

Topological insulators are tunable waveguides for hyperbolic polaritons

Jih-Sheng Wu (吳致盛), D. N. Basov, and M. M. Fogler

University of California San Diego, 9500 Gilman Drive, La Jolla, California 92093, USA

(Received 20 September 2015; published 30 November 2015)

We present a theoretical analysis showing that layered topological insulators, for example, Bi_2Se_3 are optically hyperbolic materials in the range of terahertz (THz) frequencies. As such, these topological insulators possess deeply subdiffractive, highly directional collective modes: hyperbolic phonon polaritons. We predict that in thin crystals the dispersion of these modes is split into discrete subbands and is strongly influenced by electron surface states. If the surface states are doped, then hybrid collective modes result from coupling of the phonon polaritons with surface plasmons. The strength of the hybridization can be controlled by an external gate that varies the chemical potential of the surface states. We also show that the momentum dependence of the plasmon-phonon coupling leads to a polaritonic analog of the Goos-Hänchen effect. The directionality of the polaritonic rays and their tunable Goos-Hänchen shift is observable via THz nanoimaging.

DOI: [10.1103/PhysRevB.92.205430](https://doi.org/10.1103/PhysRevB.92.205430)

PACS number(s): 73.21.-b, 78.20.-e

I. INTRODUCTION

Bismuth-based topological insulators (TIs) have attracted much interest for their unusual electron surface states (SSs), which behave as massless Dirac fermions [1,2]. However, the bulk optical response of these compounds [3–15] is also remarkable. The quintuple-layered structure of these materials causes a strong anisotropy of their phonon modes. The E_u phonons that involve atomic displacements in a plane parallel to the basal plane (henceforth, x - y or \perp plane) have lower frequencies than A_{2u} , the c -axis (henceforth, z -axis) vibrations [5]. For Bi_2Se_3 , the dominant \perp - and z -axis phonon frequencies,

$$\begin{aligned}\omega_{1,\text{to}}^\perp &= 64 \text{ cm}^{-1} = 1.9 \text{ THz}, \\ \omega_{1,\text{to}}^z &= 135 \text{ cm}^{-1} = 4.1 \text{ THz},\end{aligned}\quad (1)$$

differ more than twice. As a result, this and similar TIs can exhibit a giant anisotropy of the dielectric permittivity. There is a range of ω where the permittivity tensor is indefinite: the real part of $\epsilon^z(\omega)$ is positive, while that of $\epsilon^\perp(\omega)$ is negative. Media with such characteristics are referred to as hyperbolic [16–18] because the isofrequency surfaces of their extraordinary rays in the momentum space $\mathbf{k} = (k^x, k^y, k^z)$ are shaped as hyperboloids [Fig. 1(a)]. In the terahertz (THz) domain, the widest band of frequencies where Bi_2Se_3 behaves as a hyperbolic medium (HM) is between the aforementioned dominant frequencies, $\omega_{\text{to},1}^\perp < \omega < \omega_{\text{to},1}^z$; however, other hyperbolic bands also exist in this TI (both at THz frequencies, see Sec. II, and at visible frequencies, see Ref. [19]). It is important that the approximate equation for the extraordinary isofrequency surfaces,

$$\frac{(k^x)^2 + (k^y)^2}{\epsilon^z(\omega)} + \frac{(k^z)^2}{\epsilon^\perp(\omega)} = \frac{\omega^2}{c^2}, \quad (2)$$

is valid up to $|\mathbf{k}|$ of the order of the inverse lattice constant. Accordingly, rays of momenta $|\mathbf{k}|$ greatly exceeding the free-space photon momentum ω/c can propagate through hyperbolic materials without evanescent decay. At such \mathbf{k} , the hyperboloids can be further approximated by cones, which means that the group velocity $\mathbf{v} = \partial\omega/\partial\mathbf{k}$ of the rays makes

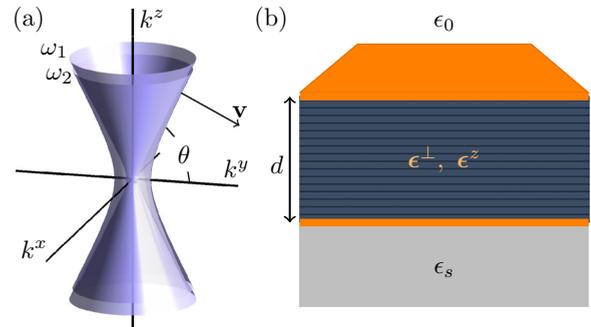


FIG. 1. (Color online) (a) Hyperboloidal isofrequency surfaces of HP^2 s for two frequencies ω_1 and ω_2 ($\omega_2 > \omega_1$). The asymptote angle θ with respect to the k^x - k^y plane is shown; the group velocity \mathbf{v} makes the same angle with respect to the k^z axis. (b) Model geometry: a TI slab of thickness d sandwiched between a substrate of permittivity ϵ_s and a superstrate of permittivity ϵ_0 . The two thin (orange) layers represent the top and the bottom surface states.

a *fixed* angle θ (or $-\theta$) with respect to the z axis, with

$$\tan \theta(\omega) = i \frac{[\epsilon^\perp(\omega)]^{1/2}}{[\epsilon^z(\omega)]^{1/2}}, \quad (3)$$

see Fig. 1(a). We refer to these deeply subdiffractive, highly directional modes as the hyperbolic phonon polaritons (HPP or HP^2 for short).

Our interest to HP^2 of TIs is stimulated by the recent discovery [20,21] and further exploration of similar collective modes in other systems such as hexagonal boron nitride [22–25] (hBN) and hBN covered by graphene [26–28] (hBN/G). There is a close analogy between these systems. In fact, except for the difference in the number of Dirac cones ($N = 1$ versus $N = 4$) and the frequency range where the hyperbolic response occurs (THz versus midinfrared), the electrodynamics of longitudinal collective modes of Bi_2Se_3 and hBN/G structures is qualitatively the same. (The analogy is the most faithful when graphene and hBN are rotationally misaligned; otherwise, their collective modes are modified by the moiré superlattice effects [28,29].)

The main goal of this paper is to investigate the interaction of HP^2 with the Dirac plasmons of the topological SS. The

latter dominate the charge (and current) density response of the system at frequencies outside the hyperbolic band where HP^2 are absent. Dirac plasmons have been extensively studied in previous literature [8,13,14,30–44] on both TI and graphene. The basic properties of the Dirac plasmons can be introduced on the example of a hypothetical TI material with a frequency-independent permittivity $\epsilon^z > 0$ and the permittivity $\epsilon^\perp(\omega)$ dominated by a single phonon mode. Such an idealized material is hyperbolic in a single frequency interval $\omega_{l0} < \omega < \omega_{h0}$, where $\epsilon^\perp(\omega) < 0$. Its Dirac plasmons exist at $\omega < \omega_{l0}$ and $\omega > \omega_{h0}$ where $\epsilon^\perp(\omega) > 0$. In the setup shown in Fig. 1(b), where the TI slab borders media of constant permittivities $\epsilon_0 > 0$ and $\epsilon_s > 0$, there are two plasmon modes. At large enough in-plane momenta $q \equiv [(k^x)^2 + (k^y)^2]^{1/2}$, these modes are confined to the opposite interfaces and electromagnetically decoupled. In the relevant range of momenta $q < q_*$, the dispersion of the plasmon bound to the top interface is given by

$$q(\omega) \simeq \frac{4}{N} \frac{\epsilon_0 + \epsilon_1}{e^2 |\mu|} (\hbar\omega)^2, \quad \hbar\omega \ll |\mu|, \quad (4)$$

where

$$\epsilon_1(\omega) = [\epsilon^\perp(\omega)]^{1/2} [\epsilon^z(\omega)]^{1/2}, \quad (5)$$

is the effective permittivity of the TI and μ is the chemical potential of the SSs measured from the Dirac point. At frequencies far below ω_{l0} or far above ω_{h0} , function $\epsilon_1(\omega)$ can be approximated by a real constant, which yields $\omega \propto \sqrt{q}$. This typical two-dimensional (2D) plasmon dispersion describes the low-frequency part of the full curve sketched in Fig. 2(a). The plasmon dispersion for the bottom interface is obtained by replacing ϵ_0 with ϵ_s (unless $\epsilon_s \gg \epsilon_0$, in which case the range $q > q_*$ is relevant where the dispersion is approximately linear, see Sec. III B).

Equation (4) implies that the nature of the plasmon modes should change drastically when ω enters the hyperbolic frequency band where $\epsilon_1(\omega)$ [Eq. (5)] is imaginary and strongly ω -dependent. This equation predicts a complex q , which suggests that the Dirac plasmons become leaky modes that rapidly decay into the HP^2 bulk continuum. However, this is not quite correct. We will show that nonleaky, i.e., propagating modes can survive in thin enough TI slabs where the HP^2 continuum is broken into discrete subbands of *waveguide* modes. The latter hybridize with plasmons to form hyperbolic plasmon phonon polaritons (HPPP or HP^3 for short), the primary target of our investigation, see Figs. 2(b) and 2(c). We explore the following properties and manifestations of the collective charge modes of the TIs: (i) the mode dispersion in the momentum-frequency space, (ii) the dependence of such dispersions on the surface doping and the thickness of the slab, and (iii) the unusual real-space dynamics of the HP^3 rays, including a polaritonic analog of the Goos-Hänchen (GH) effect [45,46].

The remainder of the paper is organized as follows. In Sec. II, we specify the model and the basic equations. In Sec. III, we present our results for the dispersion of the three different types of collective modes (plasmons, HP^2 s, and HP^3 s). In Sec. IV, which is the centerpiece of this work, we discuss waveguiding and launching of the HP^2 modes and also their tunable GH shifts. We explain how these phenomena can be probed experimentally using the imaging capabilities

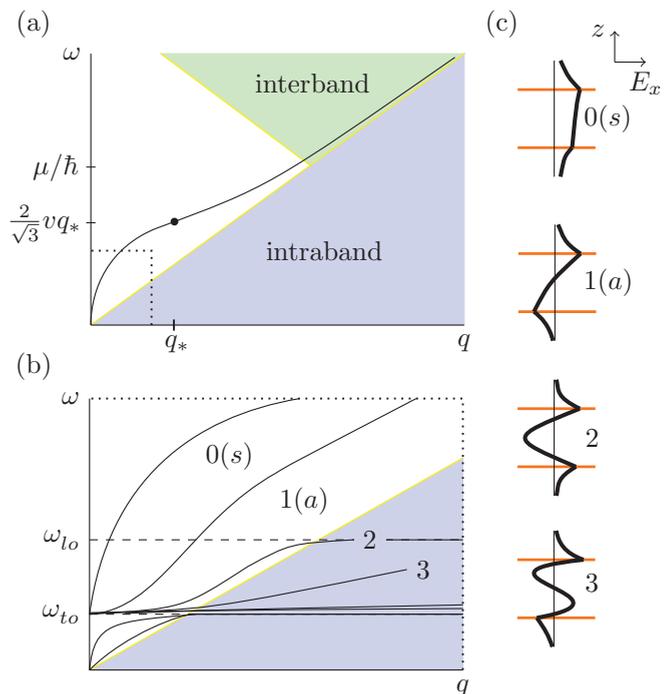


FIG. 2. (Color online) Schematic illustrations of the collective mode spectra in idealized model systems. (a) The plasmon dispersion of Dirac fermions confined to the interface of two bulk media of constant positive permittivity ϵ_0 and ϵ_s . The dispersion crosses over from $\omega \simeq v\sqrt{q}q_*/2$ to $\omega \simeq vq$ at a characteristic momentum q_* [Eq. (26)]. The shaded areas indicate the electron-hole continua where the plasmons (and any other charged collective modes) are damped. (b) The dispersion of hybrid HP^3 modes for a slab of a hypothetical TI material that has a single in-plane phonon mode at ω_{l0} and constant $\epsilon^z > 0$. Permittivity ϵ^\perp is negative at $\omega_{l0} < \omega < \omega_{h0}$ and positive at other ω . The dotted boundary corresponds to the dotted line in (a). Outside the band $\omega_{l0} < \omega < \omega_{h0}$, only plasmonic modes 0 and 1 exist. In the degenerate case $\epsilon_0 = \epsilon_s$ they correspond to the symmetric (*s*) and antisymmetric (*a*) combinations of the top and bottom interface plasmons. Inside that band, multiple branches of HP^3 are formed due to hybridization of the plasmons with the HP^2 waveguide modes. The frequencies of all the branches other than 0 and 1 tend to ω_{l0} at large momenta. (c) Schematic in-plane electric field profiles of the first few HP^3 modes (thick curves). The number of nodes in each profile (the points where they cross with the vertical lines $E_x = 0$) is equal to the modal index.

of the scattering-type scanning near-field optical microscopy (s-SNOM) [47,48]. In Sec. V, we give concluding remarks and an outlook for the future. Finally, in Appendix, we discuss signatures of the phonon-plasmon coupling measurable by the s-SNOM operating in the spectroscopic mode.

II. MODEL

Our model for the bulk permittivities of the TI is

$$\epsilon^\alpha(\omega) = \epsilon_\infty^\alpha + \sum_{j=1,2} \frac{\omega_{p,j}^{\alpha 2}}{\omega_{l0,j}^{\alpha 2} - \omega^2 - i\gamma_j^\alpha \omega}, \quad \alpha = \perp, z. \quad (6)$$

In the case of Bi_2Se_3 , we choose the parameters based on available experimental [3,4,7] and theoretical [5] literature as follows: $\epsilon_\infty^\perp = 29$, $\epsilon_\infty^z = 17.4$, $\omega_{l0,1}^\perp = 64 \text{ cm}^{-1}$, $\omega_{p,1}^\perp =$

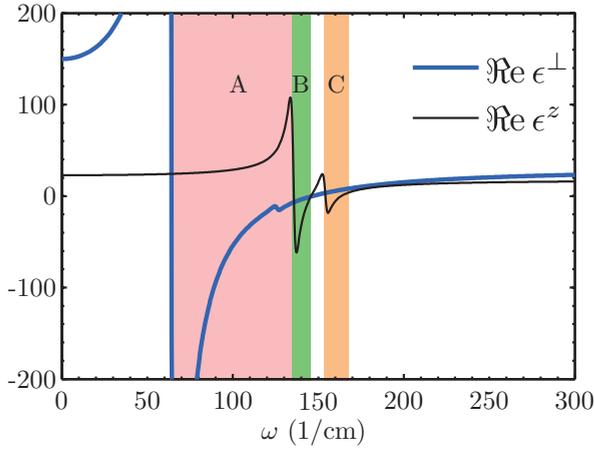


FIG. 3. (Color online) The real parts of the tangential and axial permittivities of Bi_2Se_3 . The sign changes of the permittivities are due to the E_u and A_{2u} phonons. Surface- and bulk-confined collective modes exist inside the spectral regions where at least one of the permittivities is negative. They include the type-II hyperbolic region A ($\Re \epsilon^\perp < 0, \Re \epsilon^z > 0$), the reststrahlen region B ($\Re \epsilon^\perp, \Re \epsilon^z < 0$), and the type-I hyperbolic region C ($\Re \epsilon^\perp > 0, \Re \epsilon^z < 0$).

704 cm^{-1} , $\omega_{\text{to},2}^\perp = 125 \text{ cm}^{-1}$, $\omega_{p,2}^\perp = 55 \text{ cm}^{-1}$, $\omega_{\text{to},1}^z = 135 \text{ cm}^{-1}$, $\omega_{p,1}^z = 283 \text{ cm}^{-1}$, $\omega_{\text{to},2}^z = 154 \text{ cm}^{-1}$, $\omega_{p,2}^z = 156 \text{ cm}^{-1}$, and $\gamma_f^z = 3.5 \text{ cm}^{-1}$. [Note that $\omega_{\text{to},1}^\perp$ and $\omega_{\text{to},1}^z$ were already listed in Eq. (1).] The real parts of functions $\epsilon^\perp(\omega)$ and $\epsilon^z(\omega)$ are plotted in Fig. 3. The regions where at least one of them is negative are shaded. They include region A, $\omega_{\text{to},1}^\perp < \omega < \omega_{\text{to},1}^z$, where Bi_2Se_3 is an HM of type II ($\Re \epsilon_z > 0, \Re \epsilon^\perp < 0$); region C, $\omega_{\text{to},2}^z < \omega < 163 \text{ cm}^{-1}$ where it is an HM of type I ($\Re \epsilon_z < 0, \Re \epsilon^\perp > 0$), and region B, $\omega_{\text{to},1}^z < \omega < 146 \text{ cm}^{-1}$, where it exhibits the Reststrahlen behavior ($\Re \epsilon_z < 0, \Re \epsilon^\perp < 0$). Since regions B and C are narrow, in our discussion of HP^2 and HP^3 modes we focus on region A. In this discussion, we often refer to hBN as an example of a simpler material. The type-II hyperbolic band of hBN is bounded by the frequencies [20,22]

$$\omega_{\text{to}} = 1376 \text{ cm}^{-1}, \quad \omega_{\text{lo}} = 1614 \text{ cm}^{-1}. \quad (7)$$

In this band, $\epsilon^\perp(\omega)$ of hBN can be modeled similar to Eq. (6) but using a single Lorentzian oscillator while ϵ^z can be considered ω -independent and positive.

In the case of Bi_2Se_3 , we also have to specify our assumptions about the electronic response. We consider only frequencies smaller than the bulk gap 0.3 eV of Bi_2Se_3 at which the electronic contribution to the permittivities [included in Eq. (6) via ϵ_∞^α] is purely real. Additionally, we assume that the valence bulk band is completely filled, the conduction one is empty, with no free carriers present in the bulk. However, such carriers populate the gapless SS described by the massless 2D Dirac equation. The chemical potential μ , which is located inside the bulk band gap, determines the doping of these SS. For simplicity, we ignore any virtual or real electronic transitions between the surface and the bulk states, which should not change the result qualitatively, except perhaps for the additional damping from these transitions.

The fundamental current/density response functions of the SS are the sheet conductivity σ and polarizability P , which

are related in the standard way:

$$\sigma(q, \omega) = \frac{i\omega}{q^2} e^2 P(q, \omega). \quad (8)$$

Within the random-phase approximation for Dirac fermions, $P(q, \omega)$ can be computed [49,50] analytically:

$$P(q, \omega) = -\frac{Nk_F}{2\pi\hbar v} - \frac{iN}{16\pi\hbar v} \frac{q^2}{\sqrt{q^2 - k_\omega^2}} \\ \times \left[G\left(\frac{k_\omega + 2k_F}{q}\right) - G\left(\frac{k_\omega - 2k_F}{q}\right) - i\pi \right], \\ G(x) = ix\sqrt{1 - x^2} - i \arccos x. \quad (9)$$

Here, the branch cut for the square root and logarithm functions is the negative real semi-axis, k_ω is defined by $k_\omega = (\omega + i\gamma_e)/v$, phenomenological parameter $\gamma_e > 0$ is the electron scattering rate, v is the Fermi velocity, and $k_F = |\mu|/(\hbar v)$ is the Fermi momentum. Equation (9) is a good approximation at small μ . At large doping, trigonal warping [51] and other details of realistic band structure [43] should be included. Since the above formula is a bit cumbersome, it may be helpful to mention some properties of $\sigma(q, \omega)$. For example, if $\gamma_e = +0$, the real part of $\sigma(q, \omega)$ is nonvanishing only inside the two shaded areas in Fig. 2(a), which together form the so-called electron-hole continuum [30,39]. (This real part is a measure of dissipation, i.e., Landau damping.) For a doped system at small momenta and frequencies, $q, k_\omega \ll k_F$, the expression for the conductivity can be reduced to

$$\sigma(q, \omega) \simeq \frac{Ne^2}{2\pi\hbar} \frac{k_F}{\sqrt{q^2 - k_\omega^2}} \frac{ik_\omega}{ik_\omega - \sqrt{q^2 - k_\omega^2}}. \quad (10)$$

At $q \ll \omega/v$, it further simplifies to the Drude formula

$$\sigma \simeq \frac{Ne^2}{4\pi\hbar^2} \frac{|\mu|}{\gamma_e - i\omega}, \quad \mu \neq 0. \quad (11)$$

For an undoped system, one finds instead

$$\sigma(q, \omega) = \frac{N}{16} \frac{e^2}{\hbar} \frac{ik_\omega}{\sqrt{q^2 - k_\omega^2}} \quad (12)$$

$$\simeq \frac{N}{16} \frac{e^2}{\hbar}, \quad q \ll \frac{\omega}{v}. \quad (13)$$

In order to find the dispersion of the collective modes of the TI slab, we use two computational methods. One method, which is advantageous for deriving analytical results, is to look for the poles of the response function $r_P(q, \omega)$. This function is the total P - (also known as the TM-) polarization reflectivity of the system measured when an external field is incident from the medium labeled “ ϵ_0 ” in Fig. 1(b). It must be immediately clarified that $r_P(q, \omega)$ has no poles at simultaneously real q and ω if the dissipation parameters γ and γ_e are nonzero. At least one of these arguments must be complex. Whenever one refers to the dispersion relation of a mode, one means the relation between the real parts of q and ω . The other method, which is especially convenient for numerical simulations, is to identify the sought dispersion curves with the maxima of $\Im m r_P(q, \omega)$ at *real* arguments. As long as the imaginary parts of q and ω (which give information about the propagation length and lifetime of the mode) are small, both methods give

the same dispersions. An extra benefit of working with real q and ω is that the corresponding $r_P(q, \omega)$ is the input for further calculations we discuss in Appendix where we model s-SNOM experiments for the system in hand.

Our procedure for calculating function $r_P(q, \omega)$ can be explained as follows. Taking a more general view for a moment, we regard the entire system including the substrate and superstrate as a stack of layers $j = 0, 1, \dots, M$ of thickness d_j , tangential permittivity ϵ_j^\perp , and axial permittivity ϵ_j^z . (In the present case, $M = 2$, the TI slab is layer $j = 1$ and $d_1 = d$.) Additionally, we assume that the interface of the layers j and $j + 1$ possesses the sheet conductivity $\sigma_{j,j+1}$. We observe that the P -polarization reflectivity $r_{j,j+1}$ of $j, j + 1$ interface in isolation is given by the formula (see, e.g., Ref. [27])

$$r_{j,j+1} = \frac{Q_{j+1} - Q_j + \frac{4\pi}{\omega} \sigma_{j,j+1}}{Q_{j+1} + Q_j + \frac{4\pi}{\omega} \sigma_{j,j+1}}, \quad (14)$$

$$Q_j = \frac{\epsilon_j^\perp}{k_j^z}, \quad k_j^z = \sqrt{\epsilon_j^\perp \left(\frac{\omega^2}{c^2} - \frac{q^2}{\epsilon_j^z} \right)}, \quad (15)$$

where k_j^z and q are, respectively, the axial and the tangential momenta inside layer j . Let r_j be the reflectivity of a subsystem composed of layers j, \dots, M . By this definition, $r_{M-1} = r_{M-1,M}$. The crucial point is that the desired $r_P \equiv r_0$ can be found by the backward recursion

$$r_j = r_{j,j+1} - \frac{(1 - r_{j,j+1})(1 - r_{j+1,j})r_{j+1}}{r_{j+1,j}r_{j+1} - \exp(-2ik_{j+1}d_{j+1})}, \quad (16)$$

where $r_{j+1,j}$ is the right-hand side of Eq. (14) with Q_j and Q_{j+1} interchanged. For $M = 2$, one recursion step suffices, which gives us, after some algebra [27],

$$r_P = \frac{r_{12}(r_{01} + r_{10} - 1) - r_{01} \exp(-2ik_1d_1)}{r_{10}r_{12} - \exp(-2ik_1d_1)}. \quad (17)$$

Hence function $r_P(q, \omega)$ has poles whenever

$$r_{10}(q, \omega)r_{12}(q, \omega) = \exp(-2ik_1^z d). \quad (18)$$

For large in-plane momenta $q \gg (\omega/c) \max |\epsilon_j^z|^{1/2}$, we can use the approximations $k_1^z \simeq q \tan \theta$ and

$$r_{10} \simeq \frac{\epsilon_0 - \epsilon_1 - \frac{2q}{q_{\text{top}}}}{\epsilon_0 + \epsilon_1 - \frac{2q}{q_{\text{top}}}}, \quad q_{\text{top}} \equiv \frac{i\omega}{2\pi\sigma_{\text{top}}}, \quad (19)$$

where $\sigma_{\text{top}} = \sigma_{\text{top}}(q, \omega)$ is the sheet conductivity of the SS at the top interface. Let us also define the ‘‘phase shifts’’ ϕ_{top} and ϕ_{bot} for inner reflections from the top and bottom interfaces, respectively: $r_{10} = -\exp(2i\phi_{\text{top}})$, $r_{12} = -\exp(2i\phi_{\text{bot}})$. Note that in general ϕ_{top} and ϕ_{bot} are complex numbers. Specifically,

we take

$$\phi_{\text{top}} = \arctan \left[i \frac{\epsilon_0}{\epsilon_1} \left(1 - \frac{2}{\epsilon_0} \frac{q}{q_{\text{top}}} \right) \right], \quad (20)$$

$$\phi_{\text{bot}} = \arctan \left[i \frac{\epsilon_s}{\epsilon_1} \left(1 - \frac{2}{\epsilon_s} \frac{q}{q_{\text{bot}}} \right) \right], \quad (21)$$

where the standard definition of $\arctan z$ is assumed, with the branch cuts $(-i\infty, -i), (i, i\infty)$ in the complex- z plane; q_{bot} is defined analogously to q_{top} but with the sheet conductivity σ_{bot} of the bottom SS instead of σ_{top} . Equation (18) can now be transformed to

$$q_n = -\frac{2}{\delta}(n\pi + \phi_{\text{top}} + \phi_{\text{bot}}), \quad \delta \equiv 2d \tan \theta, \quad (22)$$

where the integer subscript n labels possible multiple solutions. Admissible n must satisfy the condition $\Im m r_P(q_n, \omega) > 0$. Our numerical results for r_P computed from Eq. (17) and analytic approximations for the solutions of Eq. (22) are presented in Sec. III.

III. COLLECTIVE MODE DISPERSIONS

The false color maps of function $\Im m r_P(q, \omega)$ provide a convenient visualization of the collective mode spectra. Examples of such maps computed for Bi_2Se_3 slabs are presented in the bottom row of Fig. 4. Their counterparts for graphene-hBN-graphene (G/hBN/G) structures are shown in the top row to facilitate the interpretation. The bright lines in Fig. 4 are the dispersion curves of the collective modes. The apparent widths of those lines give an idea how damped the modes are. Below we discuss these results in more detail.

A. Hyperbolic waveguide modes

Figures 4(a) and 4(d) depict the $\Im m r_P$ maps for, respectively, G/hBN/G and Bi_2Se_3 slabs, when they are undoped, $\mu = 0$. No Dirac plasmons exist in such systems, so that the collective modes are limited to HP²s. In Fig. 4(a), we see a single family of such modes whereas in 4(d) one can actually distinguish three of them. Let us start with the former, simpler case. The key to understanding the nature of these modes is that inside the hyperbolic band $\omega_{\text{to}} < \omega < \omega_{\text{lo}}$ the z -axis momentum $k_1^z \simeq q \tan \theta$ of the modes is nearly real. Hence, the HP²s form standing waves inside the slab. The integer n in Eq. (22) corresponds to the number of nodes of these waves, see Fig. 2(c). For G/hBN/G, the requisite condition $\Im m r_P > 0$ is satisfied by all nonnegative integers n due to the fact that $\Im m \tan \theta > 0$. This inequality also ensures that $\Im m q > 0$. An analytical approximation for the dispersion curves of an undoped slab is obtained by neglecting the fractions $q/q_{\text{top}}, q/q_{\text{bot}}$ in Eqs. (20) and (21), in which case Eq. (22) yields $q(\omega)$ directly. Within this approximation, momenta q_n at given ω are equidistant:

$$q_{n+1} - q_n \simeq -\frac{2\pi}{\delta} = -\frac{\pi}{d \tan \theta(\omega)}. \quad (23)$$

The dispersion of the HP² waveguide modes is dominated by the factor $1/\tan \theta(\omega)$ in Eqs. (22) and (23), which, if all damping is neglected, changes from zero to infinity as ω increases from ω_{to} to ω_{lo} . This is precisely what we see in

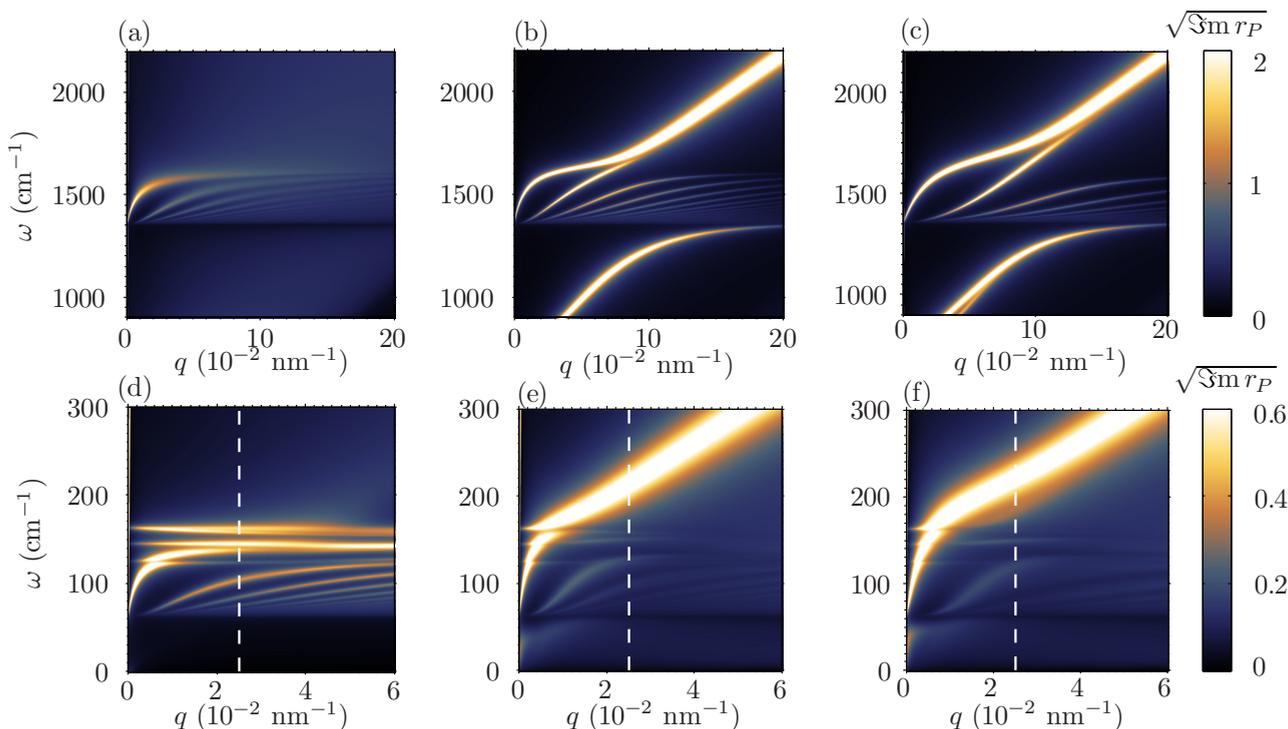


FIG. 4. (Color online) Collective mode dispersions of graphene-hBN-graphene (G/hBN/G) and Bi_2Se_3 slabs rendered using the false color maps of $\sqrt{\text{Im } r_P}$. The parameters of the calculation for G/hBN/G are: (a) $d = 60$ nm, $\mu = 0$, (b) $d = 60$ nm, $\mu = 0.29$ eV, (c) $d = 30$ nm, $\mu = 0.29$ eV. The other parameters are $v = 1.00 \times 10^8$ cm/s, $\gamma_e = 1.00$ THz, $\epsilon_0 = 1$, and $\epsilon_s = 1.5$. The parameters of the calculation for Bi_2Se_3 are (d) $d = 120$ nm, $\mu = 0$, (e) $d = 120$ nm, $\mu = 0.29$ eV, (f) $d = 60$ nm, $\mu = 0.29$ eV. In these three plots, $v = 0.623 \times 10^8$ cm/s, $\gamma_e = 1$ THz, $\epsilon_0 = 1$, and $\epsilon_s = 10$. Equal doping of the top and bottom SS is assumed. The vertical dashed lines indicate a characteristic momentum probed by the s-SNOM experiments simulated in Fig. 7 below.

Fig. 4(a): all the dispersion curves start at ω_{t0} at $q = 0$ and increase toward ω_{t0} at large q .

Equation (23) is general and it applies to Bi_2Se_3 as well. The three families of collective modes seen in Fig. 4(d), belong to the spectral regions A, B, and C of Fig. 3. In region A, which is the widest of the three, we see a set of HP^2 modes very similar to those in Fig. 4(a). They start at $\omega_{t0,1} = 64$ cm^{-1} at $q = 0$ and monotonically increase toward $\omega_{t0,2} = 135$ cm^{-1} at large q . In region C, $154 < \omega(\text{cm}^{-1}) < 163$, we again find a family of HP^2 modes but this time with a negative dispersion. This behavior is typical of type I HM ($\Re \epsilon^\perp > 0, \Re \epsilon^z < 0$). The shape of the dispersion can be understood noticing that the real part of $1/\tan \theta(\omega)$ is positive, varying from ∞ to 0 (if the phonon damping γ_j^α is neglected) while admissible n are now $n \leq 0$. [In hBN, this type I behavior is also realized [22,24,27] but the corresponding frequency range is below the axis cutoff in Fig. 4(a).] Lastly, in region B, $135 < \omega(\text{cm}^{-1}) < 146$, function $\tan \theta(\omega)$ is almost purely imaginary, which implies that the collective modes do not form standing waves but are exponentially confined to the interfaces. Also, there are only two such modes, $n = 0$ and 1. In this respect, these surface-bound HP^2 modes are similar to the Dirac plasmons, see Sec. I above and Sec. III B below. However, their dispersion is completely different from those of the plasmons, e.g., the dispersion of the upper ($n = 1$) mode has a negative slope, see Fig. 4(d). Similar collective excitations have been studied in literature devoted to other systems, e.g., anisotropic superconductors [52], which can

be consulted for details and references. Due to narrowness of regions B and C, some of the described features may be difficult to see in Fig. 4(a) and probably challenging to observe in experiments. For this reasons, we will mostly refrain from discussing regions B and C further.

One implication of Eq. (23) is that the HP^2 dispersion is widely tunable: the scaling law $q_n \propto d^{-1}$ provides a practical way to engineer a desired wavelength of the waveguide modes simply by tailoring the slab thickness d , as has been previously demonstrated using hBN slabs [20].

B. Surface plasmons

Examples of the collective mode spectra at finite doping are shown in Figs. 4(b) and 4(c) for G/hBN/G and Figs. 4(e) and 4(f) for Bi_2Se_3 . The spectra are dramatically different inside and outside the hyperbolic frequency bands. A key to understanding this difference is again the value of the momentum $k_1^z \simeq q \tan \theta(\omega)$. Outside the hyperbolic bands, it is almost purely imaginary, and so the collective excitations are exponentially confined to the surfaces of the slab. These surface modes are the Dirac plasmons introduced in Sec. I. Having in mind applications to near-field experiments, we are particularly interested in momenta q of the order of a few times 10^5 cm^{-1} , i.e., the region nearby the dashed lines $q = 0.0025$ nm^{-1} in Fig. 4. If ϵ_1 is real, there are at most two solutions of Eq. (22), one for $n = 0$ and the other for $n = 1$. However, the distinct $n = 1$ dispersion curves are visible only

in Figs. 4(b) and 4(c) for G/hBN/G and none of them is close enough to the range of q we are interested in. Therefore we focus on the $n = 0$ branch.

The shape of the plasmon dispersion curves in TI slabs and double-layer graphene systems was a subject of many previous theoretical studies [31,37,40,43,44] whose basic conclusions are reproduced by the following analysis. To the right of the dashed lines in Figs. 4(b) and 4(e) and for $d \sim 100$ nm, the dimensionless product $2k_F^z d = q\delta$ is typically large by absolute value and almost purely imaginary. This implies that the plasmons of the two interfaces are decoupled. Taking into account that $\epsilon_0 < \epsilon_s$ and $q_{\text{top}} = q_{\text{bot}}$ in Fig. 4, one can show that the dispersion of the $n = 0$ mode is controlled by the properties of the top interface. In the first approximation, this dispersion can be obtained setting $\phi_{\text{top}} \rightarrow -i\infty$, which yields

$$q_0 \approx \frac{\epsilon_0 + \epsilon_1}{2} q_{\text{top}}, \quad q_0 \gg |\delta|^{-1}. \quad (24)$$

For $\mu = 0$, momentum $q_{\text{top}} = q_{\text{top}}(q_0, \omega)$ is imaginary, cf. Eqs. (13) and (19). Hence, for real ϵ_1 , Eq. (24) has no real solutions: as already mentioned, undoped SSs do not support plasmons. Indeed, Figs. 4(a) and 4(d) contain no bright lines outside the hyperbolic bands. On the other hand, if $\mu \neq 0$, we can use Eq. (11) to transform Eq. (24) to Eq. (4), which predicts a parabolic dispersion curve $\omega \propto \sqrt{q}$ if ϵ_1 is constant. Such parabolas are seen in the upper halves of Figs. 4(b), 4(c), and 4(e), 4(f) although they appear rectilinear because of the restricted range of q .

As smaller momenta Eq. (24) no longer holds. The correct approximation for the $n = 0$ mode is obtained by setting the left-hand side of Eq. (22) to zero. This yields $\phi_{\text{top}} = -\phi_{\text{bot}}$ and

$$q_0 \simeq \frac{\epsilon_0 + \epsilon_s}{2} \frac{1}{q_{\text{top}}^{-1} + q_{\text{bot}}^{-1}} \simeq \frac{2}{N} \frac{\epsilon_0 + \epsilon_s}{e^2 |\mu|} (\hbar\omega)^2. \quad (25)$$

Thus both the low- q and high- q parts of the $n = 0$ dispersion curve are parabolic but with different curvatures. The crossover between these two parabolas occurs via a rapid increase of $\epsilon^\perp(\omega)$, and so, $\epsilon_1(\omega)$ at frequencies immediately above the hyperbolic bands. It takes place at $\omega > 1614 \text{ cm}^{-1}$ for G/hBN/G and $\omega > 163 \text{ cm}^{-1}$ for Bi₂Se₃, which generates the inflection points seen on the curves in, respectively, Figs. 4(b), 4(c) and 4(e), 4(f).

As indicated schematically in Fig. 2(a), at very large q the plasmon dispersion should have another inflection point. Using the more accurate Eq. (10) instead of Eq. (11), we find the following analytical result for the frequency of the $n = 0$ mode as a function of q :

$$\omega(q) \simeq v \frac{q + q_*}{\sqrt{1 + (2q_*/q)}}, \quad q_* = \frac{2e^2}{\hbar v} \frac{Nk_F}{\epsilon_0 + \epsilon_1}. \quad (26)$$

This equation predicts a crossover from the parabolic to the linear dispersion $\omega \simeq vq$ above $q = q_*$. However, this occurs far outside the plot range of Fig. 4.

Returning to Eq. (25), we notice that it does not contain the bulk permittivities. Hence, it should continue to hold for a range of ω inside the hyperbolic bands. A physical picture of this mode [“0(s)” in Fig. 2(c)] is in-phase oscillations of the charges of both Dirac fermion layers, i.e., the system behaving as a single 2D layer with the combined oscillator strength. As ω decreases further into the hyperbolic bands, the length

scale $|\delta|$ increases. The strength of the inequality $q_0|\delta| \ll 1$ and so the accuracy of Eq. (25) becomes progressively lower [in fact, Eq. (27) below gives a better approximation]. At $\omega = \omega_{\text{to}}$ for G/hBN/G and similarly, at $\omega = \omega_{\text{to},1}^\perp$ for Bi₂Se₃, this inequality is violated completely, which is consistent with the termination of these branches at $q = 0$ in Figs. 4(b) and (e). Similar analysis can be applied to Figs. 4(c) and 4(f) where d is twice smaller than in, respectively, Figs. 4(b) and 4(e). Because of that, the plasmon dispersion in the region $q|\delta| < 1$ is shifted to smaller q . The dispersions in the large- q regions are virtually unaffected since the stronger surface confinement of the plasmons makes them insensitive to d .

One qualitative difference between G/hBN/G and Bi₂Se₃ is the richer phonon spectrum of the latter. This leads to the avoided crossings of the plasmon branch with the dispersion lines of the HP² modes in regions B and C of Bi₂Se₃, cf. Figs. 4(b), 4(c) and 4(e), 4(f). The small shifts caused by those crossings are somewhat masked by the considerable linewidth of the $n = 0$ line due to relatively stronger phonon damping. In turn, higher electronic damping rate $\gamma_e \sim \omega_{\text{to},1}^\perp$ due to disorder scattering in Bi₂Se₃ effectively eliminates the plasmon excitations in the lower spectral region $\omega < \omega_{\text{to},1}^\perp$, see Figs. 4(e) and 4(f). Therefore we do not discuss it here.

C. Hybrid modes

From now on we turn to the subject of our primary interest, the hyperbolic collective modes of a doped TI. In this short section we address their dispersion law. Comparing Fig. 4(d) for $\mu = 0$ with Figs. 4(e) and 4(f) for $\mu > 0$, we observe significant shifts in the dispersion of the $n = 0$ mode in the upper half of the hyperbolic band $\omega_{\text{to},1}^\perp < \omega < \omega_{\text{to},1}^\perp$ of Bi₂Se₃. Similar shifts are seen in hBN near ω_{to} , cf. Fig. 4(a) with Figs. 4(b) and 4(c). These shifts result from hybridization of HP² and Dirac plasmons into combined HP³ waveguide modes. In general, calculation of these shifts requires solving Eq. (22) numerically. However, near the bottom of the hyperbolic band where these shifts become small, they can be also found analytically. Thus, Eq. (25) gets replaced by

$$q_0 \simeq \frac{\epsilon_0 + \epsilon_s}{\epsilon^\perp d + 2q_{\text{top}}^{-1} + 2q_{\text{bot}}^{-1}}, \quad |\epsilon_1| \gg \epsilon_0, \epsilon_s, \quad (27)$$

which shows explicitly that q_0 goes to zero as ω approaches $\omega_{\text{to},1}^\perp$ where ϵ^\perp sharply increases.

Unlike in Figs. 4(a) and 4(d), in Figs. 4(b), 4(c), 4(e), and 4(f), the higher-order $n > 1$ modes are more difficult to see because of their lower relative intensity compared to those of the plasmon $n = 0$ (and $n = 1$) modes. Nevertheless, these modes remain well defined (underdamped). Near the bottoms of the respective hyperbolic bands their momenta q_n still form an equidistant sequence except with a spacing

$$q_{n+1} - q_n \simeq \frac{2\pi}{l - \delta}, \quad (28)$$

which is modified compared to Eq. (23). This result can be obtained from Eq. (22) by approximating the finite differences such as $\phi_{\text{top}}(q_{n+1}) - \phi_{\text{top}}(q_n)$ by means of the derivative. Parameter l is defined by

$$l = -2 \frac{\partial \phi_{\text{top}}}{\partial q} - 2 \frac{\partial \phi_{\text{bot}}}{\partial q}. \quad (29)$$

The physical meaning of this quantity is clarified in the next section.

IV. GOOS-HÄNCHEN EFFECT

In this section, we consider the problem of the plasmon-polariton mixing from the point of view of real-space trajectories of the HP^2 excitations. The question we consider is how polariton wave packets propagate inside the slab and, in particular, how they reflect off its interfaces. As mentioned in Sec. I, for a given ω , the angle θ between the z axis and the group velocity \mathbf{v} vector of HP^2 s is nearly independent of q . Therefore, monochromatic HP^2 wave packets propagate as highly directional rays. Naively, one would then expect that the polariton rays should zigzag inside the slab returning to each interface periodically with the repeat distance of $2d |\tan \theta| = |\delta|$. Although such geometrical optics picture is adequate for insulating hyperbolic materials [53], it is not quite correct for TI with gapless doped SS. The geometrical optics neglects a lateral shift or displacement of the rays after each reflection [compare Figs. 5(a) and 5(b)], which is analogous to the GH effect of light. The GH effect was first discussed in the context of the total internal reflection of light. As explained below, it can be understood from two complementary points of view. In the wave picture, it originates from the momentum dependence of the reflection phase shift. In the particle picture, the GH effect is due to the quasi-classical tunneling (excitation of evanescent waves) along the interface. To define such a displacement one usually considers a wave packet with a

smooth envelope (for example, a Gaussian), in which case the displacement is the shift in the position of its maximum.

While the GH effect [45] was discovered measuring the reflection of light off an air-metal interface, the displacement \mathbf{l} of the reflected ray is a general wave phenomenon [46] that arises due to the dependence of the reflection phase shift ϕ on the lateral momentum $\mathbf{q} = (k^x, k^y)$. For example, the GH effect should also occur for surface plasmons [54]. The expression for \mathbf{l} has the form [55]

$$\mathbf{l} = -\Re e \frac{\partial \phi}{\partial \mathbf{q}}. \quad (30)$$

It seems to be another general rule that the momentum dependence of ϕ is significant only if the interface supports electromagnetic modes with either a large propagation length or a long decay length if such modes are evanescent. In the original photonic GH effect, this is the case under the conditions of the total internal reflection. The magnitude $|\mathbf{l}|$ of the GH displacement can be interpreted as the decay length of the evanescent transmitted wave. Alternatively, a large GH shift can occur if the interface supports surface plasmons or polaritons [56–58]. Experimental demonstration of the GH effect enhanced by surface plasmons of the air-metal interface has been reported [59].

Comparing Eqs. (29) and (30), we recognize the length scale l in the former as the sum of the GH shifts due to the top and the bottom interfaces. Therefore we conclude that the Dirac plasmons must act as the transient interface modes for the HP^2 rays bouncing inside the TI slab. Using Eqs. (20) and (30), and taking into account that $\Re e \epsilon_1 \ll \Im m \epsilon_1$, we find the GH shift at the top interface to be

$$l_{\text{top}} = \frac{4}{q_{\text{top}}} \frac{\Im m \epsilon_1}{(\epsilon_0 - \frac{2q}{q_{\text{top}}})^2 + |\epsilon_1|^2}. \quad (31)$$

A few comments on this result can be made. First, the GH shift is positive in our case, which means the displacement is in the same direction as the in-plane group velocity of the ray. Second, l_{top} depends on the permittivity of the environment. For example, at fixed q , it vanishes if ϵ_0 is very large. Conversely, for fixed ϵ_0 , the GH shift reaches its maximum

$$l_{\text{max}} = \frac{2}{\pi} \frac{\lambda_p \epsilon_0 \Im m \epsilon_1}{(\Re e \epsilon_1)^2 + (\Im m \epsilon_1)^2}, \quad \lambda_p \equiv \frac{2\pi}{\epsilon_0 q_{\text{top}}}, \quad (32)$$

at $q = \pi/\lambda_p$. Finally, l_{max} depends linearly on the characteristic size λ_p of the Dirac plasmon wavelength and inversely on the absolute value $|\epsilon_1| \approx \Im m \epsilon_1$ of the effective permittivity of the hyperbolic medium.

In Fig. 6, we show l_{max} for Bi_2Se_3 and G/hBN/G systems as a function of ω spanning their respective hyperbolic bands. The relative shift, l_{max}/λ_p , is greater in G/hBN/G because $|\epsilon_1|$ is smaller. Yet the absolute l_{max} at the same $\mu = 0.3$ eV is greater in Bi_2Se_3 (where it is ~ 200 nm) because it is hyperbolic at lower frequencies and λ_p is larger at smaller ω .

One possible setup for experimental detection of the GH effect in TI is shown in Fig. 5(c). It differs from Fig. 1(b) in the addition of a split gate between the TI slab and the substrate. If this gate is made of a good conductor with large permittivity, it would suppress the GH shift at the bottom surface. However,

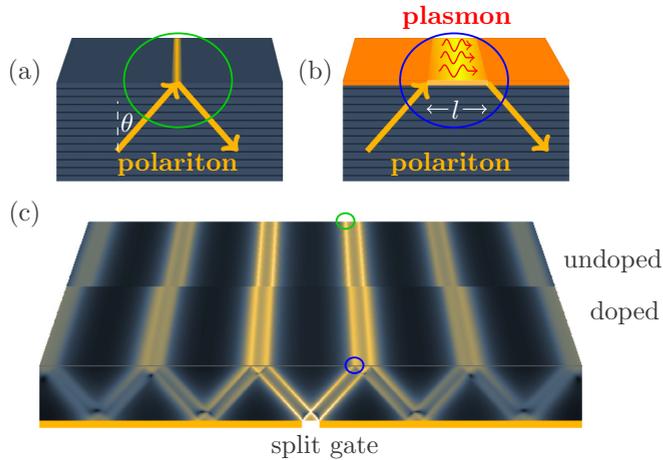


FIG. 5. (Color online) Polaritonic GH effect in TI slabs. (a) Schematics of the HP^2 ray reflection in the absence of the SS. (b) The same with the SS. The wavy lines symbolize virtual Dirac plasmons. The GH shift l is indicated. (c) The electric field distribution inside and/or at the upper surfaces of two slabs with equal $\delta = -2.2d$ but different doping. The lower (“doped”) and the upper (“undoped”) parts of the image are computed for $\lambda_p = a$ and 0, respectively. The split gate—a pair of metallic half-planes separated by a distance $2a$ —launches highly directional HP^2 rays that bounce inside the slabs creating periodic “hot stripes” at their upper surfaces. The period is larger in the “doped” slab. The two small circles, one in the undoped and one in doped part, are the representative locations of the HP^2 reflections. Their enlarged views are shown in, respectively, (a) and (b).

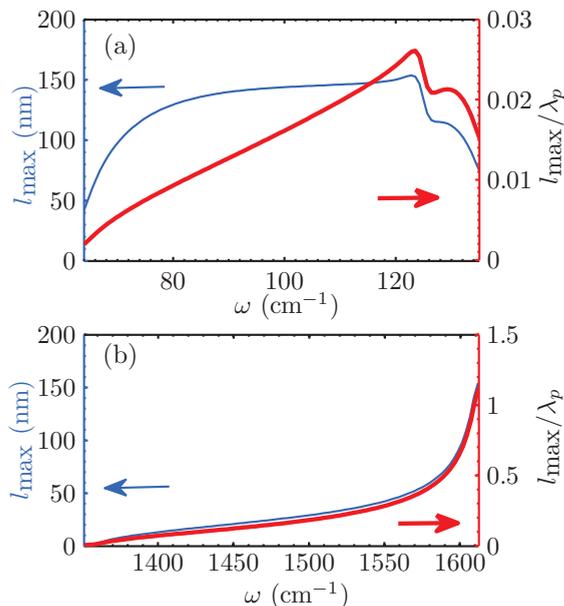


FIG. 6. (Color online) Maximum GH shift l_{\max} (in absolute units and as a fraction of λ_p) for (a) TI slab and (b) G/hBN/G structure with the same chemical potential $\mu = 0.3$ eV.

it would serve another useful purpose. Previously, it has been demonstrated [23] that in the presence of an external oscillating field, thin metallic disks or stripes can launch HP² in hBN. The split gate is to perform the same function here. The HP²s are preferentially emitted from the regions of highly concentrated field near the sharp metallic edges. We expect the rays to zigzag away from their launching points returning to the top surface with the period $l - \delta$, which is the sum of $-\delta \approx |\delta|$ due to the roundtrip inside the slab and $l = l_{\text{top}}$ due to the GH shift at the top surface. Since l depends q_{top} , which is controlled by doping, the GH effect can be detected by measuring the positions of the electric field maxima [“hot stripes” in Fig. 5(c)] as a function of μ in the experiment. Although l is quite small, the shifts accumulate after multiple reflections, which can facilitate their detection, as in the original work of Goos and Hänchen [45].

To model the response of the system shown in Fig. 5(c) quantitatively, we proceed as follows. We approximate the half-planes of the split gate by perfect conductors in the $z = 0$ plane with the edges at $x = \pm a$. Let $V(x, 0)$ be the scalar potential at $z = 0$ due to the external uniform field and all the charges induced on the gate. (Here and below the common factor $e^{-i\omega t}$ is omitted.) Let $\tilde{V}(k^x)$ be the Fourier transform of $V(x, 0)$. Using the notations for the reflection coefficients introduced in Sec. II, we express the potential $V(x, z)$ inside the slab $0 \leq z \leq d$ by the integral

$$V(x, z) = \int \frac{dk^x}{2\pi} \tilde{V}(k^x) t(k^x, z) e^{ik^x x}, \quad (33)$$

$$t(k^x, z) = \frac{e^{i|k^x|z \tan \theta} - r_{10}(k^x) e^{i|k^x| \tan \theta (2d-z)}}{1 - r_{10}(k^x) r_{12}(k^x) e^{i|k^x| \delta}}. \quad (34)$$

For a consistency check, we can consider the large $-x$ behavior of this inverse Fourier transform, which should be dictated by the poles of the integrand. These poles can be

recognized as the HP³ momenta q_n [Eq. (22)]. Since q_n form the equidistant sequence [Eq. (28)], their superposition should indeed create beats of period $l - \delta$, in agreement with our ray trajectories picture, Fig. 5(b).

Explicit calculation of $V(x, 0)$ requires a self-consistent solution of the Maxwell equations for our complicated multilayer system, which is computationally intensive. Fortunately, very similar results for $V(x, z)$ are obtained with little effort by approximating the true $V(x, 0)$ with the “bare” potential that would exist in the TI is removed, that is, if $d = \lambda_p = 0$. At distances less than c/ω from the gap in the gate, this bare potential has the simple analytical form,

$$V(x, 0) = \frac{V_0}{2} \times \begin{cases} +1, & x \leq -a, \\ -\frac{2}{\pi} \arcsin(x/a), & |x| < a, \\ -1, & x \geq a, \end{cases} \quad (35)$$

familiar from classical electrostatics. Its Fourier transform is given by

$$\tilde{V}(k^x) = \frac{iV_0}{k^x} J_0(k^x a), \quad (36)$$

where $J_0(x)$ is the Bessel function of the first kind and V_0 is potential difference between the two parts of the gate. The tangential electric field corresponding to this potential,

$$E_x = \frac{V_0}{\pi \sqrt{a^2 - x^2}}, \quad (37)$$

exhibits an inverse square-root divergence at the edges, which enables the localized HP² emission.

Carrying out the quadrature in Eq. (33) numerically, we have calculated the components and also the amplitude of the electric field $E = \sqrt{E_x^2 + E_z^2}$ over an interval of x a few $|\delta|$ in length and z varying from 0 to d . Our results for $E = E(x, z)$ for two doping levels, corresponding to $\lambda_p = 0$ (undoped SS) and $\lambda_p = a$ (doped SS) are illustrated by the false color plots in Fig. 5(c). These plots are superimposed on perspective projections of the two slabs (doped and undoped), which are placed next to each other for easy comparison. The remaining parameters of the calculations are $\delta = -2.2d$ and $a = 0.1d$. We see that a finite shift of the “hot stripes” at the top surface $z = d$ exists in the doped case. This seems to vindicate our intuition but actually the situation is a bit more subtle. The problem is that the momentum distribution of our source [Eq. (36)] is very different from what we assumed it to be in the beginning of our discussion of the GH effect. This distribution is not narrow and not centered at some finite k^x . Instead, it has positive and negative k^x harmonics of equal strength and a long power-law tail at $|k^x| \gg 1/a$. The reason why the GH shift persists in our case is the spatial separation of the k^x harmonics: due to the directionality of the HP² propagation, the stripes to the left (right) of the launching points are created predominantly by negative (positive) k_x . Since \mathbf{l} has the same direction as $\mathbf{q} = (k^x, 0)$, the stripes shift away from the origin on both sides of the y axis. A formal derivation of this result can be done by splitting the integral in Eq. (33) into the $k_x > 0$ and the $k_x < 0$ parts and identifying the relevant poles $k_x = q_n$ using contour integration methods.

From numerical experiments with different a , we found that the largest shift of the stripes is obtained for $a \sim \lambda_p$. This can be explained by arguing that the shift is maximized when the

characteristic $k^x \sim \pi/a$ contributing to the integral in Eq. (33) is close to the momentum π/λ_p at which $l = l_{\max}$ in Eq. (31).

Experimental detection of the “hot stripes” and their doping-dependent GH shift is possible via the s-SNOM imaging. This technique involves measuring the light scattered by the tip of an atomic force microscope brought to the sample and scanned along its surface [47,48]. Using clever signal processing methods, it is possible to isolate the genuine near-field component of this scattered light, which originates from conversion of evanescent electromagnetic waves emanating from the sample into free-space photons. In the proposed experiment, the evanescent waves are due to the HP^2 modes launched by the split gate. The spatial resolution of the s-SNOM imaging is set by the tip curvature radius R . For typical $R = 20\text{--}40$ nm, it is barely sufficient to observe the predicted GH shifts in hBN/G, Fig. 6(b). Nevertheless, detecting the cumulative shift after several stripe periods should be feasible. The prior success of s-SNOM imaging experiments of surface plasmons and polaritons in graphene and hBN structures [20,23–25,27,28,33,35] gives us a firm confidence in this approach. Note that if a doped graphene layer only partially covers the top surface of hBN, one literally gets the situation depicted in Fig. 5(c), where the doped and undoped regions are positioned side by side.

In the case of Bi_2Se_3 where the GH shift ~ 200 nm [Fig. 6(a)] is much larger, the spatial resolution of the s-SNOM is even less of an issue. The main obstacle is the scant availability of suitable THz sources. We are optimistic that in a near future this problem can be overcome as well.

V. SUMMARY AND OUTLOOK

Recent experiments [8,14] have shown that coupling between Dirac plasmons and bulk phonons of bismuth-based TIs should be strong. In this paper, we have studied this interaction taking into account the anisotropic phonon spectrum of such TIs. We have predicted that a TI slab can act as a tunable waveguide for phonon polaritons, with the doping of the surface states being the tuning parameter. In addition to the change in dispersion, the phonon-plasmon coupling can cause measurable real-space shifts of the polariton rays. Similar phenomena have been recently studied in artificial structures made by stacking graphene layers on top of hBN. The present work indicates that the TIs are a promising alternative platform for realizing highly tunable, strongly confined, low-loss electromagnetic modes in a *natural* material. Additionally, while hBN/G waveguides operate in midinfrared frequencies, Bi_2Se_3 and similar compounds extend the same functionality to the technologically important THz domain.

We envision several directions for further work in this field. One is to attempt a multisource coherent control of polariton emission and propagation using ultrafast laser pulses. A variety of such techniques has been developed [60] in the context of THz polaritonics of $LiNbO_3$ and $LiTiO_3$. (Incidentally, a theoretical proposal [61] of integrating graphene into such materials would lead to polariton waveguides similar in functionality and perhaps also tunability to those studied in the present work.) Another intriguing direction is to explore oscillating spin currents, which were predicted to accompany charge density currents produced by Dirac plasmons [32].

It may be also interesting to study the effect of optical hyperbolicity [19] on the high-energy bulk plasmons of the TIs [62,63]. Finally, it may be worthwhile to investigate new applications that can be enabled by tunable hyperbolic polaritons. Harnessing such types of modes for hyperlensing [64–66] or focusing [23,24] has been widely discussed. The present work shows that the GH effect and its dependence on doping and dielectric environment of the TI can be another avenue for applications, for example, THz chemical sensing or characterization of spatially inhomogeneous TI samples. We hope our work can stimulate these and other future studies.

ACKNOWLEDGMENTS

This work is supported by the University of California Office of the President and by the ONR.

APPENDIX: NEAR-FIELD SPECTRA

A fully realistic modeling of the s-SNOM imaging experiments proposed in Sec. IV is an unwieldy task requiring a repeated solution of the Maxwell equations for a system with complicated material properties, a hierarchy of widely different length scales, and no special symmetries. In this Appendix, we present some results of less ambitious calculations that simulate a simpler structure depicted in Fig. 1(b). Although no split gate is present in this structure, the measured signal is still expected to reveal characteristics of the collective modes. In this case, these modes are excited by the sharp tip itself. Hence the tip plays the role of both the launcher and the detector of the HP^3 modes. Unfortunately, this implies that only the local response can be measured, which is a superposition of responses due to a distribution of momenta up to $q_t \sim 1/R$ rather than one specific q .

We assume that the TI slab and the substrate are infinite and uniform in x and y coordinates, so that the imaging capability of the s-SNOM is irrelevant. Instead, the quantity of interest is the frequency dependence of the measured near-field signal $s(\omega)$. A few more explanations about our calculational scheme are in order. We model the tip as a metallic spheroid with the curvature radius $R = 40$ nm and total length 720 nm. We use the quasistatic approximation but include the radiative corrections included perturbatively. This model [67,68] has been successful for simulating many recent s-SNOM experiments, and should be especially suitable in the THz domain where no antenna resonances or other strong retardation effects [69] should appear. Our calculations incorporate the so-called far-field factors [67–69], which are expressed in terms of $r_P(q, \omega)$ at $q \sim \omega/c$. These factors account for the fact that the incident wave is originally created by a far-field source and the scattered wave is ultimately measured by a far-field detector. Finally, what we compute is not the full scattering amplitude s but its third harmonic s_3 , which is what experimentalists typically report. The idea is that in the experiment the tip is made to oscillate at some low frequency Ω , so that s is periodic with this fundamental tapping frequency. The third Fourier harmonic of s , which is s_3 , gives a good representation of the genuine near-field signal.

Naively, one can think of $s_3(\omega)$ as a weighted average of the surface reflectivity $r_P(q, \omega)$ over q . The weighting

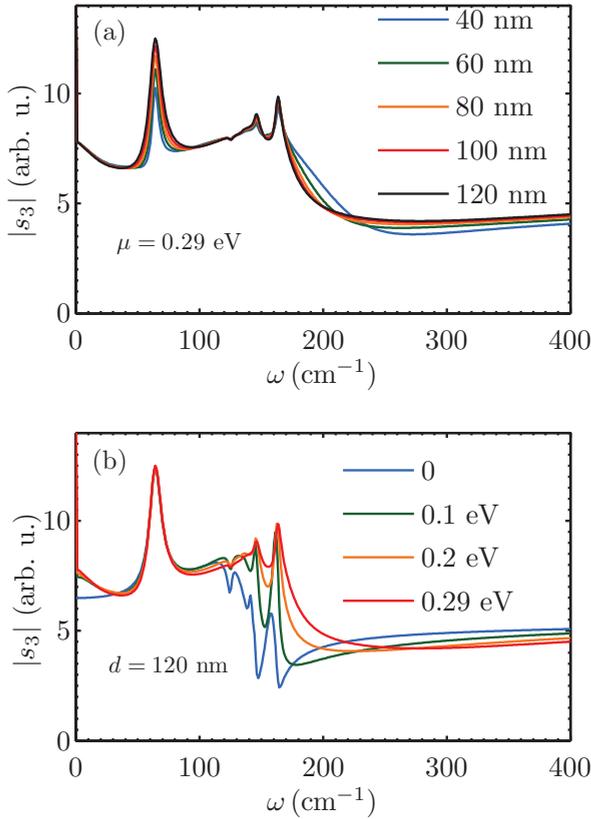


FIG. 7. (Color online) Simulation of the s-SNOM signal s_3 for Bi_2Se_3 slabs on a substrate with $\epsilon_s = 10$. (a) Fixed $\mu = 0.29$ eV and different d . (b) Fixed $d = 120$ nm and different μ .

function has a broad maximum near $q = q_t$, which in this case is equal to $q_t = 0.025$ nm $^{-1}$ [the dashed lines in Fig. 4]. The presence of strong maxima of $\Im m r_P$ due to collective modes with momenta $q \lesssim q_t$ tends to enhance $s_3(\omega)$. In a more rigorous picture [68], the maxima of $s_3(\omega)$ correspond not to the resonances of the sample alone but to those of the coupled tip-sample system. The coupling can decrease the resonance frequencies by as much as [67–69] 10–20 cm $^{-1}$ compared to those seen in $\Im m r_P$ maps.

Our results for Bi_2Se_3 slabs of various thickness d and chemical potential μ are shown in Fig. 7. Pairs of distinct peaks as well as smaller additional features are readily seen. In each trace, the stronger and sharper peak is located close to $\omega_{10,1}^\perp = 64$ cm $^{-1}$. The height of this peak decreases as d decreases [Fig. 7(a)]. However, its position is independent of d [Fig. 7(a)] or μ [Fig. 7(b)], which suggests that it is not related to the dispersive HP 3 modes. Indeed, we have verified that this prominent peak is almost entirely due to the far-field factor $|1 + r_P|^2$, which has a narrow maximum at $\omega_{10,1}^\perp$ where $r_P \approx 1$.

Each of the doped samples also produces smaller peaks in $s_3(\omega)$, of which the most prominent ones are those located near $\omega = 146$ and 163 cm $^{-1}$, the upper boundaries of regions B and C of Fig. 3. The position and especially the strength of the peaks is μ -dependent. As μ increases, the peaks grow in height and gradually shift to higher frequencies, see Fig. 7(b). These peaks are due to the surface modes: the $n = 1$ mode of region B and the $n = 0$ mode just above region C, see Figs. 4(d) and 4(e). The increase of the peak heights with μ can be qualitatively explained by the increase of the absolute value of r_P . The shift in position is unfortunately more difficult to interpret without a better understanding of the effective weighting function that relates $\Im m r_P(q, \omega)$ to $s_3(\omega)$.

While $\mu > 0$ traces are due to combined action of plasmons and phonon polaritons, the $\mu = 0$ one is expected to reveal the phonon-polariton response. Interestingly, that trace exhibits a sharp dip at $\omega = 163$ cm $^{-1}$, see Fig. 7(b). We have checked that this dip is not caused by the far-field factor. However, its relation to the HP 2 modes of Fig. 4(d) is not obvious to us.

The thickness dependence of s_3 is illustrated in Fig. 7(a). As one can see, the near-field peak at $\omega = 163$ cm $^{-1}$ has a broad high-frequency side, which systematically expands as d decreases. This trend reflects the blue shift of the $n = 0$ mode dispersion in thinner slabs, compare Figs. 4(e) and 4(f).

Overall, our simulations predict that the near-field response of Bi_2Se_3 slabs should exhibit systematic spectral changes with doping and thickness that are measurable by the s-SNOM. Such experiments may provide insights into properties of tunable HP 3 modes of these novel systems.

-
- [1] M. Z. Hasan and C. L. Kane, *Rev. Mod. Phys.* **82**, 3045 (2010).
- [2] X.-L. Qi and S.-C. Zhang, *Rev. Mod. Phys.* **83**, 1057 (2011).
- [3] W. Richter, H. Köhler, and C. R. Becker, *Phys. Stat. Sol. (b)* **84**, 619 (1977).
- [4] A. D. LaForge, A. Frenzel, B. C. Pursley, T. Lin, X. Liu, J. Shi, and D. N. Basov, *Phys. Rev. B* **81**, 125120 (2010).
- [5] W. Cheng and S.-F. Ren, *Phys. Rev. B* **83**, 094301 (2011).
- [6] A. Akrap, M. Tran, A. Ubaldini, J. Teyssier, E. Giannini, D. van der Marel, P. Lerch, and C. C. Homes, *Phys. Rev. B* **86**, 235207 (2012).
- [7] P. Di Pietro, F. M. Vitucci, D. Nicoletti, L. Baldassarre, P. Calvani, R. Cava, Y. S. Hor, U. Schade, and S. Lupi, *Phys. Rev. B* **86**, 045439 (2012).
- [8] P. Di Pietro, M. Ortolani, O. Limaj, A. Di Gaspare, V. Giliberti, F. Giorgianni, M. Brahlek, N. Bansal, N. Koirala, S. Oh, P. Calvani, and S. Lupi, *Nat. Nano.* **8**, 556 (2013).
- [9] L. Wu, M. Brahlek, R. Valdés Aguilar, A. V. Stier, C. M. Morris, Y. Lubashevsky, L. S. Bilbro, N. Bansal, S. Oh, and N. P. Armitage, *Nat. Phys.* **9**, 410 (2013).
- [10] K. W. Post, B. C. Chapler, L. He, X. Kou, K. L. Wang, and D. N. Basov, *Phys. Rev. B* **88**, 075121 (2013).
- [11] B. C. Chapler, K. W. Post, A. R. Richardella, J. S. Lee, J. Tao, N. Samarth, and D. N. Basov, *Phys. Rev. B* **89**, 235308 (2014).
- [12] A. A. Reijnders, Y. Tian, L. J. Sandilands, G. Pohl, I. D. Kivlichan, S. Y. Frank Zhao, S. Jia, M. E. Charles, R. J. Cava, N. Alidoust, S. Xu, M. Neupane, M. Z. Hasan, X. Wang, S. W. Cheong, and K. S. Burch, *Phys. Rev. B* **89**, 075138 (2014).
- [13] M. Autore, H. Engelkamp, F. D’Apuzzo, A. D. Gaspare, P. D. Pietro, I. L. Vecchio, M. Brahlek, N. Koirala, S. Oh, and S. Lupi, *ACS Photon.* **2**, 1231 (2015).
- [14] M. Autore, F. D’Apuzzo, A. Di Gaspare, V. Giliberti, O. Limaj, P. Roy, M. Brahlek, N. Koirala, S. Oh, F. J. García de Abajo, and S. Lupi, *Adv. Opt. Mat.* **3**, 1257 (2015).

- [15] K. W. Post, B. C. Chapler, M. K. Liu, J. S. Wu, H. T. Stinson, M. D. Goldflam, A. R. Richardella, J. S. Lee, A. A. Reijnders, K. S. Burch, M. M. Fogler, N. Samarth, and D. N. Basov, *Phys. Rev. Lett.* **115**, 116804 (2015).
- [16] Y. Guo, C. L. Newman, W. Cortes, and Z. Jacob, *Adv. OptoElectron.* **2012**, 452502 (2012).
- [17] A. Poddubny, I. Iorsh, P. Belov, and Y. Kivshar, *Nat. Photon.* **7**, 948 (2013).
- [18] J. Sun, N. M. Litchinitser, and J. Zhou, *ACS Photon.* **1**, 293 (2014).
- [19] M. Esslinger, R. Vogelgesang, N. Talebi, W. Khunsin, P. Gehring, S. de Zuani, B. Gompf, and K. Kern, *ACS Photon.* **1**, 1285 (2014).
- [20] S. Dai, Z. Fei, Q. Ma, A. S. Rodin, M. Wagner, A. S. McLeod, M. K. Liu, W. Gannett, W. Regan, K. Watanabe, T. Taniguchi, M. Thiemens, G. Dominguez, A. H. Castro Neto, A. Zettl, F. Keilmann, P. Jarillo-Herrero, M. M. Fogler, and D. N. Basov, *Science* **343**, 1125 (2014).
- [21] Z. Jacob, *Nat. Mater.* **13**, 1081 (2014).
- [22] J. D. Caldwell, A. V. Kretinin, Y. Chen, V. Giannini, M. M. Fogler, Y. Francescato, C. T. Ellis, J. G. Tischler, C. R. Woods, A. J. Giles, M. Hong, K. Watanabe, T. Taniguchi, S. A. Maier, and K. S. Novoselov, *Nat. Commun.* **5**, 5221 (2014).
- [23] S. Dai, Q. Ma, T. Andersen, A. S. McLeod, Z. Fei, M. K. Liu, M. Wagner, K. Watanabe, T. Taniguchi, M. Thiemens, F. Keilmann, P. Jarillo-Herrero, M. M. Fogler, and D. N. Basov, *Nat. Commun.* **6**, 6963 (2015).
- [24] P. Li, M. Lewin, A. V. Kretinin, J. D. Caldwell, K. S. Novoselov, T. Taniguchi, K. Watanabe, F. Gaussmann, and T. Taubner, *Nat. Commun.* **6**, 7507 (2015).
- [25] Z. Shi, H. A. Bechtel, S. Berweger, Y. Sun, B. Zeng, C. Jin, H. Chang, M. C. Martin, M. B. Raschke, and F. Wang, *ACS Photon.* **2**, 790 (2015).
- [26] V. W. Brar, M. S. Jang, M. Sherrott, J. J. Lopez, and H. A. Atwater, *Nano Lett.* **13**, 2541 (2013).
- [27] S. Dai, Q. Ma, M. K. Liu, T. Andersen, Z. Fei, M. D. Goldflam, M. Wagner, K. Watanabe, T. Taniguchi, M. Thiemens, F. Keilmann, G. C. A. M. Janssen, S.-E. Zhu, P. Jarillo-Herrero, M. M. Fogler, and D. N. Basov, *Nat. Nanotechnol.* **10**, 682 (2015).
- [28] G. X. Ni, H. Wang, J. S. Wu, Z. Fei, M. D. Goldflam, F. Keilmann, B. Özyilmaz, A. H. Castro Neto, X. M. Xie, M. M. Fogler, and D. N. Basov, *Nat. Mater.* (to be published).
- [29] A. Tomadin, F. Guinea, and M. Polini, *Phys. Rev. B* **90**, 161406 (2014).
- [30] A. H. Castro Neto, F. Guinea, N. M. R. Peres, K. S. Novoselov, and A. K. Geim, *Rev. Mod. Phys.* **81**, 109 (2009).
- [31] E. H. Hwang and S. Das Sarma, *Phys. Rev. B* **80**, 205405 (2009).
- [32] S. Raghu, S. B. Chung, X.-L. Qi, and S.-C. Zhang, *Phys. Rev. Lett.* **104**, 116401 (2010).
- [33] Z. Fei, G. O. Andreev, W. Bao, L. M. Zhang, A. S. McLeod, C. Wang, M. K. Stewart, Z. Zhao, G. Dominguez, M. Thiemens, M. M. Fogler, M. J. Tauber, A. H. Castro-Neto, C. N. Lau, F. Keilmann, and D. N. Basov, *Nano Lett.* **11**, 4701 (2011).
- [34] Z. Fei, A. S. Rodin, G. O. Andreev, W. Bao, A. S. McLeod, M. Wagner, L. M. Zhang, Z. Zhao, M. Thiemens, G. Dominguez, M. M. Fogler, A. H. Castro Neto, C. N. Lau, F. Keilmann, and D. N. Basov, *Nature (London)* **487**, 82 (2012).
- [35] J. Chen, M. Badioli, P. Alonso-González, S. Thongrattanasiri, F. Huth, J. Osmond, M. Spasenović, A. Centeno, A. Pesquera, P. Godignon, A. Z. Elorza, N. Camara, F. J. García de Abajo, R. Hillenbrand, and F. H. L. Koppens, *Nature (London)* **487**, 77 (2012).
- [36] A. N. Grigorenko, M. Polini, and K. S. Novoselov, *Nat. Photon.* **6**, 749 (2012).
- [37] R. E. V. Profumo, R. Asgari, M. Polini, and A. H. MacDonald, *Phys. Rev. B* **85**, 085443 (2012).
- [38] F. J. García de Abajo, *ACS Photon.* **1**, 135 (2014).
- [39] D. N. Basov, M. M. Fogler, A. Lanzara, F. Wang, and Y. Zhang, *Rev. Mod. Phys.* **86**, 959 (2014).
- [40] T. Stauber, G. Gómez-Santos, and L. Brey, *Phys. Rev. B* **88**, 205427 (2013).
- [41] R. Schütky, C. Ertler, A. Trügler, and U. Hohenester, *Phys. Rev. B* **88**, 195311 (2013).
- [42] J. Qi, H. Liu, and X. C. Xie, *Phys. Rev. B* **89**, 155420 (2014).
- [43] M. Li, Z. Dai, W. Cui, Z. Wang, F. Katmis, J. Wang, P. Le, L. Wu, and Y. Zhu, *Phys. Rev. B* **89**, 235432 (2014).
- [44] T. Stauber, *J. Phys.: Condens. Matter* **26**, 123201 (2014).
- [45] F. Goos and H. Hänchen, *Ann. Phys.* **436**, 333 (1947).
- [46] K. Y. Bliokh and A. Aiello, *J. Opt.* **15**, 014001 (2013).
- [47] F. Keilmann and R. Hillenbrand, *Phil. Trans. Roy. Soc. London, Ser. A* **362**, 787 (2004).
- [48] J. M. Atkin, S. Berweger, A. C. Jones, and M. B. Raschke, *Adv. Phys.* **61**, 745 (2012).
- [49] B. Wunsch, T. Stauber, F. Sols, and F. Guinea, *New J. Phys.* **8**, 318 (2006).
- [50] E. H. Hwang and S. Das Sarma, *Phys. Rev. B* **75**, 205418 (2007).
- [51] J. P. F. LeBlanc and J. P. Carbotte, *Phys. Rev. B* **89**, 035419 (2014).
- [52] H. T. Stinson, J. S. Wu, B. Y. Jiang, Z. Fei, A. S. Rodin, B. C. Chapler, A. S. McLeod, A. Castro Neto, Y. S. Lee, M. M. Fogler, and D. N. Basov, *Phys. Rev. B* **90**, 014502 (2014).
- [53] Z. Sun, A. Gutiérrez-Rubio, D. N. Basov, and M. M. Fogler, *Nano Lett.* **15**, 4455 (2015).
- [54] F. Huerkamp, T. A. Leskova, A. A. Maradudin, and B. Baumeier, *Opt. Express* **19**, 15483 (2011).
- [55] K. Artmann, *Ann. Phys.* **437**, 87 (1948).
- [56] T. Tamir and A. Oliner, *Proc. IEEE* **51**, 317 (1963).
- [57] T. Tamir and H. L. Bertoni, *J. Opt. Soc. Am.* **61**, 1397 (1971).
- [58] S. L. Chuang, *J. Opt. Soc. Am. A* **3**, 593 (1986).
- [59] X. Yin, L. Hesselink, Z. Liu, N. Fang, and X. Zhang, *Appl. Phys. Lett.* **85**, 372 (2004).
- [60] T. Feurer, N. S. Stoyanov, D. W. Ward, J. C. Vaughan, E. R. Statz, and K. A. Nelson, *Ann. Rev. Mater. Res.* **37**, 317 (2007).
- [61] D. Jin, A. Kumar, K. Hung Fung, J. Xu, and N. X. Fang, *Appl. Phys. Lett.* **102**, 201118 (2013).
- [62] J. J. Cha, K. J. Koski, K. C. Y. Huang, K. X. Wang, W. Luo, D. Kong, Z. Yu, S. Fan, M. L. Brongersma, and Y. Cui, *Nano Lett.* **13**, 5913 (2013).

- [63] J.-Y. Ou, J.-K. So, G. Adamo, A. Sulaev, L. Wang, and N. I. Zheludev, *Nat. Commun.* **5**, 5139 (2014).
- [64] Z. Jacob, L. V. Alekseyev, and E. Narimanov, *Opt. Express* **14**, 8247 (2006).
- [65] A. Salandrino and N. Engheta, *Phys. Rev. B* **74**, 075103 (2006).
- [66] Z. Liu, H. Lee, Y. Xiong, C. Sun, and X. Zhang, *Science* **315**, 1686 (2007).
- [67] L. M. Zhang, G. O. Andreev, Z. Fei, A. S. McLeod, G. Dominguez, M. Thiemens, A. H. Castro-Neto, D. N. Basov, and M. M. Fogler, *Phys. Rev. B* **85**, 075419 (2012).
- [68] B.-Y. Jiang, L. M. Zhang, A. H. Castro Neto, D. N. Basov, and M. M. Fogler, [arXiv:1503.00221](https://arxiv.org/abs/1503.00221).
- [69] A. S. McLeod, P. Kelly, M. D. Goldflam, Z. Gainsforth, A. J. Westphal, G. Dominguez, M. H. Thiemens, M. M. Fogler, and D. N. Basov, *Phys. Rev. B* **90**, 085136 (2014).