Upon resurrecting inter- κ -valley scattering in mTBG, disorder is instead described by $H_{\text{dis}}^{\text{smooth}}$. Here the theory localizes in the thermodynamic limit, and the conductivity accordingly approaches zero even at charge neutrality [54, 55]. Nevertheless, the localization length is expected to be extremely long since the scaling theory of Anderson localization indicates a lower critical dimension of $d = 2$ [56, 57]. The conductance thus only vanishes logarithmically with system size, suggesting that the localization length may be exponentially long, at least for a typical metal. Fradkin studied the fate of a system featuring N_f massless Dirac cones in the large- N_f limit [58, 59]. Denoting the disorder strength for all processes simply by g , he obtained an exponentially large mean free path,

$$\ell_{\text{mfp}} \sim a_M \exp \left[\frac{\pi}{2} \left(\frac{\hbar v_F}{g \xi_{\text{dis}}} \right)^2 \right], \quad (11)$$

and a still-larger localization length $\xi_{\text{loc}} \sim \ell_{\text{mfp}} \exp(64N_f^2/9)$.

IV. INTERACTIONS IN THE CLEAN LIMIT

Turning away from the question of disorder, we now investigate the effect of interactions in a homogeneous sample (though we occasionally allude to disorder effects). We begin with a discussion of the general form and magnitude of the interactions. Drawing on numerical results, experimental observations, and symmetry considerations, we then argue that the quantum valley Hall state is energetically competitive in interacting mTBG at charge neutrality.

A. Coulomb interaction

The Coulomb Hamiltonian $H_{C,\text{tot}} = \frac{1}{2} \int_{\mathbf{q}} V(\mathbf{q}) \rho(\mathbf{q}) \rho^\dagger(\mathbf{q})$ encodes the leading interaction. Here $V(\mathbf{q})$ is the Fourier-transform of the long-range Coulomb potential (which technically depends on both the layer and sublattice, but these microscopic corrections can be ignored for the purpose of our discussion). The operator $\rho(\mathbf{q})$ represents the Fourier transform of the full microscopic density. Specifically, we write $\rho(\mathbf{q}) = \sum_{\ell} \int_{\mathbf{k}} \tilde{f}_{\ell}^\dagger(\mathbf{k}) \tilde{f}_{\ell}(\mathbf{k} + \mathbf{q})$, where $\tilde{f}_{\ell}(\mathbf{k})$ denotes the annihilation operator corresponding to one of the decoupled graphene monolayers, with ℓ a combined index labelling both layer and sublattice and \mathbf{k} taking values across the full microscopic BZ. As explained in Sec. II A, to a high degree of accuracy the flat-band wavefunctions are composed entirely of states originating proximate to the Dirac cones of the decoupled monolayers. We focus on these important momenta by introducing operators $f_{\ell,n=\pm}(\mathbf{k}) \equiv \tilde{f}_{\ell}(\mathbf{k} \pm \mathbf{K})$ that are defined for $|\mathbf{k}| \ll |\mathbf{K}|$; note that this ‘‘small \mathbf{k} ’’ condition does not necessarily imply that \mathbf{k} resides within the moiré BZ. It follows that only the density operators $\rho(\mathbf{q})$ and $\rho(\mathbf{q} \pm \mathbf{K})$ with \mathbf{q} small are physically relevant to the flat-band physics:

$$\begin{aligned} \rho(\mathbf{q}) &\cong \sum_{\ell,n} \int_{\mathbf{k} \text{ small}} f_{\ell,n}^\dagger(\mathbf{k}) f_{\ell,n}(\mathbf{k} + \mathbf{q}), \\ \rho(\mathbf{q} + \mathbf{K}) &\cong \sum_{\ell} \int_{\mathbf{k} \text{ small}} f_{\ell,+}^\dagger(\mathbf{k}) f_{\ell,-}(\mathbf{k} + \mathbf{q}) \\ &= \rho^\dagger(-\mathbf{q} - \mathbf{K}). \end{aligned} \quad (12)$$

Inserting these definitions into our expression for $H_{C,\text{tot}}$ we find $H_{C,\text{tot}} \cong H_C + H'_C$ where

$$\begin{aligned} H_C &= \frac{1}{2} \int_{\mathbf{q} \text{ small}} V(\mathbf{q}) \rho(\mathbf{q}) \rho^\dagger(\mathbf{q}), \\ H'_C &= \int_{\mathbf{q} \text{ small}} V(\mathbf{q} + \mathbf{K}) \rho(\mathbf{q} + \mathbf{K}) \rho^\dagger(\mathbf{q} + \mathbf{K}). \end{aligned} \quad (13)$$

There is a vast separation of energy scales between H_C and H'_C . Since $V(\mathbf{q}) \propto 1/|\mathbf{q}|$, the largest contribution to H_C comes from momenta \mathbf{q} within the moiré BZ, *i.e.* $|\mathbf{q}| \lesssim |\kappa|$. On the other hand, in H'_C , the smallness of the internal momentum \mathbf{q} implies $V(\mathbf{q} + \mathbf{K}) \approx V(\mathbf{K})$. It follows that the relative strength of H_C and H'_C is $V(\mathbf{K} + \mathbf{q})/V(\mathbf{q}) \lesssim V(\mathbf{K})/V(\kappa) \sim |\kappa|/|\mathbf{K}| \sim \theta \ll 1$. Hereafter we focus our attention on the dominant term, H_C . Reverting to real space, the Coulomb potential is $V(\mathbf{r}) = e^2/(4\pi\epsilon|\mathbf{r}|)$. For graphene on hBN, we estimate the dielectric constant to be $\epsilon \sim 8\epsilon_0$ with ϵ_0 denoting the permittivity of free space. Using the moiré lattice spacing, $a_M = a/(2\sin(\theta/2))$, where $a = 0.246$ nm is the lattice constant of monolayer graphene as a typical length scale, one finds a characteristic interaction energy $V(a_M) \sim 14$ meV at the magic angle $\theta \sim 1.05^\circ$.

Theory estimates the bandwidth of the flat bands to be

about 10 meV and the splitting between van Hove peaks within those bands to be ~ 5 meV [43, 60]. The above Coulomb-interaction scale thus raises natural questions regarding the validity of our expansion about the Dirac cones at $\pm\kappa$ in Sec. II B. It appears that the entirety of the flat bands and perhaps even neighboring energy bands should be considered. Non-interacting simulations of mTBG systems with twist angle disorder, however, have been shown to increase the bandwidth with little change to the Dirac character at charge neutrality [20]. Moreover, STM measurements of the fully filled flat bands (*i.e.*, in a regime where correlations are presumably less important) measure van Hove peak splittings of $\sim 10 - 20$ meV [15, 16]—several times larger than the above theoretical estimate. The full bandwidth of the flat bands may therefore significantly exceed $V(a_M)$, supporting our use of the Dirac theory.

B. Preferred ground state of single-flavour theory

Before turning to the full theory, it is useful to examine interaction effects at charge neutrality in a minimal, single-flavour model that includes only one spin and one \mathbf{K} -valley. References 14, 16, 21, and 22 addressed this problem numerically via self-consistent Hartree-Fock calculations. Liu *et al.* [21] incorporated Coulomb interactions in the continuum model while Choi *et al.* [16] studied a 10-band lattice model [44] with a simplified local interaction. Both analyses find a $C_2\mathcal{T}$ -breaking gapped state with Chern number $C = \pm 1$ as the lowest-energy solution.

References 14 and 22 also predict an interaction-induced gapped phase at charge neutrality. However, while certain parameter regimes again return a $C_2\mathcal{T}$ -breaking state with nonzero Chern number, other regimes yield a $C_2\mathcal{T}$ -preserving, trivial insulator. The latter statement may seem at odds with our assertion in Sec. II A that $C_2\mathcal{T}$ protects the masslessness of the Dirac cones, but this protection only holds when the flat bands are energetically isolated. In the calculations of Refs. 14 and 22, interactions close the gap separating the flat and dispersive bands, thus negating the protection conferred upon the Dirac cones by $C_2\mathcal{T}$ symmetry. It is worth noting that STM measurements show that the flat bands indeed remain isolated as a function of filling [15–18], and yet still resolve correlation effects. We therefore view the formation of the $C_2\mathcal{T}$ -symmetric insulator as a less likely outcome.

Returning to the $C_2\mathcal{T}$ -breaking gapped states, we remark that from the perspective of the Dirac theory, it is natural to expect the phase with $C = \pm 1$ to be energetically favourable relative to the $C_2\mathcal{T}$ -breaking trivial insulator with $C = 0$. Recall from Sec. II B that the mean-field order parameter for the $C = \pm 1$ state is $\psi^\dagger \eta^z \psi$, which in principle can arise from a momentum-independent microscopic perturbation. The trivial $C = 0$ phase instead corresponds to an order parameter $\psi^\dagger \eta^z \tau^z \psi$ that yields opposite-sign masses for the Dirac cones at $\pm\kappa$ —and hence cannot arise from a

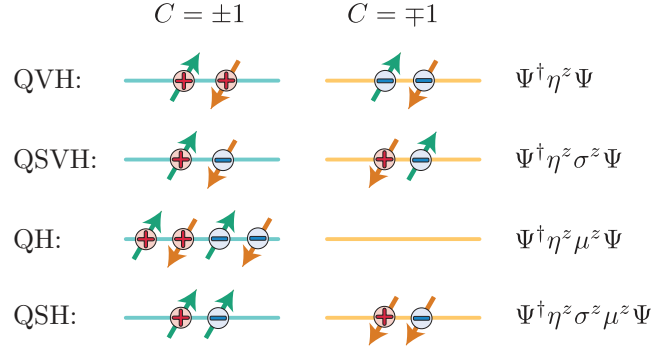


FIG. 4. Four natural $C_2\mathcal{T}$ -breaking insulators at charge neutrality. In order from top to bottom: quantum valley Hall (QVH), quantum spin-valley Hall (QSVH), quantum Hall (QH), and quantum spin Hall (QSH). The direction of the arrow indicates spin, while the sign, ‘+’ or ‘-,’ labelling the arrow indicates the \mathbf{K} -valley.

momentum-independent microscopic perturbation. In monolayer graphene the converse situation arises: the momentum-independent staggered sublattice potential generates a trivial insulator whereas the relatively baroque, momentum-dependent Haldane mass [61] is instead required to enter a $C = \pm 1$ phase. (This distinction reflects the fact that the Dirac cones at $\pm\kappa$ exhibit the same chirality in mTBG, while the Dirac cones at $\pm\mathbf{K}$ in monolayer graphene have opposite chirality [46].) Spontaneously generating a Haldane mass in monolayer graphene is thus unnatural—see, *e.g.*, Ref. 62—and it is analogously difficult to spontaneously enter the $C = 0$ phase in mTBG.

C. Inclusion of spin and \mathbf{K} -valley flavours

We have so far argued that the single-flavour version of interacting, charge-neutral mTBG prefers to enter a gapped phase with Chern number $C = \pm 1$. Inclusion of spin and \mathbf{K} -valley degrees of freedom not only allows for many distinct possible phases depending on the Chern numbers assigned to each sector, but further allows for additional phases that do not naturally descend from the single-flavour theory. Let us begin by discussing the former.

We specifically focus on four natural candidate insulators that we refer to as quantum valley Hall (QVH), quantum spin-valley Hall (QSVH), quantum Hall (QH), and quantum spin Hall (QSH) phases. Figure 4 depicts these states along with their corresponding mass terms. These insulators carry different symmetry properties as summarized in Table I. While all four phases break C_2 while preserving time reversal \mathcal{T} , whereas the converse is true of the QH and QSH states. They are further distinguished by the action of the $SU(2)_s$ spin symmetry, which is preserved (broken) by the QVH and QH (QSVH and QSH) states. Note that because of the additional \mathbf{K} -valley flavour

	\mathcal{T}	\mathcal{C}_2	$SU(2)_s$
QVH	✓	✗	✓
QSVH	✓	✗	✗
QH	✗	✓	✓
QSH	✗	✓	✗

TABLE I. Symmetry breaking pattern of the four topological states. Note that the QSH phase violates \mathcal{T} , but does preserve the ‘physical’ electronic time reversal operation $\mathcal{T}_{\text{elec}} = i\sigma^y\mathcal{T}$.

index, our QSH state differs from the $2d$ topological insulator realized, *e.g.*, in the Kane-Mele model [63]. We nevertheless adopt this nomenclature since the state breaks spin-rotation symmetry and preserves the ‘physical’ electronic time reversal operation $\mathcal{T}_{\text{elec}} \equiv i\sigma^y\mathcal{T}$ that obeys $\mathcal{T}_{\text{elec}}^2 = -1$.

It is useful to highlight some physical differences between these states and thus their compatibility with experimental observations. Neither the QVH nor QSVH insulator is expected to possess gapless edge modes at a sample boundary. Our discussion has made significant use of the approximate $U(1)_v$ valley symmetry, but this symmetry is violently broken by the edge itself, which naturally occurs on the microscopic length scale a of the underlying graphene monolayers. As a result, the edge modes from the two valleys scatter strongly, resulting in a purely insulating state. By contrast, the QH state hosts robust gapless edge modes that are completely immune from scattering by virtue of their chirality. Edge modes of the QSH insulator, while nonchiral, are nevertheless also robust since backscattering at a sample boundary must be accompanied by a spin flip. The sample studied by Lu *et al.* [14] displayed insulating transport with no signs of edge conduction. Among the four insulators, QH and QSH states thus appear unlikely, at least in that platform.

As a result of the separation of scales between \mathbf{K} -valleys and the $SU(2)_s$ symmetry, all four insulating states have very similar energies. In Appendix D, we compare the QVH ground-state energy against the other three insulators using a simple Hartree-Fock variational approach. We show that all four states are exactly degenerate in the chiral model [64], a version of the continuum model that possesses an exact particle-hole symmetry that renders it exactly solvable. Nevertheless, for more realistic versions of the continuum model (where particle-hole symmetry is absent), we find that the QVH state is actually disfavoured relative the other insulators. However, when computed numerically, we find the energy difference to be extremely small, less than $\sim 10^{-5}$ meV per electron, implying that the explicit breaking of particle-hole symmetry has little effect.

We turn now to alternative phases. Polarized phases—for which the flat bands of two flavours are fully occupied—represent one class of competing ground states. In general, both spin- and valley-polarized phases are degenerate at charge neutrality when H'_C [Eq. (13)] is neglected [65, 66]. Liu *et al.* [21] find that, within the chiral model, these

polarized states have identical Fock energies to the $\mathcal{C}_2\mathcal{T}$ -breaking insulators with nontrivial Chern number in each flavour. They also obtained self-consistent versions of these solutions numerically using a more realistic version of the continuum model; while no longer exactly degenerate, these states remained close in energy. Adding explicit \mathcal{C}_3 -breaking strain—as observed in multiple STM experiments [15–17]—was, however, found to promote the $\mathcal{C}_2\mathcal{T}$ -breaking insulators over the polarized states. Another proposed state is the intervalley coherent phase (IVC) [6], which spontaneously breaks $U(1)_v$ symmetry by coupling the $+\mathbf{K}$ and $-\mathbf{K}$ bands. General considerations [65] as well as calculations using the analytically tractable chiral model [21] indicate that IVC order is disfavoured at the Hartree-Fock level. Other numerics nevertheless challenge these conclusions [67].

Importantly, among the gapped phases discussed here, only the QVH order parameter directly couples to disorder that is smooth and preserves \mathcal{T} and $SU(2)_s$ spin symmetry. Time reversal and spin symmetry forbid coupling to the order parameters for QSVH, QH, QSH, and polarized phases, whereas smoothness of disorder prohibits coupling to an IVC order parameter. Hence, even if one of the latter states is energetically favourable in a perfectly clean system, the unavoidable presence of inhomogeneity in any physical sample may nevertheless stabilize the QVH phase, a possibility that we explore in Sec. VC.

V. INTERPLAY OF INTERACTIONS AND DISORDER

We are now in position to explore the fate of charge-neutral mTBG in the presence of interactions *and* smooth disorder. Let us first recapitulate the expected behavior in the disordered, non-interacting limit (Sec. III) and in the clean but strongly interacting regime (Sec. IV):

1. In the absence of interactions, disorder localizes the massless Dirac fermions when any form of inter- κ -valley scattering is present in a manner that is formally analogous to physics of monolayer graphene. However, while monolayer graphene only requires that disorder be long-ranged on the scale of the microscopic lattice to avoid localization, disorder must be long-ranged on the scale of the moiré lattice to suppress localization in twisted bilayer graphene.
2. We have argued that in the strongly interacting, clean limit, the QVH phase that spontaneously breaks \mathcal{C}_2 symmetry constitutes (at the very least) an energetically competitive state that is compatible with experimental observations. Moreover, we observed that among various other candidate ground states, QVH order uniquely couples to smooth disorder respecting spin and time-reversal symmetries.

To simultaneously incorporate interactions and disorder below, we start with the assumption that the QVH state is the

true ground state of the clean, interacting Hamiltonian. We construct an Ising formulation of the system in the presence of disorder, which allows us to systematically consider the crossover between the first and second panels of Fig. 2. We discuss the titular recovery of the massless Dirac cones before showing that even when the QVH insulator is not the true ground state in the clean theory, disorder may nevertheless tip the balance back in its favour. We close with some comments on the eventual localization of the Dirac fermions, as illustrated in the final panel of Fig. 2.

A. Ising model formulation and domain formation

Suppose that the interaction energy scale dominates the physics, preferring to spontaneously break the \mathcal{C}_2 symmetry and form a QVH insulator. Disorder terms that do not couple to the QVH order parameter can then be neglected, leaving only the random field $\mathcal{M}_0(\mathbf{r})$ that couples to $\psi^\dagger \eta^z \psi$ in Eq. (9) (or, in the full theory, a random scalar field that couples to $\Psi^\dagger \eta^z \Psi$). A random mass cannot produce localization but does compete against long-range order. In fact, we show that even when disorder is weak and uncorrelated, the system always loses long-range order in the thermodynamic limit due to the formation of domains, as sketched in the central panel of Fig. 2. Destruction of long-range order only becomes observable, however, once the linear extent of the system, L , exceeds the typical domain size, ξ_{dom} . The goal of this subsection is to demonstrate that ξ_{dom} is finite and to determine its size as a function of the physical parameters of the theory.

We approach the problem in the standard fashion, via the formulation of a Landau-Ginzburg theory. The order parameter for the \mathcal{C}_2 symmetry breaking is simply an Ising field ϕ obtained by coarse graining the bilinear $\psi^\dagger \eta^z \psi$, *i.e.*,

$$\phi(\mathbf{r}) \sim \int_{\mathbf{r}' \in R(\mathbf{r})} \psi^\dagger \eta^z \psi(\mathbf{r}') \sim \ell_{\text{UV}}^2 \psi^\dagger \eta^z \psi(\mathbf{r}), \quad (14)$$

where $R(\mathbf{r})$ is a spatial region centred at \mathbf{r} of typical size ℓ_{UV}^2 and ℓ_{UV} is an ultraviolet cutoff quantified below. Since we are interested in the physics deep within the ordered phase with $\langle \phi \rangle \neq 0$, a *classical* Ising model suffices:

$$H_{\text{Ising}} = \int_{\mathbf{r}} \left[\mathcal{K}(\nabla \phi)^2 + \frac{r}{2} \phi^2 + \frac{u}{4!} \phi^4 \right]. \quad (15)$$

The mass r is clearly assumed to be negative.

The scales of the original fermionic Hamiltonian ultimately determine parameters of the Ising model, though this assignment is not necessarily straightforward. Consider first \mathcal{K} . Since ϕ is dimensionless, \mathcal{K} has units of energy, and hence $\mathcal{K} \sim U$, with U a characteristic energy scale of the system. Both the rough estimate for the Coulomb potential, $V(a_M) \approx 14 \text{ meV}$, given at the end of Sec. IV A, and the experimentally measured transport gap at charge neutrality, $\Delta_{\text{CNP}} \approx 1 \text{ meV}$ [14], provide natural candidates for U .

Given uncertainties in our calculation of $V(a_M)$ related to screening from other bands, we view the latter option as a more reasonable and conservative estimate. We stress however that this choice has little direct bearing on the discussion that follows.

It is also important to assign a length scale to the interactions and hence the Ising theory. Since our primary goal is to describe domain-wall physics, the most natural scale is

$$\xi_{\text{int}} \sim \frac{\hbar v_F}{\Delta_{\text{CNP}}}, \quad (16)$$

which corresponds to the decay length of a Dirac fermion of mass $\Delta_{\text{CNP}}/v_F^2$. In our context, these fermions are the chiral modes that bind to the domain walls at which the Chern numbers for each flavour change sign, identifying ξ_{int} as the domain boundary width. Any physics occurring on scales smaller than ξ_{int} necessarily includes these fermionic degrees of freedom, and hence lies outside our Ising formulation's regime of validity. The interactions length scale therefore defines a UV cutoff.⁴ As a consistency check, we must verify that ξ_{int} exceeds the moiré lattice constant, $a_M \approx 12.8 \text{ nm}$. Inserting $v_F \approx 0.15 \times 10^6 \text{ m/s}$ [2] and $\Delta_{\text{CNP}} \approx 1 \text{ meV}$ [14] into Eq. (16), we indeed find $\xi_{\text{int}} \approx 100 \text{ nm} \sim 10 a_M$. We are therefore permitted to set $\ell_{\text{UV}} \sim \xi_{\text{int}}$. In turn, dimensional analysis gives $r, u \sim U/\xi_{\text{int}}^2$.

Because disorder breaks \mathcal{C}_2 , it should couple to the Ising field in a manner that breaks the \mathbb{Z}_2 Ising symmetry. In other words, disorder appears as a random ‘‘magnetic’’ field:

$$H_{\phi, \text{dis}} = \int_{\mathbf{r}} \mathcal{B}(\mathbf{r}) \phi(\mathbf{r}), \quad (17)$$

where $\mathcal{B}(\mathbf{r}) \sim \int_{\mathbf{r}' \in R(\mathbf{r})} \mathcal{M}_0(\mathbf{r}')/\xi_{\text{int}}^4$. The random field \mathcal{M}_0 is defined by the disorder strength δm , correlation length ξ_{dis} , and correlation function $K(\mathbf{r}/\xi_{\text{dis}})$ (in the notation of Sec. III B, these quantities correspond to $g_{\mathcal{M}_0}$, $\xi_{\mathcal{M}_0}$, and $K_{\mathcal{M}_0}$, respectively). We focus on the situation where the disorder is Gaussian correlated: $K(\mathbf{r}/\xi_{\text{dis}}) = e^{-\mathbf{r}^2/(2\xi_{\text{dis}}^2)}$. Our assertion that the interaction energy scale dominates the disorder energy scale can now be more precisely stated as $\delta m/U \ll 1$.

In summary, the Hamiltonian controlling the ordering of ϕ is $H_{\text{RFIM}} = H_{\text{Ising}} + H_{\phi, \text{dis}}$, which is none other than the much-studied random field Ising model (RFIM)⁵ [69–71]. As claimed, the RFIM in $2d$ is generically disordered [72, 73], and so ξ_{dom} is finite. The mechanism of domain formation depends largely on the magnitude of the ratio

$$\alpha = \frac{\delta m}{U} \frac{\xi_{\text{dis}}}{\xi_{\text{int}}}. \quad (18)$$

⁴ The definition of ξ_{int} and ℓ_{UV} is largely independent of our choice of U .

⁵ This theory and its derivation should not be confused with the fact that a *free* Dirac fermion with random mass disorder maps onto the random *bond* Ising model [49, 68].

This result and the scenarios we outline below are derived and further explained in Appendix E.

We first examine what occurs when $\alpha \gtrsim 1$. Since $U/\delta m$ is already presumed large, in order for α to be larger than unity, this limit corresponds to that of extremely smooth disorder: $\xi_{\text{dis}}/\xi_{\text{int}} \gg 1$. In this scenario, the energy gained by having ϕ align in the direction preferred by $\mathcal{B}(\mathbf{r})$ is larger than the interaction energy cost associated with the misalignment of ϕ along the domain boundary. The Ising field therefore directly tracks the disorder potential, implying that

$$\xi_{\text{dom}} \sim \xi_{\text{dis}}, \quad \alpha \gtrsim 1. \quad (19)$$

The situation is more subtle when $\alpha \lesssim 1$. With stronger interactions, we naturally expect larger domains. At some point, the domains are large enough that the correlated nature of the disorder is washed away, allowing us to treat it as white noise: $\overline{\mathcal{M}_0(\mathbf{r})\mathcal{M}_0(\mathbf{r}')} \cong \delta m^2 \xi_{\text{dis}}^2 \delta^2(\mathbf{r} - \mathbf{r}')$. In this case, the destruction of long-range order occurs through the condensation of domain walls. An evaluation of the domain-wall roughening yields a lower bound for their size of [72]

$$\xi_{\text{dom}} \lesssim \max(\xi_{\text{int}}, \xi_{\text{dis}}) e^{c/\alpha^2}, \quad \alpha \lesssim 1, \quad (20)$$

where $c \sim \mathcal{O}(1)$ is a non-universal constant. We can verify that when $\alpha \ll 1$ the domain length scale is indeed far greater than the disorder correlation length, *i.e.*, $\xi_{\text{dom}} \gg \xi_{\text{dis}}$.

B. Recovery of massless Dirac fermions

Next we discuss the physical consequences of the Ising model outlined above in the regime where the system size L exceeds the typical domain size ξ_{dom} . For now we continue to assume suppression of both inter- κ - and inter- K -valley scattering. At least close to the crossover scale ξ_{dom} , the Ising formulation should remain valid: the system is characterized by multiple domains of opposing Chern numbers with typical size ξ_{dom} , as the central panel of Fig. 2 illustrates. In this regime, the system can be described by eight independent Chalker-Coddington network models [37]—one for each of the two Dirac cones within the four spin/valley sectors. As mentioned briefly in the introduction and more fully in Sec. III C, each network model may be mapped directly onto that of a single gapless Dirac cone [40, 41], thus giving the promised restoration of massless Dirac fermions from a strongly correlated starting point.

We can alternatively motivate the recovery of massless Dirac cones without relying on network models. Let us return to the full disordered Dirac theory described by Eqs. (4) and (5), which includes all spin and valley degrees of freedom. Notably, here we additionally allow for weak intervalley scattering terms. Upon including strong correlations at a mean-field level, interactions dramatically enhance the effective strength of the random field that couples to the QVH order parameter $\Psi^\dagger \eta^z \Psi$. All other disorder fields, by

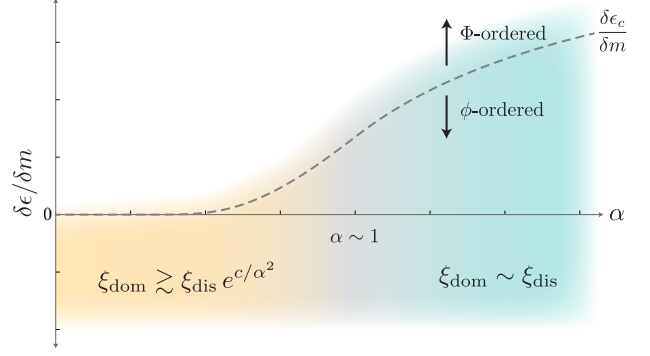


FIG. 5. Schematic phase diagram as a function of disorder, α , and ground state energy difference. The dashed line indicates the ‘critical’ energy difference $\delta\epsilon_c(\alpha)$ that characterizes the crossover from samples that are primarily ϕ -ordered to those that are primarily Φ -ordered. There are two distinct ϕ -ordered regimes. In the blue region, $\alpha \gtrsim 1$, $\langle\phi\rangle$ tracks the disorder so that domains are of the same size as the disorder correlation length, $\xi_{\text{dom}} \sim \xi_{\text{dis}}$. Conversely, in the orange region, $\alpha \lesssim 1$, the correlated nature of the disorder is unimportant, and domains are exponentially large, $\xi_{\text{dom}} \gtrsim \xi_{\text{dis}} e^{c/\alpha^2}$ [here, we assume that $\xi_{\text{dis}} > \xi_{\text{int}}$; see Eq. (20)]. In the white region above the dashed line, the competing phase prevails, and $\langle\Phi\rangle \neq 0$ throughout most of the sample.

contrast, remain weak and can be neglected to a first approximation. The problem then reduces to a set of independent Dirac cones, each governed by the far simpler disorder Hamiltonian in Eq. (10) with *only* random mass disorder. As noted earlier, the random mass is a marginally irrelevant perturbation to the clean Dirac theory when it is the sole source of disorder [49]. Massless Dirac fermions thus naturally re-emerge from this viewpoint as well.

At sufficiently low energy scales, however, the additional disorder fields neglected above eventually kick in. The dominant corrections are expected to arise from *intra- κ -valley* scattering processes, encoded by the vector- and scalar-potential terms in Eq. (10). When these terms are also present the theory is believed to flow to the IQH plateau transition, which is characterized by a finite density of states with both universal longitudinal and Hall conductances (here, valley Hall). At still lower energy scales, inter- κ -valley scattering is expected to produce localization, as Sec. V D discusses in more detail.

Nevertheless, the perspective just outlined should be viewed as a consistency check and not a proof of concept. Crucially, it cannot account for the energy scales separating the Dirac fermions of the clean, non-interacting mTBG system from the recovered Dirac cones of the interacting, disordered network model.

C. Competing phases

So far in this section, the QVH insulator has been taken as the true ground state of mTBG at charge neutrality, even in the absence of disorder. We now address the possibility that interactions prefer a different state. To simplify the problem, we consider the situation in which a single competing phase is energetically favourable relative to the QVH insulator. In accordance with the conventions of Sec. V A, this competition can be quantified through the energy difference $\delta\epsilon$ in an area of size $\ell_{UV}^2 = \xi_{\text{int}}^2$:

$$\frac{\delta\epsilon}{\xi_{\text{int}}^2} \equiv \epsilon_{\text{QVH}} - \epsilon_{\text{C}} \geq 0, \quad (21)$$

where ϵ_{QVH} and ϵ_{C} respectively denote the ground-state energy densities of the QVH state and competing phase. We further assume that the competing order may be described by an Ising field, Φ , that does not linearly couple to any disorder field; recall the discussion at the end of Sec. IV C. Generalizing our arguments to include continuous order parameters (as needed for the QSH and QSVH insulators) is straightforward, and we therefore leave the competing phase's identity unspecified. Note, however, that properties of domain walls separating the QVH order and the competing phase may depend on the precise nature of the latter.

We again work in a regime where strong interactions obviate the need to include all disorder fields save for the random mass \mathcal{M}_0 that linearly couples to the QVH order parameter via Eq. (17). Importantly, this type of disorder locally promotes the QVH state by lowering its energy relative to the competing phase, even though—as we saw earlier—it generally destroys true long-range order. When $\delta\epsilon$ is small enough, we expect the majority of the sample to realize the QVH phase and the scenario outlined in the previous section to hold. In terms of the Ising theory devised in the previous section, we can express this condition as

$$\left[\frac{1}{\text{vol}} \int_{\mathbf{r}} \langle \phi^2(\mathbf{r}) \rangle \right]^{1/2} \gtrsim \frac{1}{2}, \quad (22)$$

where ‘vol’ denotes the sample volume. When this equation holds, we say the system is ‘ ϕ -ordered’; otherwise, the system is ‘ Φ -ordered.’

Appendix F explores this problem in depth, ultimately deriving the schematic phase diagram shown in Fig. 5. We again find that the primary control parameter is the ratio α [recall Eq. (18)] corresponding to the horizontal axis. Motivated by the notion that Φ -ordered regions may be viewed as annealed ‘vacancies,’ we begin with a dilute Ising-model description. At the lattice level, the theory is conveniently formulated by promoting the Ising variables $\sigma^z = \pm 1$ to three-state spin-1 variables s , where $s = \pm 1$ correspond to the two QVH phases and $s = 0$ corresponds to the competing phase. We present a simple mean-field solution to the classical Blume-Capel model for these spin-1 degrees of freedom

[74–77] in Appendix F 1. While the phase diagram we obtain resembles the one shown in Fig. 5 in many respects, it erroneously predicts long-range ϕ -order when $\delta\epsilon < 0$ and disorder is sufficiently small, $\alpha \lesssim 1$; as discussed in Sec. V A and Appendix E, in reality, long-range order is unstable to the addition of any finite disorder. This failure of mean-field theory is not unprecedented given the low dimensionality.

In Appendix F 2, we therefore return to the Imry-Ma type arguments of Sec. V A (see also Appendix E), which allow us to derive a ‘critical’ energy difference $\delta\epsilon_c(\alpha)$ that characterizes the crossover scale separating ϕ - and Φ -ordered regimes. We plot $\delta\epsilon_c(\alpha)$ with a dashed line in Fig. 5. In the white region above the line, $\delta\epsilon \gtrsim \delta\epsilon_c(\alpha)$, the competing phase is realized throughout the majority of the system, and the network picture we propose is no longer relevant. Conversely, Eq. (22) holds in regions where $\delta\epsilon \lesssim \delta\epsilon_c(\alpha)$ (including the trivial case, $\delta\epsilon < 0$, where QVH states minimize the energy in the clean limit). Just as we found above, depending on the strength of disorder, the destruction of long-range QVH order occurs in two fashions. In Fig. 5, the parameter regime where $\langle \phi \rangle$ tracks the disorder field is shown in turquoise. The orange area indicates the opposite limit, where long-range order is eliminated by domain-wall condensation. The intermediate regime where $\alpha \sim 1$ is shown in neutral grey.

Notably, these considerations imply that disorder may not only be responsible for selecting which QVH order is locally realized, but that it may also determine whether or not QVH order is realized at all. In particular, our proposal admits a scenario in which the clean samples of Lu *et al.* [14] are Φ -ordered, while the less homogeneous samples of Cao *et al.* [1, 2] and Yankowitz *et al.* [9] realize the QVH network picture displayed in the central panel of Fig. 2—even supposing that the two sets of systems differ *solely* in the amount of disorder they present.

D. Localization

In the absence of any special symmetries, all two-dimensional systems are generically expected to localize in the thermodynamic limit, and our platform is no exception. Localization is likely irrelevant for the previously studied mTBG samples, whose linear dimensions are $\sim 2 - 8 \mu\text{m} \sim 150 - 600 a_M$. It is nevertheless instructive to briefly discuss localization within our proposed scenario. The precise manner in which localization occurs in the presence of interactions poses a notoriously difficult and subtle problem that we will not wade into in detail. Rather, our goal is to discuss some general features of the problem that can be deduced given some reasonable simplifying assumptions.

When discussing localization, one can imagine either increasing the system size or increasing the disorder strength. In the latter case, the situation rapidly becomes unwieldy: as the disorder strength approaches the interaction energy ($\delta m/U \rightarrow 1$) or the disorder correlations become ultra-short-ranged ($a_M/\xi_{\text{dis}} \rightarrow 1$), our Ising formulation breaks

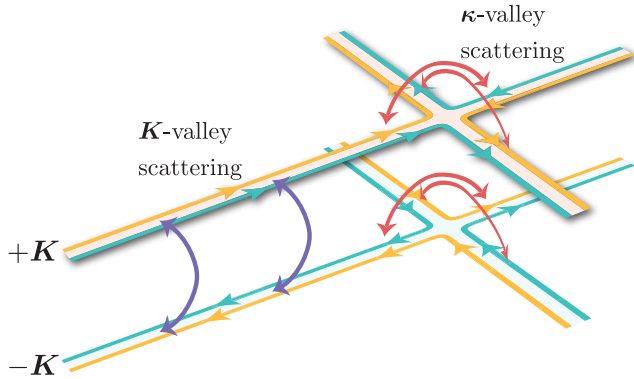


FIG. 6. Edge modes of the $+K$ and $-K$ valley sectors at a node connecting four domains. The orange and turquoise arrows represent the chiral modes at the domain boundaries. Red arrows at the node indicate $U(1)_v$ -preserving, inter- κ -valley scattering processes, which result from inhomogeneities at the moiré lattice scale, $a_M \approx 12.8$ nm. The $U(1)_v$ -breaking inter- K -valley scattering events are indicated by the purple arrows. While this type of scattering is exponentially suppressed [see Eq. (8)], it can occur at any point along a domain boundary.

down. By contrast, the Ising-model perspective remains valid when we instead consider progressively larger samples with an otherwise identical set of parameters. Interactions can still of course pose complications; for instance, in the network-model picture localization involves a network of gapless domain-wall modes that generically form Luttinger liquids [78]. We do not address such subtleties, instead postulating that the primary effect of interactions is to catalyze the spontaneous breaking of $C_2\mathcal{T}$.

The most straightforward manner by which the re-emergent Dirac fermions can localize is through inter- κ -valley scattering. Such scattering events can also localize the original Dirac cones that appear in the free-fermion band structure for mTBG, but the physics is not quite identical: the network picture underlying the re-emergent Dirac cones effectively postpones localization by renormalizing the UV scale at which it occurs. That is, if $\xi_{loc,fr}$ is the localization length in the free case, we have $\xi_{loc} \sim \xi_{dom}\xi_{loc,fr}/a_M$ with interactions. One can intuitively understand this rescaling from the perspective of the gapless domain-wall modes in the network model. As Fig. 6 illustrates, in a given K -valley, the domain-wall modes corresponding to $\pm\kappa$ co-propagate, and hence non-forward-scattering processes can only occur at nodes where multiple domain walls intersect (see red arrows).

Inter- K -valley scattering can also prompt localization [25]. Disorder coupling the two K -valleys has so far been completely ignored since it is exponentially suppressed relative to intra- κ -scattering [see Eq. (8)]. However, inter- K -valley scattering can occur at any point along the domain walls, as illustrated in Fig. 6, making it a fundamentally one-dimensional process. For very large domains, such intra-domain-wall scattering thus inevitably becomes the domi-

nant localization mechanism. The localization length is then expected to be proportional to the mean free path of the domain-wall modes [79], which is $\xi_{loc} \sim \hbar v_F/g_{KK'} \sim \hbar v_F e^{4\pi^2 \xi_{dis}^2/a^2}/g$, and hence an exponentially large function of the disorder correlation length.

VI. DISCUSSION

We have presented a theory that reconciles the seemingly conflicting experiments on charge-neutral mTBG by invoking a nontrivial interplay between strong interactions and weak disorder. In our proposed picture, uniform order (QVH or otherwise) is realized throughout ultra-homogeneous samples, like those of Lu *et al.* [14], whereas QVH domains with opposite spin/valley Chern numbers appear in systems with more disorder, like the experiments of Cao *et al.* [1, 2] and Yankowitz *et al.* [9]. In the latter samples, gapless edge modes at domain boundaries form a network that may be mapped onto a theory of massless Dirac fermions, thereby explaining their semimetallic transport measurements. By contrast, since a physical sample boundary strongly breaks the $U(1)_v$ symmetry protecting the edge modes, a uniformly ordered QVH state is an insulator at charge neutrality, in agreement with the observations of Lu *et al.* Both sample classes exhibit a local gap determined by the interaction strength—in harmony with STM experiments [15–18].

The network model outlined in this paper is somewhat reminiscent of proposals aimed at describing ‘minimally’ twisted bilayer graphene (minTBG) [80–83]. When $\theta \lesssim 1^\circ$, it becomes energetically favourable for the microscopic lattices to distort such that the AB and BA regions occupying the moiré honeycomb sites enlarge at the expense of the AA regions situated at the centre of each moiré hexagon [42, 84–89]; see Fig. 3(a) for an illustration of the undistorted case. Under the application of a displacement field, AB and BA regions develop QVH order with opposing Chern numbers [29–31], yielding four edge modes per spin at the AB/BA boundaries. While both our theory for mTBG and the theory proposed for minTBG are built on network models comprised of QVH domains, there are key qualitative distinctions that we wish to underscore. The local QVH order in minTBG arises entirely as a single-particle effect, whereas the development of QVH order in our scenario relies principally on strong interactions. Moreover, the shape and size of the AB and BA regions in minTBG are fixed; together, they comprise a single moiré unit cell. The QVH domains discussed in this paper instead result from the smooth disorder background and typically extend over many moiré unit cells.

Our proposal is supported by available experimental data and crucially can be further tested in future experiments. One natural direction is to employ large-area STM scans to locally probe both gapped domains *and* gapless domain-wall modes. (To our knowledge evidence of the latter in mTBG has not yet been reported in the literature.) Samples that are simultaneously amenable to STM and transport would

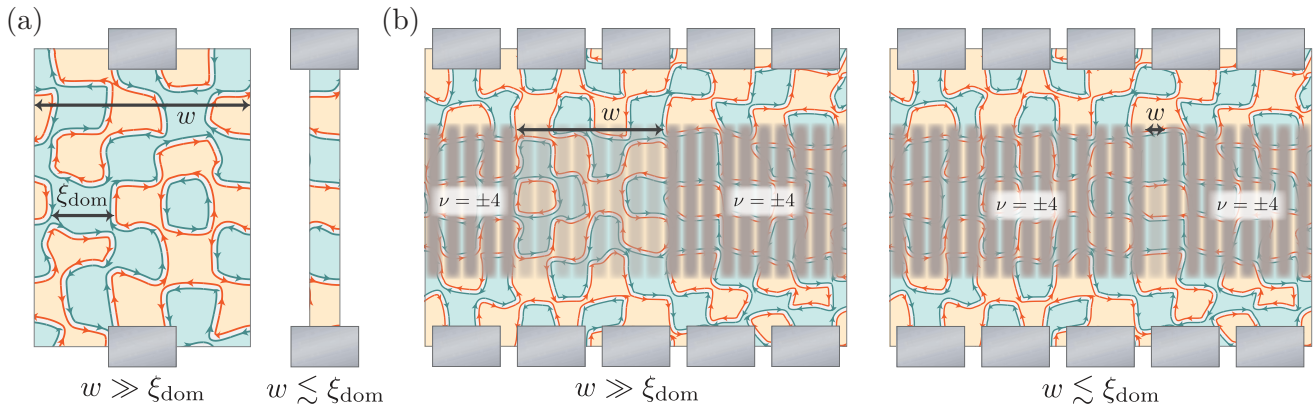


FIG. 7. Schematic illustration of a proposed transport experiment. The domain structure of one of the valley sectors is shown: regions carrying Chern number $C = +1$ are depicted in blue, whereas those carrying $C = -1$ are depicted in orange. The other valley sector is not shown explicitly. The grey rectangles above and below the sample indicate contacts through which the conductance is measured. (a) A single mTBG sample is sliced into multiple sub-systems of varying width w . When $w \gg \xi_{\text{dom}}$, as shown on the left, Dirac-like conductance is observed. Once the width is smaller than the typical domain size, $w \lesssim \xi_{\text{dom}}$, the sample appears insulating, as shown on the right. (b) An alternate experiment in which the mTBG sample remains intact. The taupe rectangles aligned in a row along the centre of the sample represent individually tunable gates through which the chemical potential may be locally varied. In regions where these gates are opaque, the chemical potential lies within the superlattice bandgap, *i.e.* the flat bands are either completely empty or full ($\nu = \pm 4$). The system is tuned to charge neutrality in all other regions (either the gates are transparent or no gates are shown). On the left, $w \gg \xi_{\text{dom}}$, and a semimetallic conductance should be observed. Conversely, since $w \lesssim \xi_{\text{dom}}$, a large resistance is expected on the right.

offer additional insight; for instance, the presence of gapless domain-wall modes should correlate with semimetallic transport, whereas such modes should be absent in homogeneous insulating samples. Some caveats are warranted, however. First, discussions of local phenomena in STM measurements are often complicated, *e.g.*, by disorder- or tip-induced localized states, and it may be difficult to unambiguously distinguish the domain physics we propose from such effects. Additionally, the samples studied by Refs. 1, 2, 9, and 14 are enclosed on both sides by hBN, preventing STM study. The nature of transport in mTBG with hBN only on the bottom side, as in the samples studied by STM in Refs. 15–18, poses an interesting open question.

One can also investigate our scenario entirely within transport [90]. Consider a single mTBG sample etched into a series of strips of varying widths w , as shown in Fig. 7(a). Transport through a given strip depends sensitively on the value of w relative to the typical domain size, ξ_{dom} . When $w \gg \xi_{\text{dom}}$ —the limit presumably relevant to the experiments of Refs. 1, 2, and 9—semimetallic transport should occur. In the opposite limit, $w \ll \xi_{\text{dom}}$, no edge modes connect the contacts and the strip should appear insulating. This experiment may be modified to preclude possible variations in the conductivity resulting from intrinsic variations between the strips, such as their local twist angle. Instead of physically cutting the sample, a ‘strip’ can be electrostatically generated through spatially varying gate voltages: Within a channel of width w , the system is locally tuned to charge neutrality, whereas elsewhere the Fermi energy is tuned to lie within the gap separating the flat and dispersive bands. One could then study the conductivity as a function of width

w for all regions within the sample. Figure 7(b) illustrates this refined version of the experiment.

Finally, although we have emphasized the validity of our proposal even in the eventuality that the clean system does not realize a QVH insulator, the observation of this state at charge neutrality in ultra-clean, insulating samples would of course serve as support for our scenario. Un-quantized valley Hall currents have been observed through non-local resistance measurements [91] in hBN-aligned graphene [26] and Bernal stacked graphene with an out-of-plane electric field [27–29], and the realization of the QVH insulator in TBG at charge neutrality could therefore be demonstrated through an analogous set of measurements. The application of a small magnetic field parallel to the sample represents a more common probe of the correlated insulating states in twisted bilayer graphene. Because the QVH state is a spin singlet, this field is unlikely to have a substantial effect. Nevertheless, ignoring both interactions and modifications of the interlayer tunnelling due to the magnetic field, a mean field perspective would suggest that the gap at charge neutrality should decrease with increasing magnetic field.

Our proposal also spotlights various other avenues for future study. The fate of the network under an applied magnetic field poses a particularly interesting problem. One possibility is that the magnetic field simply stabilizes a different competing phase, thereby destroying the network. If this transition occurs at fields strengths close to or below 1 T, where quantum oscillations are first clearly resolved, the re-emergent Dirac theory is unlikely to produce observable Landau-fan phenomena. The occurrence of such a transition is neither necessary nor expected, however. In the case

where the QVH network survives a broader magnetic-field window, there are two limits to consider. When the magnetic length ℓ_B far exceeds the typical domain size ξ_{dom} , quantum oscillations are expected to be insensitive to the re-emergent nature of the fermions, implying that a Landau fan corresponding to massless Dirac fermions should be observed at low fields. Given that the magnetic length is already quite large at 1 T, $\ell_B \approx 25 \text{ nm} \approx 2a_M$, this regime may be difficult to access experimentally (recall that the UV cutoff for our network model was $\xi_{\text{int}} \sim 100 \text{ nm}$). The opposite limit, $\ell_B \ll \xi_{\text{dom}}$, appears to be more subtle. It is conceivable that the gapless edge modes do not affect the quantum oscillations, resulting in a Landau level spectrum similar to that of *massive* Dirac cones associated with the gapped QVH domains. Alternatively, the system could exhibit physics reminiscent of the Hofstadter butterfly [92], though it seems likely that nonuniformity of the domain sizes may hinder any clear signal. Quantifying these issues could shed additional light on the experimental relevance of our scenario.

The role of interactions at the edges of the domains is another topic that we have not touched on. The edge modes may display interesting interacting phenomena that could be studied through the well-controlled bosonization formalism. In fact, Wu *et al.* [81] have analyzed this problem in the context of the minimal twist angle samples described above. Further, while disorder is generically expected to localize the edges, the inclusion of interactions may have nontrivial consequences [78].

We have said little regarding transport away from the charge neutrality point. While its semimetallic nature dictates that the conductivity σ increase with doping, it can do so in different ways. If transport is ballistic, far enough away from charge neutrality, the conductivity should essentially track the density of states: $\sigma \propto \sqrt{|n|}$, where n is the electron density [93, 94]. Provided inter- κ -scattering is the most important form of disorder, we expect the mean free path of the network model to be rescaled, implying that ballistic transport may not be unreasonable. That is, letting $\ell_{\text{mfp},fr}$

be the mean free path of the non-interacting Dirac fermions, we may postulate that the mean free path of the recovered network Dirac fermions is $\ell_{\text{mfp}} \sim \xi_{\text{dom}} \ell_{\text{mfp},fr} / a_M$. On the other hand, in monolayer graphene, the linear dependence of the conductivity on density away from charge neutrality, $\sigma \propto |n|$, is largely ascribed to long-range Coulomb scattering [95–102]. While it seems unlikely that a similar mechanism would play an important role in mTBG, it is possible that twist-angle disorder (which can also be long-range) could have a similar effect [103].

Finally, exploring the interplay between interactions and disorder at other integer fillings constitutes perhaps the most interesting future direction. The charge-neutrality regime that we examined here offers the virtue that the system is ‘almost’ insulating even at the band structure level—thereby facilitating the study of (at least certain) correlated insulators. Accessing correlated insulating states at other fillings requires a far more drastic modification of the band fillings. Generalizing our analysis to such cases could provide valuable insight into the observed phenomenology of mTBG.

ACKNOWLEDGEMENTS

We are grateful to Cory Dean, Arbel Haim, Eslam Khalaf, Stevan Nadj-Perge, Felix von Oppen, Seth Whitsitt, and Andrea Young for illuminating discussions. This work was supported by the Army Research Office under Grant Award W911NF-17-1-0323; the NSF through grant DMR-1723367; the Caltech Institute for Quantum Information and Matter, an NSF Physics Frontiers Center with support of the Gordon and Betty Moore Foundation through Grant GBMF1250; the Walter Burke Institute for Theoretical Physics at Caltech; and the Gordon and Betty Moore Foundation’s EPiQS Initiative, Grant GBMF8682 to JA. This work was performed in part at the Aspen Center for Physics, which is supported by National Science Foundation grant PHY-1607611.

A. CONTINUUM MODEL

We briefly outline the continuum model in this section. Spin indices are completely suppressed below. We first decompose the microscopic graphene operators as

$$\tilde{f}_\ell(\mathbf{r}) = e^{i\mathbf{K}\cdot\mathbf{r}} f_{+,\ell}(\mathbf{r}) + e^{-i\mathbf{K}\cdot\mathbf{r}} f_{-,\ell}(\mathbf{r}), \quad (\text{A1})$$

where ℓ indicates both layer and sublattice. As discussed in Sec. II A, the continuum model Hamiltonian decouples into \mathbf{K} -valley sectors $H_{\text{cont}} = H_+ + H_-$, where H_\pm act on $f_{\pm,\ell}$. For the moment, we consider H_+ . We express $f_{+,\ell}$ as a vector $(f_{t,A}(\mathbf{r}), f_{t,B}(\mathbf{r}), f_{\beta,A}(\mathbf{r}), f_{\beta,B}(\mathbf{r}))$, where the ‘+’ has been dropped for convenience, t, β denote

layer, and A, B denote sublattice. In this basis, H_+ acts as

$$H_+ = \begin{pmatrix} iv_0 \boldsymbol{\eta}_{\theta/2} \cdot \boldsymbol{\nabla} & T(\mathbf{r}) \\ T^\dagger(\mathbf{r}) & iv_0 \boldsymbol{\eta}_{-\theta/2} \cdot \boldsymbol{\nabla} \end{pmatrix}, \quad (\text{A2})$$

where $\boldsymbol{\eta}_\phi = e^{-i\phi\eta^z/2} (\eta^x, \eta^y) e^{i\phi\eta^z/2}$ act on the sublattice space and $\boldsymbol{\nabla} = (\partial_x, \partial_y)$. The tunnelling matrix $T(\mathbf{r})$ is given by

$$T(\mathbf{r}) = \sum_{\ell=1,2,3} t_\ell e^{-i\mathbf{q}_\ell \cdot \mathbf{r}}, \quad \mathbf{q}_\ell = \mathcal{R}_{2\pi(\ell-1)/3} [\mathbf{K}_t - \mathbf{K}_\ell] \quad (\text{A3})$$

where $\mathcal{R}_\phi[\mathbf{v}]$ rotates the vector \mathbf{v} by ϕ and the matrices t_ℓ are defined through

$$t_\ell = e^{2\pi i(\ell-1)\eta^z/3} \begin{pmatrix} w_0 & w_1 \\ w_1 & w_0 \end{pmatrix} e^{-2\pi i(\ell-1)\eta^z/3}. \quad (\text{A4})$$

The physical parameters of the model are the twist angle θ , the velocity of the microscopic graphene layers v_0 , and the tunnelling amplitudes, w_0 and w_1 . We take the angle to be close to the magic angle, $\theta = 1.05^\circ$, and the graphene velocity to be $v_0 = 9.1 \times 10^5$ m/s. The tunnelling amplitudes are typically taken to be $(w_0, w_1) = (85, 110)$ meV [4]. However, for the chiral version [see Sec. D 5] of the model, we set $w_0 = 0$, keeping $w_1 = 110$ meV [64].

The Hamiltonian corresponding to the other valley, \mathcal{H}_- , may be obtained by acting time-reversal (\mathcal{T}) or by rotating by 180° (C_2).

The continuum Hamiltonians maybe also be expressed in momentum space. Returning to second quantized notation, it may be written

$$H_\mu = \sum_{\mathbf{G}, \mathbf{G}', \ell, \ell'} \int_{\mathbf{k} \in \text{BZ}} f_{\mu, \ell}^\dagger(\mathbf{k} + \mathbf{G}) H_{\mathbf{G}, \ell; \mathbf{G}', \ell'}^{(\mu)}(\mathbf{k}) f_{\mu, \ell'}(\mathbf{k} + \mathbf{G}'), \quad (\text{A5})$$

where $\mu = +, -$ labels the \mathbf{K} -valley and the \mathbf{G} s are moiré reciprocal lattice vectors. Here, $H^{(\mu)}(\mathbf{k})$ may be thought of as an infinite matrix taking values within the moiré BZ with indices (\mathbf{G}, ℓ) . It can be diagonalized through the unitary rotation

$$c_{\mu, i}^\dagger(\mathbf{k}) = \sum_{\mathbf{G}, \ell} u_{\mu, i; \mathbf{G}, \ell}(\mathbf{k}) f_{\mu, \ell}^\dagger(\mathbf{k} + \mathbf{G}), \quad f_{\mu, \ell}^\dagger(\mathbf{k} + \mathbf{G}) = \sum_i u_{\mu, i; \mathbf{G}, \ell}^*(\mathbf{k}) c_{\mu, i}^\dagger(\mathbf{k}), \quad (\text{A6})$$

where i indexes the band. In terms of the $c_{\mu, i}(\mathbf{k})$ operators H_μ is

$$H_\mu = \sum_i \int_{\mathbf{k} \in \text{BZ}} c_{\mu, i}^\dagger(\mathbf{k}) \epsilon_i(\mathbf{k}) c_{\mu, i}(\mathbf{k}). \quad (\text{A7})$$

We note that invariance of $c_{\mu, i}(\mathbf{k})$ under shifts of \mathbf{k} by a reciprocal lattice vector, $\mathbf{k} \rightarrow \mathbf{k} + \mathbf{G}$, implies $u_{\mu, i; \mathbf{G}, \ell}(\mathbf{k} + \mathbf{G}') = u_{\mu, i; \mathbf{G} + \mathbf{G}', \ell}(\mathbf{k})$.

B. SPIN AND TIME-REVERSAL SYMMETRIC BILINEARS

Here, we enumerate some of the symmetries of the Dirac theory. It is convenient to express them in terms of a large, unphysical, SU(8) symmetry generated by the \mathbf{K} -valley, κ -valley, and spin symmetries. The generators of these symmetries are

$$\text{SU}(2)_s : (\sigma^x, \sigma^y, \sigma^z), \quad \text{SU}(2)_\kappa : (\tau^x, \tau^y, \tau^z), \quad \text{SU}(2)_K : (\mu^x \eta^y, \mu^y \eta^x, \mu^z). \quad (\text{B1})$$

Since the $\text{SU}(2)_K$ triplet does not take a particularly simple form, we define $\bar{\mu}^i = (\mu^x \eta^x, \mu^y \eta^y, \mu^z)$. Finally, the γ -matrices are $\gamma^\mu = (\mu^z \eta^z, i\eta^y, -i\mu^z \eta^x)$. By combining the γ^μ , σ^i , τ^i , $\bar{\mu}^i$, we can generate all bilinears (pairing terms are not considered).

We are interested only in those bilinears that preserve the spin and time-reversal symmetries. Clearly, spin-conservation requires that disorder not couple to any bilinear containing σ^i , so we ignore it completely, treating Ψ as a spinless fermion. Time reversal then acts as

$$\mathcal{T} : \quad \Psi \rightarrow \mu^x \tau^x \Psi, \quad i \rightarrow -i. \quad (\text{B2})$$

It follows that the $\text{SU}(2)_\kappa$, $\text{SU}(2)_K$ triplets and the γ -matrices map as

$$\begin{aligned} \mathcal{T} : \quad & (\tau^x, \tau^y, \tau^z) \rightarrow (\tau^x, \tau^y, -\tau^z), \\ & (\mu^x \eta^y, \mu^y \eta^y, \mu^z) \rightarrow -(\mu^x \eta^y, \mu^y \eta^y, \mu^z), \\ & (\gamma^0, \gamma^x, \gamma^y) \rightarrow (-\gamma^0, \gamma^x, \gamma^y) \end{aligned} \quad (\text{B3})$$

These transformation properties result in the following time-reversal invariant bilinears:

$$\begin{aligned} \bar{\Psi} M \Psi, \quad & M \in \{\bar{\mu}^i, \bar{\mu}^i \tau^{x,y}, \tau^z\}, \\ \bar{\Psi} \gamma^0 M \Psi, \quad & M \in \{\mathbb{1}, \tau^{x,y}, \tau^z \bar{\mu}^i\}, \\ \bar{\Psi} \gamma^{x,y} M \Psi, \quad & M \in \{\bar{\mu}^i, \bar{\mu}^i \tau^{x,y}, \tau^z\}, \end{aligned} \quad (\text{B4})$$

where $\bar{\Psi} = \Psi^\dagger \gamma^0$. We are most concerned the mass bilinears, shown on the first line. We note that $\bar{\Psi} \bar{\mu}^i \Psi = \Psi^\dagger (\mu^y \eta^z, -\mu^x \eta^x, \eta^z) \Psi$. The last term, $\Psi^\dagger \eta^z \Psi$, is the order parameter for the QVH state. For completeness we also list the bilinears that break time-reversal symmetry:

$$\begin{aligned} \bar{\Psi} M \Psi, \quad & M \in \{\mathbb{1}, \tau^{x,y}, \tau^z \bar{\mu}^i\}, \\ \bar{\Psi} \gamma^0 M \Psi, \quad & M \in \{\bar{\mu}^i, \bar{\mu}^i \tau^{x,y}, \tau^z\}, \\ \bar{\Psi} \gamma^{x,y} M \Psi, \quad & M \in \{\mathbb{1}, \tau^{x,y}, \tau^z \bar{\mu}^i\}. \end{aligned} \quad (\text{B5})$$

C. SUPPRESSION OF INTER- K -VALLEY SCATTERING

We briefly outline a schematic argument for the exponential suppression of inter- K -valley scattering processes. We begin by considering the operators on the microscopic graphene lattice. Suppose disorder couples as

$$H_{\text{micro}} = \sum_{\ell, \ell'} \int_{\mathbf{r}} \mathcal{R}(\mathbf{r}) \tilde{f}_{\ell}^{\dagger}(\mathbf{r}) T_{\ell\ell'} \tilde{f}_{\ell'}(\mathbf{r}) \quad (\text{C1})$$

Here, ℓ labels both the layer and sublattice of the fermion $\tilde{f}_{\ell}(\mathbf{r})$, $T_{\ell\ell'}$ is a matrix whose precise form is unimportant, and $\mathcal{R}(\mathbf{r})$ is the disorder field with values drawn from a Gaussian probability distribution:

$$\overline{\mathcal{R}(\mathbf{r})} = 0, \quad \overline{\mathcal{R}(\mathbf{r})\mathcal{R}(\mathbf{r}')^*} = g^2 e^{-(\mathbf{r}-\mathbf{r}')^2/(2\xi_{\text{dis}}^2)}. \quad (\text{C2})$$

In momentum space, we find

$$H_{\text{micro}} = \int_{\mathbf{k}} \mathcal{R}(\mathbf{q}) \tilde{f}_{\ell}^{\dagger}(\mathbf{k}) T_{\ell\ell'} \tilde{f}_{\ell'}(\mathbf{k} + \mathbf{q}), \quad (\text{C3})$$

where

$$\overline{\mathcal{R}(\mathbf{q})\mathcal{R}^*(\mathbf{q}')^*} = \delta^2(\mathbf{q} - \mathbf{q}') g^2 \xi_{\text{dis}}^2 e^{-\mathbf{q}^2 \xi_{\text{dis}}^2/2}. \quad (\text{C4})$$

We now wish to expand about the $+\mathbf{K}$ and $-\mathbf{K}$ points. Letting $\tilde{f}_{n=\pm, \ell}(\mathbf{k}) = f_{\ell}(\pm\mathbf{K} + \mathbf{k})$, the Hamiltonian divides into two pieces

$$\begin{aligned} H_{KK} &= \sum_{n=\pm} \int_{\mathbf{k}, \mathbf{q}} \mathcal{R}(\mathbf{q}) f_n^{\dagger}(\mathbf{k}) T f_n(\mathbf{k} + \mathbf{q}), \\ H_{KK'} &= \int_{\mathbf{k}, \mathbf{q}} \sum_{j=1}^3 \mathcal{R}(\mathbf{q} + \mathbf{Q}_j) f_{+}^{\dagger}(\mathbf{k}) T f_{-}(\mathbf{k} + \mathbf{q}) + h.c. \end{aligned} \quad (\text{C5})$$

where \mathbf{Q}_j are the three (smallest) momenta such that $-\mathbf{K} + \mathbf{Q}_j = +\mathbf{K}$, each of which has magnitude $|\mathbf{K}| = 4\pi/3a$, where a is the lattice constant of monolayer graphene. We have also suppressed the summation over the ℓ indices of the fermions and matrix T . Letting $\mathcal{R}_{(+)}(\mathbf{q}) = \sum_j \mathcal{R}(\mathbf{q} + \mathbf{Q}_j)$, we then see

$$\begin{aligned} \overline{\mathcal{R}_{(+)}(\mathbf{q})\mathcal{R}_{(+)}^*(\mathbf{q}')^*} &= \delta^2(\mathbf{q} - \mathbf{q}') \xi_{\text{dis}}^2 g^2 \sum_j e^{-(\mathbf{q} + \mathbf{Q}_j)^2 \xi_{\text{dis}}^2/2} \\ &= \delta^2(\mathbf{q} - \mathbf{q}') \xi_{\text{dis}}^2 g^2 e^{-\mathbf{K}^2 \xi_{\text{dis}}^2/2} e^{-\mathbf{q}^2 \xi_{\text{dis}}^2/2} \sum_j e^{-\mathbf{q} \cdot \mathbf{Q}_j \xi_{\text{dis}}^2}. \end{aligned} \quad (\text{C6})$$

Ignoring the anisotropic term on the right, the disorder field corresponding $\mathbf{K} \rightarrow -\mathbf{K}$ scattering has the same correlation length, ξ_{dis} , but with an exponentially suppressed amplitude: $g_{KK'} \sim g e^{-4\pi^2 \xi_{\text{dis}}^2/a^2}$.

These arguments may appear to carry over directly to the case of inter- κ -scattering, *i.e.*, we may wish to conclude that the typical inter- κ valley scattering amplitude $g_{\kappa\kappa'}$ is exponentially suppressed relative to the typical intra- κ

scattering amplitude g : $g_{\kappa\kappa'} \sim ge^{-4\pi^2\xi_{\text{dis}}^2/a_M^2}$. However, in this case, there are additional subtleties to take into account. While the continuum Hamiltonian does not mix the f fermions on the scale of the large BZ, $\sim 1/a$, they *are* mixed on the scale of the moiré BZ, $\sim 1/a_M$. In particular, the flat band operator c (or, equivalently, the Dirac operator ψ) at a momentum quantum number \mathbf{k} in the moiré BZ is composed of a superposition of f fermions with momenta $\mathbf{k} + \mathbf{G}$ (in the microscopic BZ), where the \mathbf{G} s are moiré reciprocal lattice vectors, as indicated in Eq. (A6). With the exception of $\mathbf{G} = 0$, all such reciprocal lattice vectors are already of order $|\kappa|$ or larger. Importantly, this mixing is responsible for the very flatness of the bands and therefore constitutes a nonnegligible effect. As a result, unless ξ_{dis} is much, much larger than a_M , these higher moments may nevertheless contribute substantially to the $\kappa \rightarrow -\kappa$ scattering processes. We therefore emphasize that the analysis and proposal presented in this paper is *not* predicated on the assumption that $g_{\kappa\kappa'}$ is small.

D. MEAN FIELD ANALYSIS OF INSULATING PHASES

Based on numerical results, we argued in Sec. IV B that the ground state of a single flavour theory with interactions is a Chern insulator. Upon including valley and flavour indices in Sec. IV C, we identified four natural insulating states distinguished by their symmetry action, as summarized in Table I. Further, we noted that only the order parameter for the QVH insulator could couple to disorder, which is vital for the scenario we propose.

Here, we discuss the circumstances under which the QVH insulator is or is not energetically preferred compared to the QSH, QH, and QSVH. We determine the band structure using the continuum model (see Appendix A), and, in spite of the concerns raised at the end of Sec. IV A, we model the interactions using H_C , as written in Eq. (13). Moreover, to further simplify the calculation, we project H_C onto the flat bands, a simplification that may admittedly neglect relevant contributions from the dispersive bands. We therefore view this exercise mainly as a guide intended to expose trends rather than provide rigorous quantitative energetics. Nevertheless, we show that within a simple mean field analysis, the Fock terms are not expected to distinguish these phases. While it appears that the Hartree terms favour the QSH, QSVH, and QH insulators over the QVH phase, we find that this preference is not the case for the chiral model [64]—they remain degenerate. We next calculate the energy difference between the QVH and other phases numerically for a more realistic set of parameters and demonstrate that while the energy difference is no longer zero, it remains negligibly small.

1. Flat band projection

The Hamiltonian H_C of Eq. (13) is still quite complicated: it includes all bands of the model, whereas we are only interested in what happens to the flat bands. Since these bands are separated from the dispersive bands by a gap E_g by assumption, the latter states can be integrated out to give an effective Hamiltonian acting only on the flat band subspace. The leading order contribution is obtained simply by projecting H_C to the flat bands:

$$H_{C,1} = \int_{\mathbf{q}_{\text{small}}} V(\mathbf{q}) \rho_{f\bar{f}}(\mathbf{q}) \rho_{f\bar{f}}(-\mathbf{q}), \quad (\text{D1})$$

where $\rho_{f\bar{f}}(\mathbf{q})$ is the density operator projected onto the flat bands.

We show that the mean field decoupling of $H_0 + H_{C,1}$ [where H_0 is given in Eq. (2)] are independent of the sign of the Dirac mass. To do so, we define the variational Hamiltonian $H_{\text{MF}}(\{M_\mu\}) = \sum_\mu H_{\text{MF}}^{(\mu)}(M_\mu)$ where $\mu = (n, \alpha)$ sums over both \mathbf{K} -valleys, $n = \pm$, and spin, $\alpha = \uparrow, \downarrow$. The individual mean field Hamiltonians are

$$H_{\text{MF}}^{(\mu)}(M_\mu) = \int_{\mathbf{k} \in \text{BZ}} c_\mu^\dagger(\mathbf{k}) \underbrace{\left[h_\mu(\mathbf{k}) + M_\mu \eta^z \right]}_{\bar{h}_\mu(\mathbf{k}; M_\mu)} c_\mu(\mathbf{k}),$$

$$h_\mu(\mathbf{k}) = h_{\mu,0}(\mathbf{k}) + h_{\mu,x}(\mathbf{k})\eta^x + h_{\mu,y}(\mathbf{k})\eta^y. \quad (\text{D2})$$

We study the dependence of $\langle \{M_\mu\} | H_0 + H_{C,1} | \{M_\mu\} \rangle$ on the signs of M_μ , where $|\{M_\mu\}\rangle$ denotes the ground-state of $H_{\text{MF}}(\{M_\mu\})$.

2. Density operator and form factors

One complication of this calculation is the presence of form factors in the definition of the densities and thus $H_{C,1}$ as well. In particular, we have

$$\rho_{fl}(\mathbf{q}) = \sum_\mu \rho_\mu(\mathbf{q}),$$

$$\rho_\mu(\mathbf{q}) = \sum_{\mu,\ell} \int_{\mathbf{k} \text{ small}} f_{\mu,\ell}^\dagger(\mathbf{k}) f_{\mu,\ell}(\mathbf{k} + \mathbf{q}), \quad (\text{D3})$$

where $f_{\mu,\ell}(\mathbf{k}) = f_{n=\pm, \alpha, \ell}(\mathbf{k})$ denotes the electron operator with spin $\alpha = \uparrow, \downarrow$ and total momentum $\pm \mathbf{K} + \mathbf{k}$. As in Sec. IV and Appendix A, ℓ labels both layer and sublattice. In what follows we omit the label “*fl.*” Recall that neither the momentum of the density operator, \mathbf{q} , nor the momentum being summed over, \mathbf{k} , is required to lie within the moiré Brillouin zone. We therefore instead write

$$\rho_\mu(\mathbf{q} + \mathbf{G}') = \int_{\mathbf{k} \in \text{BZ}} \sum_{\mathbf{G}, \ell} f_{\mu,\ell}^\dagger(\mathbf{k} + \mathbf{G}) f_{\mu,\ell}(\mathbf{k} + \mathbf{q} + \mathbf{G} + \mathbf{G}'), \quad (\text{D4})$$

where \mathbf{G} and \mathbf{G}' are moiré reciprocal lattice vectors and both \mathbf{k} and \mathbf{q} lie within the moiré BZ. Using the definition of $c_{\mu,i}$ in relation to $f_{\mu,\ell}$ given in Eq. (A6), the density may now be expressed directly in terms of the flat band creation and annihilation operators:

$$\rho_\mu(\mathbf{q} + \mathbf{G}) = \sum_{ij \in fl} c_{\mu,i}^\dagger(\mathbf{k}) \lambda_{\mu;ij}(\mathbf{k}, \mathbf{k} + \mathbf{q} + \mathbf{G}) c_{\mu,j}(\mathbf{k} + \mathbf{q}),$$

$$\lambda_{\mu;ij}(\mathbf{k}, \mathbf{k} + \mathbf{q} + \mathbf{G}) = \sum_{\mathbf{G}', \ell} u_{\mu,i; \mathbf{G}', \ell}^*(\mathbf{k}) u_{\mu,j; \mathbf{G}', \ell}(\mathbf{k} + \mathbf{q} + \mathbf{G}). \quad (\text{D5})$$

We frequently refer to the functions $\lambda_{\mu,ij}$ as ‘form factors’ in what follows. We have used the fact that the band operators are invariant under reciprocal lattice translations up to a phase, $c_{\mu,j}(\mathbf{p} + \mathbf{G}) = e^{i\phi} c_{\mu,j}(\mathbf{p})$. From the fact that $u_{\mu,i; \mathbf{G}, \ell}(\mathbf{k} + \mathbf{G}') = u_{\mu,i; \mathbf{G} + \mathbf{G}', \ell}(\mathbf{k})$, we also have $\lambda_{\mu;ij}(\mathbf{k}, \mathbf{k}' + \mathbf{G}) = \lambda_{\mu,ij}(\mathbf{k} - \mathbf{G}, \mathbf{k}')$. Finally, with this

notation, the flat-band Coulomb interaction is

$$H_{C,1} = \int_{\mathbf{q}, \mathbf{k}, \mathbf{k}'} \sum_{\mathbf{G}} \sum_{\mu, \nu} c_{\mu}^{\dagger}(\mathbf{k}) \lambda_{\mu}(\mathbf{k}, \mathbf{k} + \mathbf{q} + \mathbf{G}) c_{\mu}(\mathbf{k} + \mathbf{q}) \cdot c_{\nu}^{\dagger}(\mathbf{k}' + \mathbf{q}) \lambda_{\nu}(\mathbf{k}' + \mathbf{q} + \mathbf{G}, \mathbf{k}') c_{\nu}(\mathbf{k}'). \quad (\text{D6})$$

3. Symmetry constraints

We begin by discussing the symmetry properties of the mean-field kernel $\bar{h}_{\mu}(\mathbf{k}; M_{\mu})$. We begin with the symmetry transformations

$$\begin{aligned} \mathcal{T} : \quad & c(\mathbf{k}) \rightarrow \mu^x c(-\mathbf{k}), \\ \mathcal{C}_2 \mathcal{T} : \quad & c(\mathbf{k}) \rightarrow \eta^x c(\mathbf{k}), \end{aligned} \quad (\text{D7})$$

where μ^x acts on the \mathbf{K} -valley indices and η^x acts on the (flat) band indices. Both are anti-Hermitian, taking $i \rightarrow -i$. In terms of the mean field Hamiltonian, they imply

$$\bar{h}_{+, \alpha}(\mathbf{k}; M) = \bar{h}_{-, \alpha}^*(-\mathbf{k}; M), \quad \bar{h}_{\mu}(\mathbf{k}; M) = \eta^x \bar{h}_{\mu}^*(\mathbf{k}; -M) \eta^x. \quad (\text{D8})$$

Obviously, since $h_{\mu}(\mathbf{k}) = \bar{h}_{\mu}(\mathbf{k}; M = 0)$, these relations also hold for the non-interacting part of the flat band Hamiltonian.

We now define the projector

$$P_{\mu; ij}(\mathbf{k}; M) = \left\langle c_{\mu, j}^{\dagger}(\mathbf{k}) c_{\mu, i}(\mathbf{k}) \right\rangle_M. \quad (\text{D9})$$

The subscript M is used as a shorthand to denote which mean field Hamiltonian the ground state begin used to compute the expectation value is associated with. The equalities of Eq. (D8) then imply

$$P_{+, \alpha}(\mathbf{k}; M) = P_{-, \alpha}^T(-\mathbf{k}; M), \quad (\text{D10a})$$

$$P_{\mu}(\mathbf{k}; M) = \eta^x P_{\mu}^T(\mathbf{k}; -M) \eta^x. \quad (\text{D10b})$$

Note that $P_{\mu}^{\dagger}(\mathbf{k}; M) = P_{\mu}(\mathbf{k}; M)$. Similarly, we find that the form factors must satisfy

$$\lambda_{+, \alpha}(\mathbf{k}, \mathbf{k} + \mathbf{q}) = \lambda_{-, \alpha}^T(-\mathbf{k} - \mathbf{q}, -\mathbf{k}), \quad (\text{D11a})$$

$$\lambda_{\mu}(\mathbf{k}, \mathbf{k} + \mathbf{q}) = \eta^x \lambda_{\mu}^T(\mathbf{k} + \mathbf{q}, \mathbf{k}) \eta^x. \quad (\text{D11b})$$

4. Evaluation of mean field Hamiltonian

We wish to compute the expectation value $\langle \{M_{\mu}\} | H_0 + H_{C,1} | \{M_{\mu}\} \rangle$. This function may be separated into three pieces:

$$\langle \{M_{\mu}\} | H_0 + H_{C,1} | \{M_{\mu}\} \rangle = \langle H_0 \rangle_{\{M_{\mu}\}} + H_F(\{M_{\mu}\}) + H_H(\{M_{\mu}\}), \quad (\text{D12})$$

where H_F and H_H are the Fock and Hartree decouplings of the Coulomb interaction. These three terms are discussed in the following subsections.

a. Quadratic term: $\langle H_0 \rangle$

We write the quadratic part of the Hamiltonian as a sum over the valleys and spins, $H_0 = \sum_\mu H_0^{(\mu)}$, where

$$H_0^{(\mu)} = \int_{\mathbf{k}} c_\mu^\dagger(\mathbf{k}) h_\mu(\mathbf{k}) c_\mu(\mathbf{k}). \quad (\text{D13})$$

The kernel $h_\mu(\mathbf{k})$ is defined in Eq. D2. Taking the expectation value, we find

$$\langle H_0^{(\mu)} \rangle_{M_\mu} = \int_{\mathbf{k}} \text{tr}[P_\mu(\mathbf{k}; M_\mu) h_\mu(\mathbf{k})]. \quad (\text{D14})$$

Inserting the relations given in Eqs. (D10b) and (D11b), we arrive at

$$\begin{aligned} \langle H_0^{(\mu)} \rangle_{M_\mu} &= \int_{\mathbf{k}} \text{tr}[\eta^x P_\mu^T(\mathbf{k}; -M_\mu) \eta^x \eta^x h_\mu^T(\mathbf{k}) \eta^x] \\ &= \langle H_0^{(\mu)} \rangle_{-M_\mu}. \end{aligned} \quad (\text{D15})$$

Hence, we have verified that $\langle H_0 \rangle$ is independent of the signs of the mass terms.

b. Fock term: H_F

The Fock term is

$$\begin{aligned} H_F(\{M_\mu\}) &= \sum_\mu H_F^{(\mu)}(M_\mu), \\ H_F^{(\mu)}(M_\mu) &= - \int_{\mathbf{k}, \mathbf{p}} \sum_{\mathbf{G}} V(\mathbf{p} - \mathbf{k} + \mathbf{G}) \text{tr}[\lambda_\mu(\mathbf{k}, \mathbf{p} + \mathbf{G}) P_\mu^T(\mathbf{p}; M_\mu) \lambda_\mu(\mathbf{p} + \mathbf{G}, \mathbf{k}) P_\mu(\mathbf{k}; M_\mu)]. \end{aligned} \quad (\text{D16})$$

Inserting the relations from Eqs. (D10b) and (D11b), we find

$$\begin{aligned} H_F^{(\mu)}(M_\mu) &= - \int_{\mathbf{k}, \mathbf{p}} \sum_{\mathbf{G}} V(\mathbf{p} - \mathbf{k} + \mathbf{G}) \text{tr}[\lambda_\mu^T(\mathbf{p} + \mathbf{G}, \mathbf{k}) P_\mu^T(\mathbf{p}; -M_\mu) \lambda_\mu^T(\mathbf{k}, \mathbf{p} + \mathbf{G}) P_\mu^T(\mathbf{k}; -M_\mu)] \\ &= H_F^{(\mu)}(-M_\mu). \end{aligned} \quad (\text{D17})$$

We again conclude that the Fock contribution is independent of the sign M_μ takes.

c. *Hartree term: H_H*

The Hartree term can be written

$$H_H(\{M_\mu\}) = \sum_{\mathbf{G}} V(\mathbf{G}) \sum_{\mu,\nu} \langle \rho_\mu(\mathbf{G}) \rangle_{M_\mu} \langle \rho_\nu(-\mathbf{G}) \rangle_{M_\nu}. \quad (\text{D18})$$

We therefore begin by calculating $\langle \rho_\mu(\mathbf{G}) \rangle_M$:

$$\langle \rho_\mu(\mathbf{G}) \rangle_M = \int_{\mathbf{k}} \text{tr}[P_\mu(\mathbf{k}; M) \lambda_\mu(\mathbf{k}, \mathbf{k} + \mathbf{G})]. \quad (\text{D19})$$

We use the constraints imposed by time reversal [Eqs. (D10a) and (D11a)] to relate the expectation values of the densities of the two valleys to one another:

$$\begin{aligned} \langle \rho_{+,\alpha}(\mathbf{G}) \rangle_M &= \int_{\mathbf{k}} \text{tr}[P_{-,\alpha}^T(-\mathbf{k}; M) \lambda_{-,\alpha}^T(-\mathbf{k} - \mathbf{G}, -\mathbf{k})] = \int_{\mathbf{k}} \text{tr}[P_{-,\alpha}(\mathbf{k}; M) \lambda_{-,\alpha}(\mathbf{k}, \mathbf{k} + \mathbf{G})] \\ &= \langle \rho_{-,\alpha}(\mathbf{G}) \rangle_M. \end{aligned} \quad (\text{D20})$$

We see that the expectation value of the density operator is independent of the valley and spin degree of freedom, motivating us to define the function

$$R(M; \mathbf{G}) \equiv \langle \rho_\mu(\mathbf{G}) \rangle_M. \quad (\text{D21})$$

Note that the identity $\rho_\mu(\mathbf{G}) = \rho_\mu^\dagger(-\mathbf{G})$ implies $R(M; \mathbf{G}) = R^*(M; -\mathbf{G})$. The $\mathcal{C}_2\mathcal{T}$ symmetry [Eqs. (D10b) and (D11b)] then gives,

$$\begin{aligned} \langle \rho_\mu(\mathbf{G}) \rangle_M &= \int_{\mathbf{k}} \text{tr}[P_\mu^T(\mathbf{k}; -M) \lambda_\mu^T(\mathbf{k} + \mathbf{G}, \mathbf{k})] = \int_{\mathbf{k}} \text{tr}[P_\mu(\mathbf{k}; -M) \lambda_\mu(\mathbf{k}, \mathbf{k} - \mathbf{G})] \\ &= \langle \rho_\mu(-\mathbf{G}) \rangle_{-M} = \langle \rho_\mu(\mathbf{G}) \rangle_{-M}^*. \end{aligned} \quad (\text{D22})$$

We conclude that $R(-M; \mathbf{G}) = R^*(M; \mathbf{G})$.

The Hartree term is therefore

$$H_H(\{M_\mu\}) = \sum_{\mathbf{G}} V(\mathbf{G}) \left| \sum_{\mu} R(M_\mu; \mathbf{G}) \right|^2. \quad (\text{D23})$$

The relative signs of the mass terms of the four states under consideration are shown in Tab. II. Separating $R(M; \mathbf{G})$ into real and imaginary parts, $R(M; \mathbf{G}) = R'(M; \mathbf{G}) + iR''(M; \mathbf{G})$, we conclude that

$$\begin{aligned} H_H^{\text{QVH}} &= 16 \sum_{\mathbf{G}} V(\mathbf{G}) [R'(M; \mathbf{G})^2 + R''(M; \mathbf{G})^2], \\ H_H^{\text{QSVH}} &= H_H^{\text{QH}} = H_H^{\text{QSH}} = 16 \sum_{\mathbf{G}} V(\mathbf{G}) R'(M; \mathbf{G})^2. \end{aligned} \quad (\text{D24})$$

	$M_{+, \uparrow}$	$M_{+, \downarrow}$	$M_{-, \uparrow}$	$M_{-, \downarrow}$
QVH	1	1	1	1
QSVH	1	-1	1	-1
QH	1	1	-1	-1
QSH	1	-1	-1	1

TABLE II. Relative signs of the mass terms corresponding to the four phases depicted in Fig. 4.

It follows that the QVH state is *higher* in energy than the other three insulating states by $16 \sum_{\mathbf{G}} V(\mathbf{G}) R''(M; \mathbf{G})^2$.

We note that since $\lambda(\mathbf{k}, \mathbf{k}) = \mathbb{1}$, for $\mathbf{G} = 0$ we necessarily have $R''(M; \mathbf{0}) = 0$, implying that for this term at least, there is no difference in energy between the QVH insulator and the other three. In a typical tight-binding model, the $\mathbf{G} = 0$ term accounts for the entirety of the Hartree energy. For the continuum model, however, the internal spatial structure of the wavefunctions also affects the Hartree energy. Nevertheless, the form factors $\lambda_{\mu}(\mathbf{k}, \mathbf{k} + \mathbf{G})$ decay quite quickly as a function of \mathbf{G} [21]—implying that the spatial variation of the density within the unit cell is not too large. As we discuss in the next two sections, the contribution from $R''(M; \mathbf{G})$ is essentially negligible.

5. Chiral model

We show that in the chiral model [64], the functions $R(M; \mathbf{G})$ are purely real, implying that the Hartree terms are all degenerate. The chiral model is a particular case of the continuum model in which hopping only occurs between A and B sites both within and between graphene layers. This constraint is implemented by setting w_0 in Eq. (A4) to zero. The result is an exact particle-hole (chiral) symmetry Γ that interchanges positive and negative energy states. We follow the discussion in the Appendix of Ref. 21. Γ may be assumed to act as

$$\Gamma : \quad c(\mathbf{k}) \rightarrow \eta^z c(\mathbf{k}). \quad (\text{D25})$$

In fact, in this basis, the sublattice index of the $c(\mathbf{k})$'s can be identified with the A and B sublattices of the two layers. It's then convenient to reinterpret the wavefunctions written in Eq. (A6), $u_{\mu, i; \mathbf{G}, \ell}(\mathbf{k})$. We explicitly identify the index $i = \text{A, B}$ with the sublattice, leaving ℓ to denote the layer. It then follows that the form factor may be written

$$\lambda_{\mu, ij}(\mathbf{k}, \mathbf{k}' + \mathbf{G}) = \left[\lambda_{\mu}^{(0)}(\mathbf{k}, \mathbf{k}' + \mathbf{G}) \mathbb{1}_{2 \times 2} + i \lambda_{\mu}^{(z)}(\mathbf{k}, \mathbf{k}' + \mathbf{G}) \eta^z \right]_{ij}, \quad (\text{D26})$$

where both $\lambda_{\mu}^{(0)}$ and $\lambda_{\mu}^{(z)}$ are real functions.

An additional symmetry allows one to rotate the two layers in opposite directions. The authors of Ref. 64 use this observation to simplify the problem substantially, resulting in an exact expression for the ground state wavefunction

at the magic angle. For any angle, however, it implies that the Hamiltonian of Eq. (A2) satisfies

$$\begin{pmatrix} iv_0 \boldsymbol{\eta}_{\theta/2} \cdot \boldsymbol{\nabla} & T(\mathbf{r}) \\ T^\dagger(\mathbf{r}) & iv_0 \boldsymbol{\eta}_{-\theta/2} \cdot \boldsymbol{\nabla} \end{pmatrix} = \eta^z \tau^z \begin{pmatrix} -iv_0 \boldsymbol{\eta}_{\theta/2} \cdot \boldsymbol{\nabla} & T(\mathbf{r}) \\ T^\dagger(\mathbf{r}) & -iv_0 \boldsymbol{\eta}_{-\theta/2} \cdot \boldsymbol{\nabla} \end{pmatrix} \tau^z \eta^z, \quad (\text{D27})$$

where Pauli operators η^z and τ^z act on the sublattice (A,B) and layer (t, b) indices respectively. The continuum representation of the wavefunction given in Eq. (A6) therefore satisfies

$$u_{\mu,i;\mathbf{G},\ell}(\mathbf{k}) = e^{i\varphi_{\mathbf{k}}} \sum_{i',\ell'} \eta_{ii'}^z \tau_{\ell\ell'}^z u_{\mu,i';-\mathbf{G}\ell'}(-\mathbf{k}), \quad (\text{D28})$$

which in turn implies

$$\lambda_\mu(\mathbf{k}, \mathbf{k}' + \mathbf{G}) = \lambda_\mu(-\mathbf{k}, -\mathbf{k}' - \mathbf{G}). \quad (\text{D29})$$

Similarly, the mean field Hamiltonian must give $\bar{h}_\mu(\mathbf{k}; M) = \bar{h}_\mu(-\mathbf{k}; M)$ and therefore

$$P_\mu(\mathbf{k}; M) = P_\mu(-\mathbf{k}; M). \quad (\text{D30})$$

These relations provide an additional constraint on the form of $\langle \rho(\mathbf{G}) \rangle_M$:

$$\begin{aligned} \langle \rho_\mu(\mathbf{G}) \rangle_M &= \int_{\mathbf{k}} \text{tr} [P_\mu(\mathbf{k}; M) \lambda_\mu(\mathbf{k}, \mathbf{k} + \mathbf{G})] = \int_{\mathbf{k}} \text{tr} [P_\mu(-\mathbf{k}; M) \lambda_\mu(-\mathbf{k}, -\mathbf{k} - \mathbf{G})] \\ &= \int_{\mathbf{k}} \text{tr} [P_\mu(\mathbf{k}; M) \lambda_\mu(\mathbf{k}, \mathbf{k} - \mathbf{G})] \\ &= \langle \rho_\mu(-\mathbf{G}) \rangle_M = \langle \rho_\mu(\mathbf{G}) \rangle_M^*. \end{aligned} \quad (\text{D31})$$

That is, $R(M; \mathbf{G})$ is *real*: $R''(M; \mathbf{G}) = 0$. From Eq. (D24), we conclude that the Hartree energies corresponding to all four insulating states are fully degenerate in the chiral limit:

$$H_H^{\text{QVH}} = H_H^{\text{QSVH}} = H_H^{\text{QH}} = H_H^{\text{QSH}}. \quad (\text{D32})$$

6. Numerical evaluation of Hartree term

We now return to the non-chiral version of the model. In Fig. 8(a) we plot the energy difference per electron of the Hartree term for the model using the parameters given in Appendix A as a function of the Dirac mass M . Even for a mass $M = 3 \text{ meV}$, the energy difference is as small as $2.5 \times 10^{-6} \text{ meV}$ — certainly our rough model is not expected to be reliable for such small energy differences.

We can understand the smallness in several ways. As mentioned at the end of Appendix D4c, the form factors $\lambda(\mathbf{k}, \mathbf{k} + \mathbf{G})$ decay quite quickly as a function of \mathbf{G} . We can further show that $R''(M; \mathbf{G}) = 0$ for all \mathbf{G} such that $\mathbf{G} = m_y[\mathbf{G}]$, $\mathbf{G} = c_3 m_y[\mathbf{G}]$, or $\mathbf{G} = c_3^2 m_y[\mathbf{G}]$. To do so, we start by using the fact that a basis exists in which

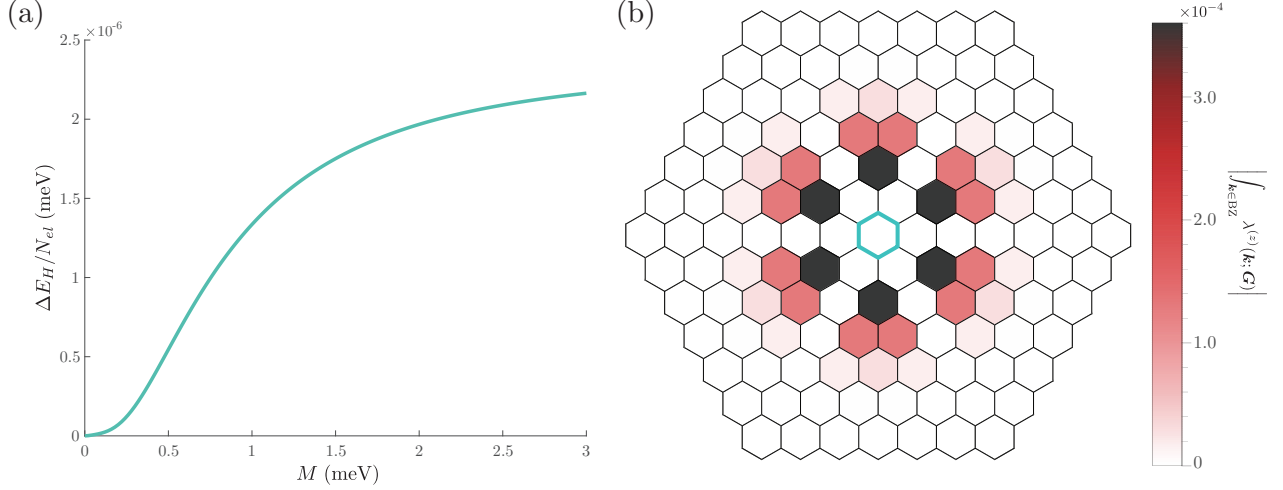


FIG. 8. (a) Energy difference as a function of the variational mass M between the QVH phase, H_H^{QVH} , and the other three phases, $H_H^{\text{other}} = H_H^{\text{QSH}} = H_H^{\text{QH}} = H_H^{\text{QSVH}}$, per electron at charge neutrality: $\Delta E_H/N_{el} = (H_H^{\text{QVH}} - H_H^{\text{other}})/N_{el}$. (b) Colour plot of $|\int_{\mathbf{k}} \lambda_{\mu}^{(z)}(\mathbf{k}; \mathbf{G})|$ as a function of the moiré reciprocal lattice vector \mathbf{G} . Each hexagon represents a different \mathbf{G} , with the central hexagon outlined in turquoise corresponding to $\mathbf{G} = 0$. Noticeably, $|\int_{\mathbf{k}} \lambda^{(z)}(\mathbf{k}; \mathbf{G})| = 0$ along all mirror axes, as we showed in the main text.

Eq. (D7) holds and the mirror symmetry acts as [46]

$$m_y : \quad c(\mathbf{k}) \rightarrow \eta^x c(m_y[\mathbf{k}]). \quad (\text{D33})$$

Since $h_{\mu}(\mathbf{k})$ satisfies the symmetry whereas the mass term $M\eta^z$ does not (e.g. $\bar{h}_{\mu}(\mathbf{k}; M) = \eta^x \bar{h}_{\mu}(m_y[\mathbf{k}]; -M)\eta^x$), we must have

$$P_{\mu}(\mathbf{k}; M) = \eta^x P_{\mu}(m_y[\mathbf{k}]; -M)\eta^x, \quad \lambda_{\mu}(\mathbf{k}, \mathbf{k} + \mathbf{G}) = \eta^x \lambda_{\mu}(m_y[\mathbf{k}], m_y[\mathbf{k} + \mathbf{G}])\eta^x. \quad (\text{D34})$$

We therefore find

$$\begin{aligned} \langle \rho_{\mu}(\mathbf{G}) \rangle_M &= \int_{\mathbf{k} \in \text{BZ}} \text{tr} [\eta^x P_{\mu}(m_y[\mathbf{k}]; -M)\eta^x \eta^x \lambda_{\mu}(m_y[\mathbf{k}], m_y[\mathbf{k} + \mathbf{G}])\eta^x] \\ &= \int_{\mathbf{k} \in \text{BZ}} \text{tr} [P_{\mu}(\mathbf{k}; -M)\lambda_{\mu}(\mathbf{k}, \mathbf{k} + m_y[\mathbf{G}])] \\ &= \langle \rho_{\mu}(m_y[\mathbf{G}]) \rangle_{-M} = \langle \rho_{\mu}(m_y[\mathbf{G}]) \rangle_M^*. \end{aligned} \quad (\text{D35})$$

It follows that $\langle \rho_{\mu}(\mathbf{G}) \rangle_M$ is *real* for all moiré reciprocal lattice vectors such that $\mathbf{G} = m_y[\mathbf{G}]: R''(M; \mathbf{G} = m_y[\mathbf{G}]) = 0$. The reflection axis chosen for m_y was actually arbitrary—by C_3 rotational symmetry, the same should hold for the two equivalent axes given by $C_3 m_y$ and $C_3^2 m_y$. Notably, this means that $R''(M; \mathbf{G}) = 0$ for the shortest set reciprocal lattice vectors.

We can quantify the size of $R''(M; \mathbf{G})$ for arbitrary \mathbf{G} through the follow set of observations. First, we note that the energies of the flat bands may be written as $E_{\mu, \pm}(\mathbf{k}) = h_{\mu, 0}(\mathbf{k}) \pm \epsilon_{\mu}(\mathbf{k})$, where $\epsilon_{\mu}^2(\mathbf{k}) = h_{\mu, x}^2(\mathbf{k}) + h_{\mu, y}^2(\mathbf{k})$.

This allows us to express the projection matrix as

$$P_\mu(\mathbf{k}; M) = \frac{1}{2} \left(\mathbb{1} - \frac{1}{\sqrt{\epsilon_\mu^2(\mathbf{k}) + M^2}} (h_{\mu,x}(\mathbf{k})\eta^x + h_{\mu,y}(\mathbf{k})\eta^y + M\eta^z) \right). \quad (\text{D36})$$

It then follows that

$$\begin{aligned} R''(M; \mathbf{G}) &= \frac{1}{2} (\langle \rho_\mu(\mathbf{G}) \rangle_M - \langle \rho_\mu(\mathbf{G}) \rangle_{-M}) = -\frac{1}{2} \int_{\mathbf{k} \in \text{BZ}} \frac{M}{\sqrt{\epsilon_\mu^2(\mathbf{k}) + M^2}} \text{tr} [\eta^z \lambda_\mu(\mathbf{k}, \mathbf{k} + \mathbf{G})] \\ &= - \int_{\mathbf{k} \in \text{BZ}} \frac{M}{\sqrt{\epsilon_\mu^2(\mathbf{k}) + M^2}} \lambda_\mu^{(z)}(\mathbf{k}; \mathbf{G}), \end{aligned} \quad (\text{D37})$$

where we've defined

$$\lambda_\mu^{(z)}(\mathbf{k}; \mathbf{G}) = -\frac{i}{2} \text{tr} [\eta^z \lambda_\mu(\mathbf{k}, \mathbf{k} + \mathbf{G})]. \quad (\text{D38})$$

We can verify through Eq. (D11b) and the identity $\lambda_\mu(\mathbf{k}, \mathbf{k} + \mathbf{G}) = \lambda_\mu^\dagger(\mathbf{k}, \mathbf{k} - \mathbf{G})$ that $\text{tr}[\eta^z \lambda_\mu(\mathbf{k}, \mathbf{k} + \mathbf{G})]$ must be imaginary. In limit that M is large, Eq. (D37) implies

$$R''(M; \mathbf{G}) \rightarrow - \int_{\mathbf{k} \in \text{BZ}} \lambda_\mu^{(z)}(\mathbf{k}; \mathbf{G}). \quad (\text{D39})$$

Assuming that $R''(M; \mathbf{G})$ is a monotonically increasing function of M (which Fig. 8(a) verifies at least for the parameters considered), we expect $\lambda_\mu^{(z)}$ to supply an upper bound on R'' :

$$|R''(M; \mathbf{G})| \leq \left| \int_{\mathbf{k} \in \text{BZ}} \lambda_\mu^{(z)}(\mathbf{k}; \mathbf{G}) \right|. \quad (\text{D40})$$

In Fig. 8(b) we plot the right hand side of the above equation as a function of \mathbf{G} . The fact that $\int_{\mathbf{k}} \lambda_\mu^{(z)}(\mathbf{k}; \mathbf{G})$ vanishes for all \mathbf{G} such that $\mathbf{G} = m_y[\mathbf{G}]$, $\mathbf{G} = c_3 m_y[\mathbf{G}]$, and $\mathbf{G} = c_3^2 m_y[\mathbf{G}]$ follows from the symmetry analysis given at the beginning of this section—as we see, the reciprocal lattice vectors with the smallest amplitudes do not contribute to $R''(M; \mathbf{G})$.

More importantly, the largest value of $\int_{\mathbf{k}} \lambda_\mu^{(z)}(\mathbf{k}, \mathbf{G})$ is already incredibly small—its maximum value is $\sim 3.6 \times 10^{-4}$. Even when multiplied by the relatively large interaction scale $V(a_M)$, the energy difference between the QVH and the other insulating phases remains small, as evinced by Fig. 8(a). We conclude that, at least within the approximation considered here, the QVH insulator is indistinguishable from its cousins, the QSVH, QH, and QSH states.

E. RANDOM FIELD ISING MODEL DOMAIN ESTIMATES

In this appendix, we discuss the Imry-Ma [69] arguments used in Sec. VA to obtain the estimates given in Eqs. (19) and (20) for the minimal domain size ξ_{dom} . We consider the regime where the homogeneous system

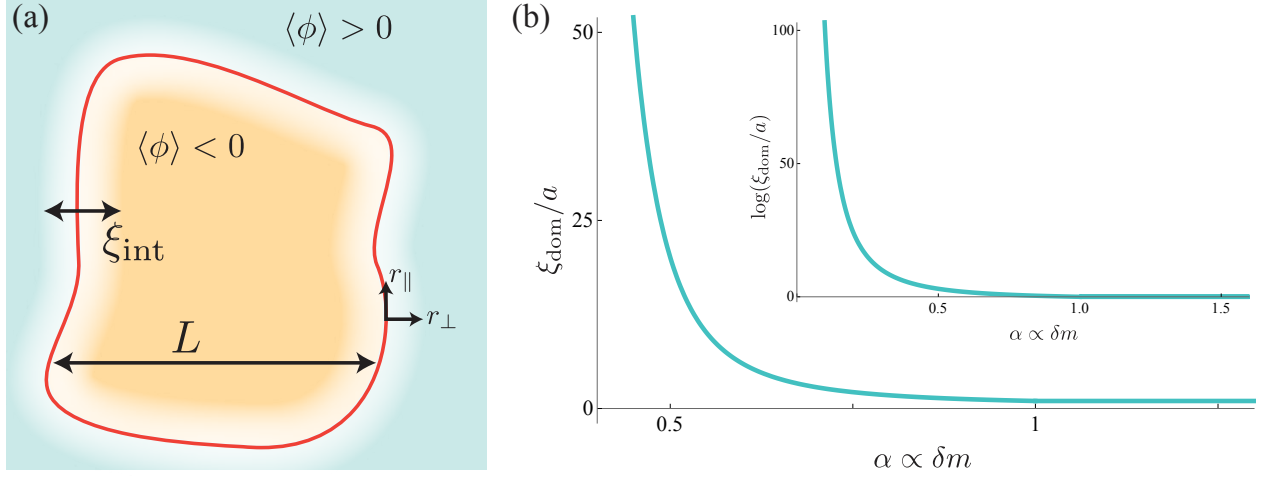


FIG. 9. (a) Illustration of a domain D of linear size $\sim L$ with $\langle \phi \rangle < 0$ (orange region) immersed within a region of $\langle \phi \rangle > 0$ (blue region). The boundary region of the domain, ∂D , is indicated in white. Its width, $\sim \xi_{\text{int}}$, is shown with an arrow. The coordinates $(r_{\perp}, r_{\parallel})$ used to estimate $\int_{\mathbf{r}} \mathcal{K} (\nabla \phi)^2$ are shown to the right of the domain. (b) Schematic plot of domain size, ξ_{dom} , as a function of α [Eq. (18)] for Gaussian-correlated disorder, Appendix E3. The inset plots the logarithm of the domain size. In both, $a = \max(\xi_{\text{dis}}, \xi_{\text{int}})$. When $\alpha \lesssim 1$, the disorder is effectively local and the domains are exponentially large, as per Eq. (E25). On the other hand, for $\alpha \gtrsim 1$, the domain size is set by the disorder correlation length ξ_{dis} . Coefficients of $\mathcal{O}(1)$ have been chosen by hand to smoothen the crossover between these two regimes. Since we assume that $\delta m \ll U$, $\alpha \gtrsim 1$ implies that $\xi_{\text{dis}} \gg \xi_{\text{int}}$.

would like to order—in this sense, we are assuming that disorder is weak compared to the interaction energy: $\delta m \ll U$. We next estimate the energy cost $E_{\text{dom}}(L)$ of changing the sign of ϕ within a domain D of linear extent $\sim L$, as depicted in Fig. 9(a). There are two contributions to E_{dom} : one from the interaction energy, $E_{\text{int}}(L)$, and another from the disorder potential, $E_{\text{dis}}(L)$. As reasoned in the main text, we assume that $|\langle \phi \rangle| \sim \mathcal{O}(1)$. Since we are primarily interested in the relative scaling of the two energy terms, coefficients of $\mathcal{O}(1)$ are not included.

The interaction energy of the domain is determined by the kinetic term of the Ising model:

$$E_{\text{int}}(L) \sim \int d^2 \mathbf{r} \mathcal{K} (\nabla \phi)^2. \quad (\text{E1})$$

The coefficient \mathcal{K} should have units of energy, and so we naturally set $\mathcal{K} \sim U$, as discussed in the main text. The Ising field ϕ changes only within the boundary region ∂D of the flipped domain D . Given our initial definition of ϕ [Eq. (14)], this change can only occur on the scale of ξ_{int} [Eq. (16)], implying that $(\nabla \phi)^2 \sim U \phi / \xi_{\text{int}}^2 \sim 1 / \xi_{\text{int}}^2$. Integrating over ∂D , including its width, contributes a factor of $\xi_{\text{int}} L$ so that the total cost is

$$E_{\text{int}}(L) = U \frac{L}{\xi_{\text{int}}}. \quad (\text{E2})$$

More concretely, this estimate can be obtained through the *ansatz* $\phi(\mathbf{r}) \sim \tanh(r_{\perp} / \xi_{\text{int}})$, where r_{\perp} is the direction perpendicular to the domain boundary, with the boundary itself occurring at $r_{\perp} = 0$ [see Fig. 9(a)]. Ignoring the effect of curvature, we again find

$$E_{\text{int}}(L) \sim U \int dr_{\parallel} \int dr_{\perp} \frac{1}{\xi_{\text{int}}^2} \text{sech}^4 \left(\frac{r_{\perp} - r_0}{\xi_{\text{int}}} \right) \sim U \cdot \frac{1}{\xi_{\text{int}}^2} \cdot L \cdot \frac{4}{3} \xi_{\text{int}} \sim U \frac{L}{\xi_{\text{int}}}. \quad (\text{E3})$$

We now consider the contribution to the energy cost of the domain due to the random field $\mathcal{B}(\mathbf{r})$ [as defined in and below Eq. (17)]. For a given realization of disorder, we have

$$E_{\text{dis}}(L) \sim \int_{\mathbf{r} \in D} \mathcal{B}(\mathbf{r}). \quad (\text{E4})$$

Depending on where the domain is placed, disorder can either increase or decrease the domain energy. For an arbitrarily chosen D , E_{dis} will average to zero, with a standard deviation given by

$$E_{\text{rms}}^2 \sim \overline{\left[\int_{\mathbf{r} \in D} \mathcal{B}(\mathbf{r}) \right]^2} = \frac{\delta m^2}{\xi_{\text{int}}^4} \int_{\mathbf{r}, \mathbf{r}' \in D} K \left(\frac{\mathbf{r} - \mathbf{r}'}{\xi_{\text{dis}}} \right). \quad (\text{E5})$$

Importantly, however, the location of the domain is not arbitrary. We can choose to place our domain in a region where this contribution is negative, taking the typical value

$$E_{\text{dis}} \sim -\sqrt{E_{\text{rms}}^2} \sim -\frac{\delta m}{\xi_{\text{int}}^2} \left[\int_{\mathbf{r}, \mathbf{r}' \in D} K \left(\frac{\mathbf{r} - \mathbf{r}'}{\xi_{\text{dis}}} \right) \right]^{1/2}. \quad (\text{E6})$$

The total cost of the domain is therefore

$$E_{\text{dom}}(L) \sim U \frac{L}{\xi_{\text{int}}} - \frac{\delta m}{\xi_{\text{int}}^2} \left[\int_{\mathbf{r}, \mathbf{r}' \in D} K \left(\frac{\mathbf{r} - \mathbf{r}'}{\xi_{\text{dis}}} \right) \right]^{1/2}. \quad (\text{E7})$$

If L_* exists such that $E_{\text{dis}}(L_*) = 0$, the formation of the domain is energetically favourable and long-range order is destroyed. This destruction occurs in all of the examples we consider.

1. Long-range disorder

The simplest example actually turns out to be the case of long-range disorder [104]:

$$K \left(\frac{\mathbf{r}}{\xi_{\text{dis}}} \right) = \frac{\xi_{\text{dis}}}{|\mathbf{r}|}. \quad (\text{E8})$$

We do not discuss this form of K in the main text since it is unlikely to describe the physical system; it nevertheless serves as a convenient example. We note that while ξ_{dis} is a lengthscale, it does not truly represent a correlation length in this context. Instead, it simply enters into the disorder strength as a multiplicative factor:

$$\overline{\mathcal{B}(\mathbf{r})\mathcal{B}(0)} = (\delta m^2 \xi_{\text{dis}}) \frac{1}{\xi_{\text{int}}^4} \frac{1}{|\mathbf{r}|}. \quad (\text{E9})$$

Inserting this definition into Eq. (E7), we find that the change in energy expected for a (judiciously-chosen) domain of size L is

$$E_{\text{dom}}(L) \sim \frac{L}{\xi_{\text{int}}} \left(U - \delta m \frac{\sqrt{\xi_{\text{dis}} L}}{\xi_{\text{int}}} \right). \quad (\text{E10})$$

For large L , it's clear that the domain energy eventually becomes negative, destabilizing the ordered phase. This destruction first occurs at the emergent length scale

$$L_* \sim \left(\frac{U \xi_{\text{int}}}{\delta m \xi_{\text{dis}}} \right)^2 \xi_{\text{dis}}. \quad (\text{E11})$$

We conclude that when the disorder is long-range, domains are expected to form once the system size is larger than L_* .

2. White noise (short-range) disorder

We now consider local, white noise disorder:

$$K \left(\frac{\mathbf{r}}{\xi_{\text{dis}}} \right) = \xi_{\text{dis}}^2 \delta^2(\mathbf{r}). \quad (\text{E12})$$

As in the long-range case, the parameter ξ_{dis} enters only as a multiplicative factor. Together with the disorder strength δm and the Fermi velocity v_F , they form a dimensionless parameter $\delta m \xi_{\text{dis}} / \hbar v_F$ discussed in Sec. III B.

Following the arguments above, an appropriately chosen domain therefore contributes an energy

$$E_{\text{dis}}(L) \sim -\delta m \frac{\xi_{\text{dis}} L}{\xi_{\text{int}}^2}. \quad (\text{E13})$$

The total energy cost of the domain is

$$E_{\text{dom}}(L) \sim U \frac{L}{\xi_{\text{int}}} - \delta m \frac{\xi_{\text{dis}} L}{\xi_{\text{int}}^2} = U \frac{L}{\xi_{\text{int}}} (1 - \alpha), \quad (\text{E14})$$

where we have defined

$$\alpha \equiv \frac{\delta m \xi_{\text{dis}}}{U \xi_{\text{int}}}, \quad (\text{E15})$$

as given in Eq. (18) of the main text. Notably, it is not $\delta m/U$ that controls the domain energy cost, but instead the ratio α . This feature is related to our remark that the true disorder strength is actually $g = \delta m \xi_{\text{dis}}$. The correct energy scale is therefore obtained in units of the UV cutoff, giving $g/\xi_{\text{int}} = \alpha U$, from which it follows that α is the appropriate tuning parameter, *not* $\delta m/U$. Equation (E14) simply tells us that when disorder is larger than the interaction scale, $\alpha \gtrsim 1$, there is no reason for the system to order. In this limit, the domain structure and fate of the theory is complicated and will not be relevant for us [105, 106].

Conversely, for $\alpha \lesssim 1$, Eq. (E14) may appear to imply that that the system should order. However, while Eq. (E15) is sufficient for large α , the analysis above omits the effect of domain roughening. This effect should be included in general, and it completely alters our conclusions when α is small.

Roughening in the context of the RFIM was first discussed in Ref. 72, and we now summarize the reasoning made there. We begin by considering a portion of a domain wall of linear extent y , displacing it by a (small) length w , and determining the change in energy, $\delta E(w, y)$. First, the displacement increases the length of the boundary

by $\delta E_{\text{int}} \sim U w / \xi_{\text{int}}$. With regards to the disorder field, we can choose to displace the boundary to either the left or the right direction, each of which has a 50% likelihood of decreasing the energy. There is therefore a 75% probability that the displacement lowers the energy by a typical amount $\delta E_{\text{dis}} \sim -\delta m \xi_{\text{dis}} \sqrt{w y} / \xi_{\text{int}}^2$. In total, the displacement results in a typical energy change

$$\delta E(w, y) \sim U \frac{w}{\xi_{\text{int}}} - \delta m \frac{\xi_{\text{dis}}}{\xi_{\text{int}}^2} \sqrt{w y}. \quad (\text{E16})$$

We now minimize δE with respect to w , to obtain

$$\begin{aligned} w_* &\sim \left(\frac{\delta m \xi_{\text{dis}}}{U \xi_{\text{int}}} \right)^2 y = \alpha^2 y, \\ \delta E_*(y) &\equiv \delta E(w_*, y) \sim -\alpha^2 U \frac{y}{\xi_{\text{int}}}. \end{aligned} \quad (\text{E17})$$

Next, we note that this procedure may be performed for segments of all sizes along the domain boundary. In particular, there are $N(y_\ell) = L/y_\ell$ segments of size $y_\ell = e^{-\ell} L$, each of which contributes an energy $\delta E_*(y_\ell)$. Summing over all scales returns the total energy contribution from domain wall roughening:

$$\begin{aligned} \delta E_{\text{tot}}(L, a) &= \int_0^{\log(L/a)} d\ell N(y_\ell) \delta E_*(y_\ell) \sim - \int_a^L \frac{dy}{y} \frac{L}{y} \alpha^2 U \frac{y}{\xi_{\text{int}}} \\ &\sim -\alpha^2 U \frac{L}{\xi_{\text{int}}} \log \left(\frac{L}{a} \right). \end{aligned} \quad (\text{E18})$$

Here, a is the smallest scale at which roughening may occur; in this context, $a \sim \xi_{\text{int}}$, though we will find otherwise in the next section. (Note that this ‘ a ’ should *not* be confused with the microscopic lattice constant of monolayer graphene.) Throughout this derivation, we have assumed that a is significantly smaller than L . Finally, the total energy cost of the domain is

$$E_{\text{dom}}(L) \sim \frac{L}{\xi_{\text{int}}} \left[U - \alpha^2 U \log \left(\frac{L}{\xi_{\text{int}}} \right) \right]. \quad (\text{E19})$$

Solving for $E_{\text{dom}}(L_*) = 0$, we find

$$L_* \sim \xi_{\text{int}} e^{c/\alpha^2}, \quad (\text{E20})$$

where we have introduced the non-universal constant $c \sim \mathcal{O}(1)$ to account for the imprecise nature of our scaling arguments. Once more, for systems larger than L_* , multiple domains should be apparent.

As we mentioned below Eq. (E18), our integration was predicated on the assumption that the domain size L was much larger than ξ_{int} . It is clear that this is only satisfied provided the disorder is weak: $\alpha \ll 1$. When the disorder is stronger, the situation is more complicated.

3. Gaussian-correlated disorder

We now consider the situation considered in the main text, that of Gaussian correlated disorder:

$$K\left(\frac{\mathbf{r}}{\xi_{\text{dis}}}\right) = e^{-\frac{r^2}{2\xi_{\text{dis}}^2}}. \quad (\text{E21})$$

Unlike the previous two cases, the scale ξ_{dis} is a true correlation length in this scenario, as is clear from the form of the disorder-induced energy reduction:

$$E_{\text{dis}}(L) \sim -\delta m \frac{\xi_{\text{dis}} L}{\xi_{\text{int}}^2} \sqrt{1 - e^{-L^2/2\xi_{\text{dis}}^2}}. \quad (\text{E22})$$

While the domain size appeared as a ratio of the UV cutoff $\ell_{\text{UV}} = \xi_{\text{int}}$ in the previous two examples, here $E_{\text{dis}}(L)$ is also a function of L/ξ_{dis} .

There are two natural limits to consider. In the first, we take the domain size to be small enough relative to ξ_{dis} that the smoothness of the disorder is still important, *i.e.* we cannot simply ignore the exponential in Eq. (E22). As an extreme example, when $L \ll \xi_{\text{dis}}$,

$$E_{\text{dis}}(L) \sim -\delta m \frac{L^2}{\xi_{\text{int}}^2}. \quad (\text{E23})$$

That is, the change in energy is proportional to the *volume* of the domain. This observation makes sense given that $\mathcal{B}(\mathbf{r})$ should be essentially constant for two points within a distance ξ_{dis} of one another. In fact, it seems clear that an energetically favourable domain should be at least ξ_{dis} in extent: $L_* \gtrsim \xi_{\text{dis}}$. We therefore examine the threshold scenario given by $L = \xi_{\text{dis}}$. We conclude that domain formation is favourable when

$$E_{\text{dom}}(\xi_{\text{dis}}) \sim U \frac{\xi_{\text{dis}}}{\xi_{\text{int}}} \left(1 - \frac{\delta m \xi_{\text{dis}}}{U \xi_{\text{int}}}\right) = U \frac{\xi_{\text{dis}}}{\xi_{\text{int}}} (1 - \alpha) \lesssim 1. \quad (\text{E24})$$

The parameter α that appeared in the white noise case, Eq. (E15), has showed up again. When it is greater than unity, $\alpha \gtrsim 1$, the disorder destroys long-range order, resulting in domains of typical size $\xi_{\text{dom}} \sim \xi_{\text{dis}}$.

When $\alpha \lesssim 1$, the interaction energy cost associated with the boundary of a domain of linear extent ξ_{dis} is greater than the gain associated with aligning with the random field. For domains larger than ξ_{dis} , the random field within the domain is only weakly correlated. The exponential under the square root may therefore be neglected, resulting in an expression identical to our original estimate for the domain energy with white noise disorder in Eq. (E14). As we discussed there, this expression was not complete: the roughening of the domain walls must also be taken into account, resulting in the contribution given in Eq. (E18). The arguments made in Sec. E2 follow through for weak, Gaussian-correlated disorder in all respects save for one minor caveat. Unlike the white noise disorder case, the roughening cutoff for Gaussian-correlated disorder is not necessarily ξ_{int} . Instead, only scales down to *at most* ξ_{dis} should be included, since this is where our omission of the exponential ceases to be valid, *i.e.* $a = \max(\xi_{\text{int}}, \xi_{\text{dis}})$. Setting the domain energy to zero, we find

$$\xi_{\text{dom}} \lesssim \max(\xi_{\text{int}}, \xi_{\text{dis}}) e^{c/\alpha^2}, \quad (\text{E25})$$

where $c \sim \mathcal{O}(1)$ is again a non-universal constant. In Fig. 9(b), we show ξ_{dom} for Gaussian-correlated disorder for both regimes, $\alpha \lesssim 1$ and $\alpha \gtrsim 1$.

F. COMPETING ORDERS

We now address the possibility considered in Sec. VC that the QVH state is not the ground state of the clean theory at charge neutrality—either one of the other three $C_2\mathcal{T}$ -breaking insulators (QSH, QH, or QSVH) or a completely different order minimizes the energy of the homogeneous theory.

We are interested in studying the conditions under which the QVH phase is realized. To simplify the analysis, we assume that there is a single competing phase whose order parameter does not couple to disorder, but whose ground state energy density, $\mathcal{E}_{\text{comp}}$, is lower than the energy density of the QVH phase, \mathcal{E}_{QVH} , by a small amount. We measure this distinction in terms of the energy difference $\delta\epsilon$ within a region of area $\ell_{\text{UV}}^2 = \xi_{\text{int}}^2$:

$$\frac{\delta\epsilon}{\xi_{\text{int}}^2} = \mathcal{E}_{\text{QVH}} - \mathcal{E}_{\text{C}} \geq 0 \quad (\text{F1})$$

Throughout this section, we assume that $\delta\epsilon \ll U$. While this ground state energy difference implies that the competing phase is realized in a perfectly clean sample, disorder exclusively favours the local realization of the QVH phase. We therefore expect the majority of the sample to be in the QVH phase when $\delta\epsilon$ is sufficiently small. Using the Ising notation of Sec. VA and Appendix E, we quantify this expectation as

$$\left[\frac{1}{\text{vol}} \int_{\mathbf{r}} \langle \phi^2(\mathbf{r}) \rangle \right]^{1/2} \gtrsim \frac{1}{2}, \quad (\text{F2})$$

where ‘vol’ denotes the sample volume.

We approach the problem in two complementary fashions. The question of an Ising order parameter competing with another phase may bring to mind dilute Ising physics, where here ‘vacancies’ represent regions where the Ising ϕ field is not ordered. In Appendix F1, we describe a mean field solution of a classical $2d$ lattice model formulated to tackle this type of question.

While useful, because of the low-dimensionality of the problem, mean field theory is not particularly reliable. In particular, we are free to take the limit $\delta\epsilon \rightarrow -\infty$, effectively removing the ‘competing’ phase from the problem. In this limit, our results should agree with those of Sec. VA and Appendix E. There, we found that any disorder was sufficient to destroy long-range order. In contrast, the mean field calculation falsely finds long-range order in this limit. We therefore devise an Ising formulation of the problem in Appendix F2, which allows us to make Imry-Ma arguments similar to those of Appendix E.

1. Blume-Capel description

In keeping with the Ising description of the QVH insulator, we view the ordering of the competing phase as the presence of an annealed ‘vacancy.’ At finite temperature, this physics is known to give rise to the tricritical Ising fixed point, though this observation is not relevant for our discussion. While continuum descriptions do exist,

for our purposes, it is most convenient to employ a lattice model. We therefore consider the Blume-Capel model [74, 75] on an (unspecified) lattice of coordination number z with quenched random-field disorder:

$$H_{\text{BC}} = -\frac{J}{z} \sum_{\langle r, r' \rangle} s_r s_{r'} + \mu \sum_r s_r^2 + \sum_r h_r s_r, \quad (\text{F3})$$

where the classical spins may take three values: $s_r \in \{+1, -1, 0\}$. As above, the quenched disorder is represented through a random ‘magnetic field’ h_r . For simplicity, we assume that h_r satisfies Gaussian white noise disorder. The corresponding probability distribution reads

$$\mathcal{P}(h_r) = \frac{e^{-\frac{h_r^2}{2h_0^2}}}{\sqrt{2\pi h_0^2}}. \quad (\text{F4})$$

The use of this distribution is equivalent to our previous definitions of the disorder distribution, entirely in terms of moments:

$$\overline{h_r} = 0, \quad \overline{h_r h_{r'}} = h_0^2 \delta_{r, r'}. \quad (\text{F5})$$

We associate $s_r = \pm 1$ with the realization of the QVH phase, *i.e.* $\langle \phi \rangle \sim \pm 1$, and ‘vacancies’ $s_r = 0$ with competing phase. The exchange energy J corresponds to the Coulomb interaction strength, $J \sim U$, while the random field strength h_0 , should be mapped to the disorder strength $h_0 \sim \delta m \xi_{\text{dis}} / \xi_{\text{int}}$ in units of the UV cutoff ξ_{int} [see the discussion below Eq. (E15)]. Finally, the so-called ‘crystal field,’ μ , can be related to the energy splitting $\delta\epsilon$ by establishing when the competing phase (all $s_r = 0$) and QVH phase (all $s_r = +1$ or -1) are degenerate, indicating that $\mu = \delta\epsilon + J/2 \sim \delta\epsilon + U/2$.

As discussed, we analyze this model in mean field theory [76, 77]. Letting $m \equiv \langle s_r \rangle$ be the average magnetization, the mean-field free energy is

$$\begin{aligned} f_{\text{BC}}(m) &= \frac{1}{2} J m^2 - \overline{\log(1 + e^{-\beta\mu} 2 \cosh[\beta(Jm + h)])} \\ &= \frac{1}{2} J m^2 - \int \frac{dh}{\sqrt{2\pi} h_0} e^{-h^2/2h_0^2} \log(1 + e^{-\beta\mu} 2 \cosh[\beta(Jm + h)]), \end{aligned} \quad (\text{F6})$$

where β is the inverse temperature and we explicitly average over the Gaussian distribution of Eq. (F4) in the second line. Taking the zero temperature limit, $\beta \rightarrow \infty$, the integral can be evaluated exactly, giving

$$\begin{aligned} f_{\text{BC}}(m) &= \frac{1}{2} J m^2 + \frac{1}{2} \left[(\mu - Jm) \text{Erfc}\left(\frac{\mu - Jm}{\sqrt{2}h_0}\right) + (\mu + Jm) \text{Erfc}\left(\frac{\mu + Jm}{\sqrt{2}h_0}\right) \right] \\ &\quad - \frac{h_0}{\sqrt{2\pi}} \left(e^{-(\mu - Jm)^2/2h_0^2} + e^{-(\mu + Jm)^2/2h_0^2} \right), \end{aligned} \quad (\text{F7})$$

where $\text{Erfc}(x)$ is the complementary error function. The magnetization is determined by extremizing f_{BC} , result-

ing in the self-consistency equation

$$m = \frac{1}{2} \left[\text{Erfc} \left(\frac{\mu - Jm}{\sqrt{2}h_0} \right) - \text{Erfc} \left(\frac{\mu + Jm}{\sqrt{2}h_0} \right) \right]. \quad (\text{F8})$$

The expectation value of the spin squared, $q \equiv \sqrt{\langle s_r^2 \rangle}$, is directly analogous to the expression on the right-hand side of Eq. (F2), *i.e.*, when $q \gtrsim 1/2$, QVH order prevails. It is calculated by taking the derivative of f_{BC} with respect to μ :

$$q^2 = \frac{\partial}{\partial \mu} f_{\text{BC}} = \frac{1}{2} \left[\text{Erfc} \left(\frac{\mu - Jm}{\sqrt{2}h_0} \right) + \text{Erfc} \left(\frac{\mu + Jm}{\sqrt{2}h_0} \right) \right]. \quad (\text{F9})$$

In Figs. 10(a) and (c), we plot m and q as functions of $\delta\epsilon/U$ and α , respectively. To make contact with the phase diagram in the main text, Fig. 5, we also plot m and q with the y -axis given by $\gamma \delta\epsilon/\delta m$, where $\gamma = \xi_{\text{int}}/\xi_{\text{dis}}$, in Figs. 10(b) and (d).

Figures 10(a) and (b) indicate that m orders for $\delta\epsilon \lesssim 0$ when disorder is sufficiently small. While these calculations agree with our expectations when $\delta m = 0$, we showed in Appendix E that any nonzero disorder destroys long-range order. The presence of regions with $m \neq 0$ is therefore an artifact of the mean field theory; given the low dimension, the failure of mean field theory in this regard is not surprising. Nevertheless, we take it as a good sign that m approaches zero close to $\alpha \sim 0.8 \sim 1$ for $\delta\epsilon < 0$ since this condition defines the crossover regime identified in Appendix E. We therefore optimistically associate mean field ordered regions with those that in reality possess exponentially large domains.

The density plots in Figs. 10(c) and (d) display q . Obviously, when our mean field prescription indicates that m is ordered, q is non-zero as well, as a quick comparison with (a) and (b) clearly shows. Outside of these regions, however, we find that q only vanishes exactly when $\delta m \rightarrow 0$ (equivalently, $h_0 \rightarrow 0$) as well. From Eq. (F9), we verify that when $m = 0$,

$$q(m=0) = \sqrt{\text{Erfc} \left(\frac{\mu}{\sqrt{2}h_0} \right)}, \quad (\text{F10})$$

implying that contours of constant q are represented by straight lines extending from the $\mu = 0$ origin (not to be confused with $\delta\epsilon = 0$ origin), as shown in Fig. 10(c). More precisely, we can numerically solve for the line along which $q = 1/2$:

$$\frac{1}{2} = \sqrt{\text{Erfc} \left(\frac{\eta_{1/2}}{\sqrt{2}} \right)}, \quad (\text{F11})$$

to obtain $\eta_{1/2} \cong 1.15$. Then, provided $\delta\epsilon/U$ and α are such that $m = 0$, we find that $q = 1/2$ along the line

$$\frac{\delta\epsilon}{U} = \eta_{1/2}\alpha - \frac{1}{2}. \quad (\text{F12})$$

We plot this contour with a pink dashed line in Fig. 10(c). It follows that the system is primarily in the QVH phase

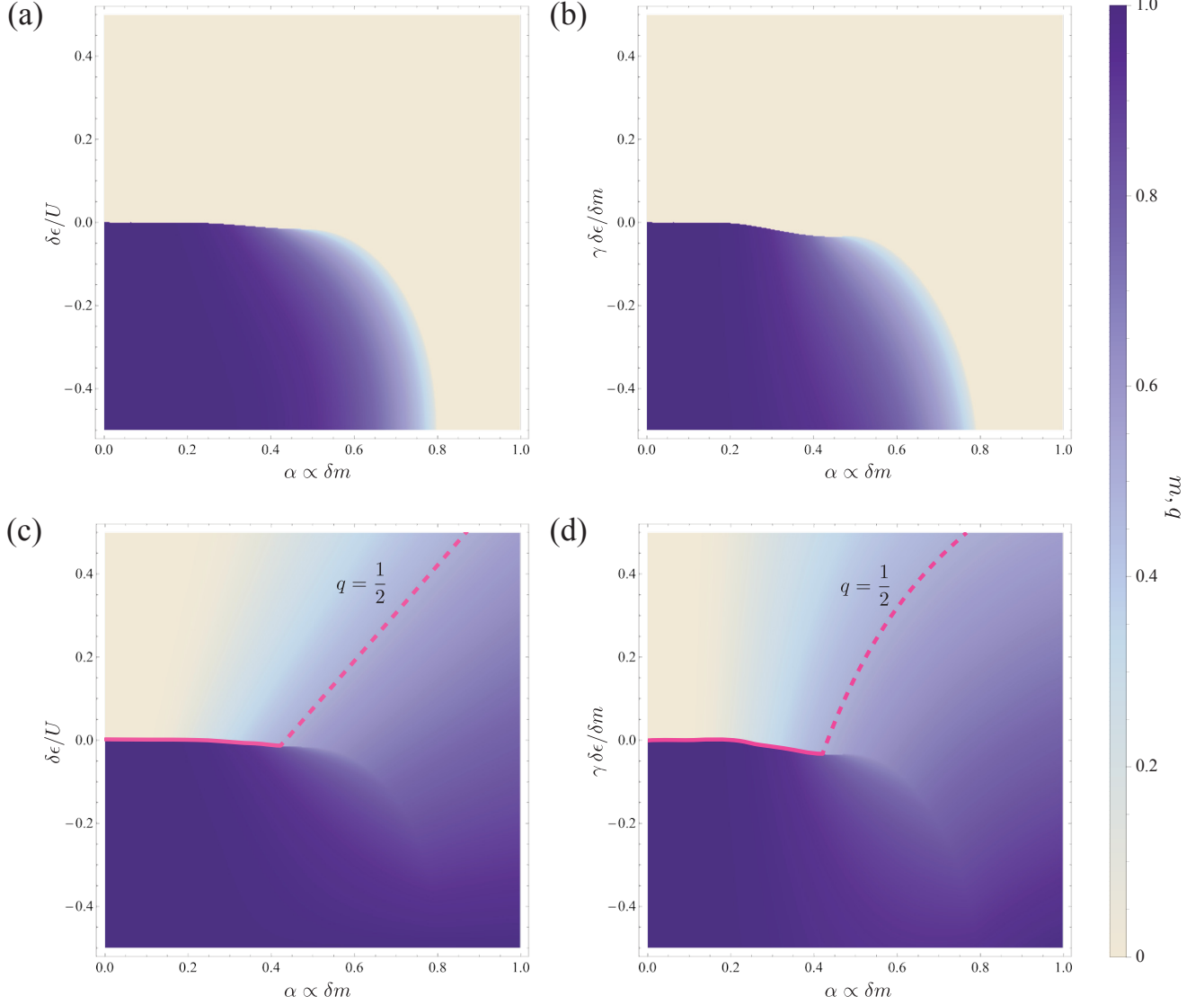


FIG. 10. (a), (b) Density plot of the (absolute value of the) magnetization, obtained by minimizing $f_{BC}(m)$ in Eq. (F7). (c), (d) Density plots of q , as given in Eq. (F9). The colour scheme for all plots, (a)-(d), is shown on the right, and, in (b) and (d), $\gamma = \xi_{\text{int}}/\xi_{\text{dis}}$. The solid pink line in (c) and (d) indicates the first order phase transition between regions with q small and regions with $q \sim 1$ (as follows from having $m \sim \pm 1$ in that region). The dashed pink line, on the other hand, is the contour along which $q = 1/2$ and $m = 0$; we view it as demarcating a crossover between regions where the competing phase percolates and regions where the QVH insulator percolates. It follows that for both (c) and (d), the network scenario we propose should be valid in the regions below and to the right of the pink lines.

when either

$$\delta\epsilon \lesssim U \left(\eta_{1/2} \alpha - \frac{1}{2} \right) \quad \text{or} \quad \delta\epsilon \lesssim 0. \quad (\text{F13})$$

The modification needed to obtain the crossover lines drawn in (d) is straightforward:

$$\delta\epsilon \lesssim \delta m \frac{\xi_{\text{dis}}}{\xi_{\text{int}}} \left(\eta_{1/2} - \frac{1}{2\alpha} \right) \quad \text{or} \quad \delta\epsilon \lesssim 0. \quad (\text{F14})$$

2. Competing Ising field description

The mean field theory discussed above had the advantage of simplicity, but did not correctly capture the absence of long-range order. We therefore employ an Imry-Ma description, similar to the analysis of Appendix E. The ordering of both phases is now modelled by two distinct Ising fields. As above, we associate ϕ with the QVH insulator (*i.e.*, C_2 symmetry breaking) and Φ with the competing phase. The total energy is given by $H_{\text{Ising}} + H'_{\text{Ising}} + H_{\phi\Phi} + H_{\phi,\text{dis}}$ where

$$\begin{aligned} H_{\text{Ising}} &= \int d^2\mathbf{r} \left[\mathcal{K} (\nabla\phi)^2 - \frac{|r|}{2}\phi^2 + \frac{u}{4!}\phi^4 \right], \\ H'_{\text{Ising}} &= \int d^2\mathbf{r} \left[\mathcal{K}' (\nabla\Phi)^2 - \frac{|r'|}{2}\Phi^2 + \frac{u'}{4!}\Phi^4 \right], \\ H_{\phi\Phi} &= \int d^2\mathbf{r} \lambda \phi^2 \Phi^2, \\ H_{\text{dis}} &= \int d^2\mathbf{r} \mathcal{B}(\mathbf{r})\phi(\mathbf{r}). \end{aligned} \tag{F15}$$

Since both ϕ and Φ are dimensionless, $\mathcal{K}, \mathcal{K}'$ have dimensions of energy. We assume that the interaction scales of the QVH and competing phases are similar, prompting us to set both to $\sim U$. Similarly, the remaining parameters describing H_{Ising} and H'_{Ising} , $r, r', u,$ and u' , have units of energy over length squared. Their natural scale is therefore U/ℓ_{UV}^2 where ℓ_{UV} is the UV cutoff, which should in turn be approximately given by $\xi_{\text{int}} = \hbar v_F / \Delta_{\text{CNP}}$, as discussed in Sec. V A. However, this assignment of energy scales cannot be the entire story since the difference in ground state energies, Eq. (F1), has not yet been included. Because $\delta\epsilon$ is assumed to be much smaller than U , and we ignore coefficients of $\mathcal{O}(1)$, the exact implementation is unimportant. Nevertheless, to be concrete, we note that if one wishes to ensure that Eq. (F1) holds while also requiring the magnitudes of ϕ and Φ to be identical in their respective ordered phases, the following choice is sufficient:

$$|r'| \sim |r| + \frac{2|r|}{3u} \frac{\delta\epsilon}{\xi_{\text{int}}^2}, \quad u' \sim u + \frac{2}{3} \frac{\delta\epsilon}{\xi_{\text{int}}^2}. \tag{F16}$$

The parameter λ in $H_{\phi\Phi}$ is assumed to be larger than the other scales of the theory in order to guarantee that $\langle\phi\rangle \neq 0$ and $\langle\Phi\rangle \neq 0$ do not occur within the same region. Finally, the last term, H_{dis} , describes the behaviour of disorder. We will consider both white noise and Gaussian-correlated, as defined in Eqs. (E12) and (E21) respectively.

We examine this system in several steps. Using Imry-Ma type arguments similar to those of Appendix E, we begin by studying the formation of a ϕ -ordered domain within a uniformly Φ -ordered system for both white noise and Gaussian-correlated disorder. As we did in Appendix E, coefficients of $\mathcal{O}(1)$ are ignored. Next, we argue that if the physical parameters favour the formation of a single ϕ -ordered domain, a macroscopically large fraction of the system should also ϕ -order. Our final result is a function of the ratio α [see Eq. (E15)], $\delta\epsilon_c(\alpha)$, that parametrizes a crossover between the two regimes of interest: when $\delta\epsilon \lesssim \delta\epsilon_c(\alpha)$, the system is primarily ϕ -ordered, whereas when $\delta\epsilon \gtrsim \delta\epsilon_c(\alpha)$, the system is primarily Φ -ordered. Figure 5 shows the resulting phase diagram.

a. *Single ϕ -domain formation: white noise disorder*

To make contact with the mean field theory of Appendix F 1, we begin by considering white noise disorder. We assume that the competing phase is realized, $\langle \Phi \rangle \neq 0$, and examine the energy cost associated with the formation of a ϕ -ordered domain. As in Appendix E, there are energy contributions from interactions along the domain boundary and from the random field $\mathcal{B}(\mathbf{r})$. Since we assume that $\mathcal{K} \sim \mathcal{K}' \sim U$, the interaction energy cost E_{int} is identical to the expression given in Eq. (E2)⁶. Similarly, the contribution from disorder, E_{dis} , follows from the expression in Eq. (E6), giving the same result as in Eq. (E13). Unlike Appendix E, there is an important additional cost associated with the difference in ground state energy. On general grounds, the cost must increase with the domain *area*:

$$E_{\text{comp}}(L) \sim \delta\epsilon \frac{L^2}{\xi_{\text{int}}^2}. \quad (\text{F17})$$

We could also have obtained this result from the Hamiltonian defined in Eq. (F15) with the coefficients defined in Eq. (F16). The total energy cost of a ϕ -ordered domain is given by the sum of this expression with E_{int} and E_{dis} :

$$E_{\phi\text{-dom}}(L) \sim \delta\epsilon \frac{L^2}{\xi_{\text{int}}^2} + U \frac{L}{\xi_{\text{int}}} - \delta m \frac{\xi_{\text{dis}} L}{\xi_{\text{int}}^2} = \delta m \frac{\xi_{\text{dis}} L}{\xi_{\text{int}}^2} \left(\frac{\delta\epsilon}{\delta m} \frac{L}{\xi_{\text{dis}}} + \frac{1}{\alpha} - 1 \right). \quad (\text{F18})$$

This result is the analogue of Eq. (E14). There, we concluded that when $\alpha \gtrsim 1$, disorder was ‘‘large’’ and the system would not order. While this expression also indicates that $\alpha \gtrsim 1$ is necessary to destroy the local order (here, Φ -order instead a different type of ϕ -order), the energy cost of the ϕ -domain is also dependent on its size, L : the smaller the domain size, the more favourable it is. A threshold value of $\delta\epsilon$ can therefore be defined by the condition $E_{\phi\text{-dom}}(a) < 0$, where a is the smallest possible domain size. (Again, ‘ a ’ should not be confused with the microscopic lattice constant of monolayer graphene here or below.) For the current situation, clearly $a \sim \xi_{\text{int}}$; nevertheless, with an eye to the subsequent section, it is convenient to leave a unspecified. That is, $E_{\phi\text{-dom}}(a) < 0$ provided

$$\delta\epsilon \lesssim \delta\epsilon_c(\alpha), \quad \delta\epsilon_c(\alpha) \equiv \delta m \frac{\xi_{\text{dis}}}{a} \left(1 - \frac{1}{\alpha} \right), \quad \text{when } \alpha \gtrsim 1. \quad (\text{F19})$$

Here, we have defined the ‘critical’ energy difference $\delta\epsilon_c(\alpha)$ in the region where $\alpha \gtrsim 1$ for white noise disorder with a minimal domain size $a = \xi_{\text{int}}$. We generalize this definition to smaller values of α below.

We note that up to coefficients of $\mathcal{O}(1)$, this inequality has the same dependence on α as our mean field result in Eq. (F14)! At least in the simple regime, the Blume-Capel and Imry-Ma descriptions are in agreement.

As we saw in Appendix E 2, once $\alpha \lesssim 1$, the effects domain wall roughening become important and must be included. Because roughening does not change the domain area significantly, the roughening contribution Eq. (E18) remains valid⁷. We note that this situation is similar to what occurs in the absence of a competing order

⁶ One might argue that it is more honest to define $\mathcal{K}' \sim \mathcal{K} + \delta\epsilon \sim U + \delta\epsilon$ in analogy with the definitions of Eq. (F16). However, since $\delta\epsilon \ll U$ by assumption, this difference is negligible.

⁷ Alternatively, we can argue that since the displacement is equally likely to increase or decrease the domain area, Eq. (E16) remains valid on average

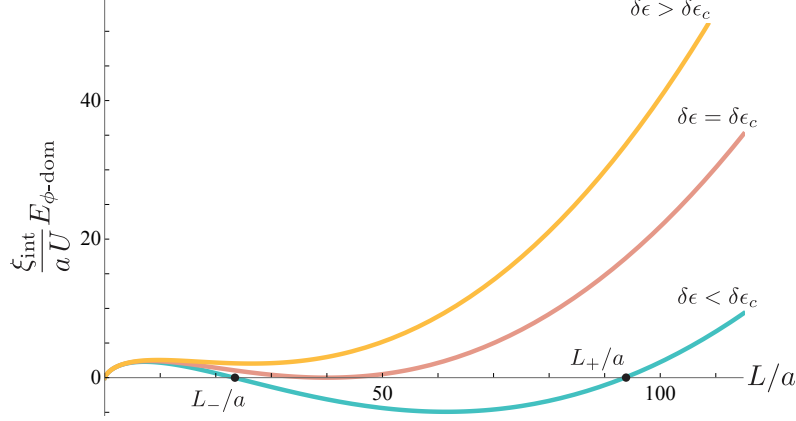


FIG. 11. Plot of the energy cost associated with adding a ϕ -ordered domain to a uniformly Φ -ordered system when $\delta\epsilon > \delta\epsilon_c$ (orange), $\delta\epsilon = \delta\epsilon_c$ (pink), and $\delta\epsilon < \delta\epsilon_c$ (turquoise). For $\delta\epsilon < \delta\epsilon_c$, we see that domain formation is energetically favourable, $E_{\phi\text{-dom}} < 0$, for domains with linear extent L satisfying $L_- < L < L_+$. Here, we have set $\alpha \sim 0.6$, for which $\delta\epsilon_c \sim 0.015 \delta m \xi_{\text{dis}}/a$.

when a small, uniform magnetic field is applied [72, 106]. The resulting cost of a ϕ domain is

$$\begin{aligned} E_{\phi\text{-dom}}(L) &\sim \delta\epsilon \left(\frac{L}{\xi_{\text{int}}} \right)^2 + U \frac{L}{\xi_{\text{int}}} - U \left(\frac{\delta m \xi_{\text{dis}}}{U \xi_{\text{int}}} \right)^2 \frac{L}{\xi_{\text{int}}} \log \left(\frac{L}{a} \right) \\ &= U \frac{L}{\xi_{\text{int}}} \left[\frac{\delta\epsilon}{U} \frac{L}{\xi_{\text{int}}} + 1 - \alpha^2 \log \left(\frac{L}{a} \right) \right]. \end{aligned} \quad (\text{F20})$$

Again, a is the minimal domain size, which is equivalent to ξ_{int} in this case. We can now define a critical energy difference in the small α regime. We find that there exists a solution $E_{\text{dom}}(L) = 0$ provided $\delta\epsilon$ satisfies

$$\delta\epsilon \lesssim \delta\epsilon_c(\alpha), \quad \delta\epsilon_c(\alpha) \equiv \frac{\xi_{\text{dis}}}{a} \delta m \alpha e^{-c(\frac{1}{\alpha^2}+1)}, \quad \text{when } \alpha \lesssim 1. \quad (\text{F21})$$

In Fig. 11, we plot $E_{\phi\text{-dom}}(L)$ as a function L for several values of $\delta\epsilon$. As indicated in the figure, when $\delta\epsilon < \delta\epsilon_c$, there is an entire region where $E_{\phi\text{-dom}} < 0$ for $L_- < L < L_+$. Naturally, as $\delta\epsilon \rightarrow 0$, $L_- \rightarrow L_*$ [as defined in Eq. (E20)] while $L_+ \rightarrow \infty$.

b. Single ϕ -domain formation: Gaussian correlated disorder

We now repeat the exercise above for Gaussian-correlated disorder. The energy cost of inserting a ϕ -ordered domain into a uniformly Φ -ordered system is on average

$$E_{\phi\text{-dom}}(L) \sim \delta\epsilon \left(\frac{L}{\xi_{\text{int}}} \right)^2 + U \frac{L}{\xi_{\text{int}}} - \delta m \frac{\xi_{\text{dis}} L}{\xi_{\text{int}}^2} \sqrt{1 - e^{-L^2/2\xi_{\text{dis}}^2}}. \quad (\text{F22})$$

We first study the regime where the smoothness of the disorder is important, *i.e.* the exponential under the square root is important. In this case, we expect the ϕ -domains to track the disorder potential and therefore be of the same

size as the disorder correlation length ξ_{dis} . In order for this to be energetically favourable, we must have

$$0 > E_{\phi\text{-dom}}(\xi_{\text{dis}}) \sim \delta\epsilon \frac{\xi_{\text{dis}}^2}{\xi_{\text{int}}^2} + U \frac{\xi_{\text{dis}}}{\xi_{\text{int}}} - \delta m \frac{\xi_{\text{dis}}^2}{\xi_{\text{int}}^2} = U \left(\frac{\xi_{\text{dis}}}{\xi_{\text{int}}} \right)^2 \left[\frac{\delta\epsilon}{\delta m} + \frac{1}{\alpha} - 1 \right], \quad (\text{F23})$$

It follows that ϕ -ordered domains of linear extent ξ_{dis} should form once

$$\delta\epsilon \lesssim \delta\epsilon_c(\alpha), \quad \delta\epsilon_c(\alpha) \equiv \delta m \left(1 - \frac{1}{\alpha} \right), \quad \text{when } \alpha \gtrsim 1. \quad (\text{F24})$$

This critical energy difference is nearly identical to the analogous expression obtained for white noise disorder in Eqs. (F14) and (F19). The most notable difference between the two inequalities is the prefactor $\xi_{\text{dis}}/\xi_{\text{int}}$ multiplying the right-hand side. Going back to the previous section, we see that this coefficient originates from setting the minimal domain size to ξ_{int} . In contrast, for Gaussian-correlated disorder, the smallest allowed domains are expected to be ξ_{dis} , and so $\delta\epsilon_c(\alpha)$ contains no such prefactor.

As we saw in Appendix E, once $\alpha \lesssim 1$, Gaussian-correlated disorder can be treated as local white-noise disorder, which necessitates a treatment that includes the effects of domain wall roughening. The relevant expression for $E_{\phi\text{-dom}}(L)$ is therefore identical to the one given in Eq. (F20), save that the smallest domain size is given by $a = \max(\xi_{\text{int}}, \xi_{\text{dis}})$. The inequality describing the favourability of domain formation is now

$$\delta\epsilon \lesssim \delta\epsilon_c(\alpha), \quad \delta\epsilon_c(\alpha) \equiv \frac{\xi_{\text{dis}}}{a} \delta m \alpha e^{-c(\frac{1}{\alpha^2}+1)}, \quad \text{when } \alpha \lesssim 1. \quad (\text{F25})$$

c. Multiple ϕ -domains

The formation of a single domain does not necessarily imply the network model we propose as a description for mTBG at charge neutrality. Instead, we want the ϕ -ordered regions to percolate throughout the sample, as implied by the condition given in Eq. (F2). We argue that ϕ -order should start dominating at a crossover set by the scale $\delta\epsilon_c(\alpha)$. As discussed in Appendix F2 a, within our approximation, domain boundaries between different ϕ orientations have the same cost as domains between Φ - and ϕ -ordered regions. As a result, we can imagine ‘tiling’ the ϕ -ordered regions into domains of some size ξ_* . For instance, when $\alpha \lesssim 1$, the energy difference between a uniformly Φ -ordered system and a (non-uniformly) ϕ -ordered system is

$$\begin{aligned} \Delta E &\sim \delta\epsilon \left(\frac{L}{\xi_{\text{int}}} \right)^2 + \left(\frac{L}{\xi_*} \right)^2 \left[U \frac{\xi_*}{\xi_{\text{int}}} - U \left(\frac{\delta m \xi_{\text{dis}}}{U \xi_{\text{int}}} \right)^2 \frac{\xi_*}{\xi_{\text{int}}} \log \left(\frac{\xi_*}{a} \right) \right] \\ &= U \frac{L^2}{\xi_{\text{int}} \xi_*} \left[\frac{\delta\epsilon}{U} \frac{\xi_*}{\xi_{\text{int}}} + 1 - \alpha^2 \log \left(\frac{\xi_*}{a} \right) \right] \\ &= \frac{L^2}{\xi_*^2} E_{\phi\text{-dom}}(\xi_*), \end{aligned} \quad (\text{F26})$$

where this expression is the same for both white noise and Gaussian-correlated disorder provided we recall that $a = \xi_{\text{int}}$ in the former case while $a = \max(\xi_{\text{dis}}, \xi_{\text{int}})$ in the latter. It follows that when the typical domain size ξ_* is such that $E_{\phi\text{-dom}}(\xi_*) < 0$ (*i.e.* $L_- < \xi_* < L_+$), a wholly (but non-uniformly) ϕ -ordered sample may be

considered energetically favourable. An identical argument holds for $\alpha \gtrsim 1$ with $\xi_* = \xi_{\text{dis}}$. If we now imagine fixing α and increasing $\delta\epsilon$, we expect Eq. (F2) to hold up to some value, $\delta\tilde{\epsilon}_c(\alpha)$, of the same order as $\delta\epsilon_c(\alpha)$. Given the general lack of precision throughout this appendix, we assume that $\delta\tilde{\epsilon}_c(\alpha) \sim \delta\epsilon_c(\alpha)$. This identity sets the dashed line in Fig. 5.

-
- [1] Y. Cao, V. Fatemi, S. Fang, K. Watanabe, T. Taniguchi, E. Kaxiras, and P. Jarillo-Herrero, *Nature (London)* **556**, 43 (2018), arXiv:1803.02342 [cond-mat.mes-hall].
- [2] Y. Cao, V. Fatemi, A. Demir, S. Fang, S. L. Tomarken, J. Y. Luo, J. D. Sanchez-Yamagishi, K. Watanabe, T. Taniguchi, E. Kaxiras, R. C. Ashoori, and P. Jarillo-Herrero, *Nature (London)* **556**, 80 (2018), arXiv:1802.00553 [cond-mat.mes-hall].
- [3] J. M. B. Lopes Dos Santos, N. M. R. Peres, and A. H. Castro Neto, *Phys. Rev. Lett.* **99**, 256802 (2007), arXiv:0704.2128 [cond-mat.mtrl-sci].
- [4] R. Bistritzer and A. H. MacDonald, *Proceedings of the National Academy of Science* **108**, 12233 (2011), arXiv:1009.4203 [cond-mat.mes-hall].
- [5] Y. Kim, B. J. Wieder, C. L. Kane, and A. M. Rappe, *Phys. Rev. Lett.* **115**, 036806 (2015), arXiv:1504.03807 [cond-mat.mtrl-sci].
- [6] H. C. Po, L. Zou, A. Vishwanath, and T. Senthil, *Physical Review X* **8**, 031089 (2018), arXiv:1803.09742 [cond-mat.str-el].
- [7] R. de Gail, M. O. Goerbig, F. Guinea, G. Montambaux, and A. H. Castro Neto, *Phys. Rev. B* **84**, 045436 (2011), arXiv:1103.3172 [cond-mat.mes-hall].
- [8] W.-Y. He, Z.-D. Chu, and L. He, *Phys. Rev. Lett.* **111**, 066803 (2013), arXiv:1301.7573 [cond-mat.mes-hall].
- [9] M. Yankowitz, S. Chen, H. Polshyn, Y. Zhang, K. Watanabe, T. Taniguchi, D. Graf, A. F. Young, and C. R. Dean, *Science* **363**, 1059 (2019), arXiv:1808.07865 [cond-mat.mes-hall].
- [10] A. L. Sharpe, E. J. Fox, A. W. Barnard, J. Finney, K. Watanabe, T. Taniguchi, M. A. Kastner, and D. Goldhaber-Gordon, *Science* **365**, 605 (2019), arXiv:1901.03520 [cond-mat.mes-hall].
- [11] M. Serlin, C. L. Tschirhart, H. Polshyn, Y. Zhang, J. Zhu, K. Watanabe, T. Taniguchi, L. Balents, and A. F. Young, *Science* **367**, 900 (2020).
- [12] N. Bultinck, S. Chatterjee, and M. P. Zaletel, arXiv:1901.08110 [cond-mat.str-el] (2019).
- [13] Y.-H. Zhang, D. Mao, and T. Senthil, arXiv:1901.08209 [cond-mat.str-el] (2019).
- [14] X. Lu, P. Stepanov, W. Yang, M. Xie, M. A. Aamir, I. Das, C. Urgell, K. Watanabe, T. Taniguchi, G. Zhang, A. Bachtold, A. H. MacDonald, and D. K. Efetov, arXiv:1903.06513 [cond-mat.str-el] (2019).
- [15] A. Kerelsky, L. J. McGilly, D. M. Kennes, L. Xian, M. Yankowitz, S. Chen, K. Watanabe, T. Taniguchi, J. Hone, C. Dean, A. Rubio, and A. N. Pasupathy, *Nature* **572**, 95 (2019), arXiv:1812.08776 [cond-mat.mes-hall].
- [16] Y. Choi, J. Kemmer, Y. Peng, A. Thomson, H. Arora, R. Polski, Y. Zhang, H. Ren, J. Alicea, G. Refael, F. von Oppen, K. Watanabe, T. Taniguchi, and S. Nadj-Perge, *Nature Physics* **10.1038/s41567-019-0606-5** (2019), arXiv:1901.02997 [cond-mat.mes-hall].
- [17] Y. Jiang, X. Lai, K. Watanabe, T. Taniguchi, K. Haule, J. Mao, and E. Y. Andrei, *Nature* **573**, 91 (2019), arXiv:1904.10153 [cond-mat.mes-hall].
- [18] Y. Xie, B. Lian, B. Jäck, X. Liu, C.-L. Chiu, K. Watanabe, T. Taniguchi, B. A. Bernevig, and A. Yazdani, *Nature* **572**, 101 (2019), arXiv:1906.09274 [cond-mat.mes-hall].
- [19] A. Uri, S. Grover, Y. Cao, J. A. Crosse, K. Bagani, D. Rodan-Legrain, Y. Myasoedov, K. Watanabe, T. Taniguchi, P. Moon, M. Koshino, P. Jarillo-Herrero, and E. Zeldov, arXiv:1908.04595 [cond-mat.mes-hall] (2019).
- [20] J. H. Wilson, Y. Fu, S. Das Sarma, and J. H. Pixley, *1908.02753 [cond-mat.dis-nn]* (2019).
- [21] S. Liu, E. Khalaf, J. Y. Lee, and A. Vishwanath, arXiv:1905.07409 [cond-mat] (2020).
- [22] M. Xie and A. H. MacDonald, arXiv:1812.04213 [cond-mat.str-el] (2018).
- [23] F. Zhang, J. Jung, G. A. Fiete, Q. Niu, and A. H. MacDonald, *Phys. Rev. Lett.* **106**, 156801 (2011), arXiv:1010.4003 [cond-mat.str-el].
- [24] I. Martin, Y. M. Blanter, and A. F. Morpurgo, *Phys. Rev. Lett.* **100**, 036804 (2008), arXiv:0709.3522 [cond-mat.mes-hall].
- [25] Z. Qiao, J. Jung, Q. Niu, and A. H. MacDonald, *Nano Letters* **11**, 3453 (2011), arXiv:1107.4550 [cond-mat.mes-hall].
- [26] R. V. Gorbachev, J. C. W. Song, G. L. Yu, A. V. Kretinin, F. Withers, Y. Cao, A. Mishchenko, I. V. Grigorieva, K. S. Novoselov, L. S. Levitov, and A. K. Geim, *Science* **346**, 448 (2014).
- [27] M. Sui, G. Chen, L. Ma, W.-Y. Shan, D. Tian, K. Watanabe, T. Taniguchi, X. Jin, W. Yao, D. Xiao, and Y. Zhang, *Nature Physics* **11**, 1027 (2015).
- [28] Y. Shimazaki, M. Yamamoto, I. V. Borzenets, K. Watanabe, T. Taniguchi, and S. Tarucha, *Nature Physics* **11**, 1032 (2015).
- [29] L. Ju, Z. Shi, N. Nair, Y. Lv, C. Jin, J. Velasco, C. Ojeda-Aristizabal, H. A. Bechtel, M. C. Martin, A. Zettl, J. Analytis, and F. Wang, *Nature (London)* **520**, 650 (2015).
- [30] J. Li, K. Wang, K. J. McFaul, Z. Zern, Y. Ren, K. Watanabe, T. Taniguchi, Z. Qiao, and J. Zhu, *Nature Nanotechnology* **11**, 1060 (2016), arXiv:1509.03912 [cond-mat.mes-hall].
- [31] L.-J. Yin, H. Jiang, J.-B. Qiao, and L. He, *Nature Communications* **7**, 11760 (2016), arXiv:1511.06498 [cond-mat.mes-hall].
- [32] D. A. Abanin, S. A. Parameswaran, S. A. Kivelson, and S. L. Sondhi, *Phys. Rev. B* **82**, 035428 (2010), arXiv:1003.1978 [cond-mat.mes-hall].
- [33] S. A. Parameswaran and B. E. Feldman, *Journal of Physics Condensed Matter* **31**, 273001 (2019), arXiv:1809.09616 [cond-mat.str-el].
- [34] D. F. Mross, Y. Oreg, A. Stern, G. Margalit, and M. Heiblum, *Phys. Rev. Lett.* **121**, 026801 (2018).
- [35] C. Wang, A. Vishwanath, and B. I. Halperin, *Phys. Rev. B* **98**, 045112 (2018).

- [36] B. Lian and J. Wang, *Phys. Rev. B* **97**, 165124 (2018).
- [37] J. T. Chalker and P. D. Coddington, *Journal of Physics C: Solid State Physics* **21**, 2665 (1988).
- [38] D. K. K. Lee and J. T. Chalker, *Phys. Rev. Lett.* **72**, 1510 (1994), [arXiv:cond-mat/9311050 \[cond-mat\]](#).
- [39] D. K. K. Lee, J. T. Chalker, and D. Y. K. Ko, *Phys. Rev. B* **50**, 5272 (1994).
- [40] C. M. Ho and J. T. Chalker, *Physical Review B* **54**, 8708 (1996), [arXiv:cond-mat/9605073 \[cond-mat\]](#).
- [41] D.-H. Lee, *Phys. Rev. B* **50**, 10788 (1994), [cond-mat/9404011 \[cond-mat\]](#).
- [42] N. N. T. Nam and M. Koshino, *Phys. Rev. B* **96**, 075311 (2017), [arXiv:1706.03908 \[cond-mat.mtrl-sci\]](#).
- [43] M. Koshino, N. F. Q. Yuan, T. Koretsune, M. Ochi, K. Kuroki, and L. Fu, *Physical Review X* **8**, 031087 (2018), [arXiv:1805.06819 \[cond-mat.mes-hall\]](#).
- [44] H. C. Po, L. Zou, T. Senthil, and A. Vishwanath, *Phys. Rev. B* **99**, 195455 (2019), [arXiv:1808.02482 \[cond-mat.str-el\]](#).
- [45] J. Kang and O. Vafek, *Physical Review X* **8**, 031088 (2018), [arXiv:1805.04918 \[cond-mat.str-el\]](#).
- [46] L. Zou, H. C. Po, A. Vishwanath, and T. Senthil, *Physical Review B* **98**, 085435 (2018), [arXiv:1806.07873 \[cond-mat.str-el\]](#).
- [47] M. R. Zirnbauer, *Journal of Mathematical Physics* **37**, 4986 (1996), [arXiv:math-ph/9808012 \[math-ph\]](#).
- [48] A. Altland and M. R. Zirnbauer, *Phys. Rev. B* **55**, 1142 (1997), [arXiv:cond-mat/9602137 \[cond-mat\]](#).
- [49] A. W. W. Ludwig, M. P. A. Fisher, R. Shankar, and G. Grinstein, *Phys. Rev. B* **50**, 7526 (1994).
- [50] P. M. Ostrovsky, I. V. Gornyi, and A. D. Mirlin, *Phys. Rev. Lett.* **98**, 256801 (2007), [arXiv:cond-mat/0702115 \[cond-mat.mes-hall\]](#).
- [51] A. F. Morpurgo and F. Guinea, *Phys. Rev. Lett.* **97**, 196804 (2006), [arXiv:cond-mat/0603789 \[cond-mat.mes-hall\]](#).
- [52] A. Pruisken, *Nuclear Physics B* **235**, 277 (1984).
- [53] A. M. M. Pruisken, Field theory, scaling and the localization problem, in *The Quantum Hall Effect*, edited by R. E. Prange and S. M. Girvin (Springer New York, New York, NY, 1990) pp. 117–173.
- [54] I. L. Aleiner and K. B. Efetov, *Phys. Rev. Lett.* **97**, 236801 (2006), [arXiv:cond-mat/0607200 \[cond-mat.dis-nn\]](#).
- [55] A. Altland, *Phys. Rev. Lett.* **97**, 236802 (2006), [arXiv:cond-mat/0607247 \[cond-mat.mes-hall\]](#).
- [56] E. Abrahams, P. W. Anderson, D. C. Licciardello, and T. V. Ramakrishnan, *Phys. Rev. Lett.* **42**, 673 (1979).
- [57] P. A. Lee and T. V. Ramakrishnan, *Rev. Mod. Phys.* **57**, 287 (1985).
- [58] E. Fradkin, *Phys. Rev. B* **33**, 3257 (1986).
- [59] E. Fradkin, *Phys. Rev. B* **33**, 3263 (1986).
- [60] S. Carr, S. Fang, Z. Zhu, and E. Kaxiras, *Phys. Rev. Research* **1**, 013001 (2019), [arXiv:1901.03420 \[cond-mat.mes-hall\]](#).
- [61] F. D. M. Haldane, *Phys. Rev. Lett.* **61**, 2015 (1988).
- [62] C. Weeks and M. Franz, *Phys. Rev. B* **81**, 085105 (2010).
- [63] C. L. Kane and E. J. Mele, *Phys. Rev. Lett.* **95**, 226801 (2005).
- [64] G. Tarnopolsky, A. J. Kruchkov, and A. Vishwanath, *Phys. Rev. Lett.* **122**, 106405 (2019), [arXiv:1808.05250 \[cond-mat.str-el\]](#).
- [65] Y.-H. Zhang, D. Mao, Y. Cao, P. Jarillo-Herrero, and T. Senthil, *Phys. Rev. B* **99**, 075127 (2019), [arXiv:1805.08232 \[cond-mat.str-el\]](#).
- [66] J. Y. Lee, E. Khalaf, S. Liu, X. Liu, Z. Hao, P. Kim, and A. Vishwanath, [arXiv:1903.08685 \[cond-mat.str-el\]](#) (2019).
- [67] N. Bultinck, E. Khalaf, S. Liu, S. Chatterjee, A. Vishwanath, and M. P. Zaletel, *Phys. Rev. X* **10**, 031034 (2020).
- [68] V. S. Dotsenko and V. S. Dotsenko, *Advances in Physics* **32**, 129 (1983).
- [69] Y. Imry and S.-k. Ma, *Phys. Rev. Lett.* **35**, 1399 (1975).
- [70] T. Nattermann and J. Villain, *Phase Transitions* **11**, 5 (1988).
- [71] T. Nattermann, *Spin Glasses And Random Fields. Series: Series on Directions in Condensed Matter Physics* **12**, 277 (1997), [arXiv:cond-mat/9705295 \[cond-mat.stat-mech\]](#).
- [72] K. Binder, *Zeitschrift für Physik B Condensed Matter* **50**, 343 (1983).
- [73] M. Aizenman and J. Wehr, *Phys. Rev. Lett.* **62**, 2503 (1989).
- [74] M. Blume, *Phys. Rev.* **141**, 517 (1966).
- [75] H. Capel, *Physica* **32**, 966 (1966).
- [76] M. Kaufman and M. Kanner, *Phys. Rev. B* **42**, 2378 (1990).
- [77] R. Vasseur and T. Lookman, *Phys. Rev. B* **81**, 094107 (2010).
- [78] Y.-Z. Chou, R. M. Nandkishore, and L. Radzihovsky, *Phys. Rev. B* **99**, 165108 (2019), [arXiv:1901.05464 \[cond-mat.str-el\]](#).
- [79] F. Evers and A. D. Mirlin, *Reviews of Modern Physics* **80**, 1355 (2008), [arXiv:0707.4378 \[cond-mat.mes-hall\]](#).
- [80] D. K. Efimkin and A. H. MacDonald, *Phys. Rev. B* **98**, 035404 (2018), [arXiv:1803.06404 \[cond-mat.mes-hall\]](#).
- [81] X.-C. Wu, C.-M. Jian, and C. Xu, *Phys. Rev. B* **99**, 161405 (2019), [arXiv:1811.08442 \[cond-mat.str-el\]](#).
- [82] Y.-Z. Chou, Y.-P. Lin, S. Das Sarma, and R. M. Nandkishore, *Phys. Rev. B* **100**, 115128 (2019).
- [83] C. Chen, A. H. Castro Neto, and V. M. Pereira, *Phys. Rev. B* **101**, 165431 (2020).
- [84] C. R. Woods, L. Britnell, A. Eckmann, R. S. Ma, J. C. Lu, H. M. Guo, X. Lin, G. L. Yu, Y. Cao, R. V. Gorbachev, A. V. Kretinin, J. Park, L. A. Ponomarenko, M. I. Katsnelson, Y. N. Gornostyrev, K. Watanabe, T. Taniguchi, C. Casiraghi, H. J. Gao, A. K. Geim, and K. S. Novoselov, *Nature Physics* **10**, 451 (2014), [arXiv:1401.2637 \[cond-mat.mes-hall\]](#).
- [85] M. M. van Wijk, A. Schuring, M. I. Katsnelson, and A. Fasolino, *2D Materials* **2**, 034010 (2015), [arXiv:1503.02540 \[cond-mat.mes-hall\]](#).
- [86] S. Dai, Y. Xiang, and D. J. Srolovitz, *Nano Letters* **16**, 5923 (2016).
- [87] N. Y. Kim, H. Y. Jeong, J. H. Kim, G. Kim, H. S. Shin, and Z. Lee, *ACS Nano* **11**, 7084 (2017).
- [88] F. Gargiulo and O. V. Yazyev, *2D Materials* **5**, 015019 (2018), [arXiv:1711.08647 \[cond-mat.mes-hall\]](#).
- [89] K. Zhang and E. B. Tadmor, *Journal of Mechanics Physics of Solids* **112**, 225 (2018).
- [90] We thank Cory Dean for sharing these ideas with us.
- [91] D. A. Abanin, A. V. Shytov, L. S. Levitov, and B. I. Halperin, *Physical Review B* **79**, 035304 (2009).
- [92] D. R. Hofstadter, *Phys. Rev. B* **14**, 2239 (1976).
- [93] J. Tworzyno, B. Trauzettel, M. Titov, A. Rycerz, and C. W. J. Beenakker, *Phys. Rev. Lett.* **96**, 246802 (2006), [arXiv:cond-mat/0603315 \[cond-mat.mes-hall\]](#).
- [94] N. M. R. Peres, A. H. Castro Neto, and F. Guinea, *Phys. Rev. B* **73**, 195411 (2006), [arXiv:cond-mat/0512476 \[cond-mat.mes-hall\]](#).
- [95] S. Das Sarma, S. Adam, E. H. Hwang, and E. Rossi, *Reviews of Modern Physics* **83**, 407 (2011), [arXiv:1003.4731 \[cond-mat.mes-hall\]](#).

- [96] T. Ando, *Journal of the Physical Society of Japan* **75**, 074716 (2006).
- [97] V. V. Cheianov and V. I. Fal'ko, *Phys. Rev. Lett.* **97**, 226801 (2006), [arXiv:cond-mat/0608228](#) [cond-mat.mes-hall].
- [98] K. Nomura and A. H. MacDonald, *Phys. Rev. Lett.* **96**, 256602 (2006), [arXiv:cond-mat/0604113](#) [cond-mat.mes-hall].
- [99] E. H. Hwang, S. Adam, and S. D. Sarma, *Phys. Rev. Lett.* **98**, 186806 (2007), [arXiv:cond-mat/0610157](#) [cond-mat.mes-hall].
- [100] K. Nomura and A. H. MacDonald, *Phys. Rev. Lett.* **98**, 076602 (2007), [arXiv:cond-mat/0606589](#) [cond-mat.mes-hall].
- [101] M. Trushin and J. Schliemann, *EPL (Europhysics Letters)* **83**, 17001 (2008), [arXiv:0802.2794](#) [cond-mat.mes-hall].
- [102] M. I. Katsnelson, F. Guinea, and A. K. Geim, arXiv e-prints, [arXiv:0901.1398](#) (2009), [arXiv:0901.1398](#) [cond-mat.mes-hall].
- [103] M. I. Katsnelson and A. K. Geim, *Philosophical Transactions of the Royal Society of London Series A* **366**, 195 (2008), [arXiv:0706.2490](#) [cond-mat.mes-hall].
- [104] T. Nattermann, *Journal of Physics C: Solid State Physics* **16**, 6407 (1983).
- [105] E. T. Seppälä, V. Petäjä, and M. J. Alava, *Phys. Rev. E* **58**, R5217 (1998).
- [106] E. T. Seppälä and M. J. Alava, *Physical Review E* **63**, 10.1103/physreve.63.066109 (2001).