

Localization and Delocalization Errors in Density Functional Theory and Implications for Band-Gap Prediction

Paula Mori-Sánchez, Aron J. Cohen, and Weitao Yang

Department of Chemistry, Duke University, Durham, North Carolina 27708, USA

(Received 27 August 2007; published 7 April 2008)

The band-gap problem and other systematic failures of approximate exchange-correlation functionals are explained from an analysis of total energy for fractional charges. The deviation from the correct intrinsic linear behavior in finite systems leads to delocalization and localization errors in large and bulk systems. Functionals whose energy is convex for fractional charges such as the local density approximation display an incorrect apparent linearity in the bulk limit, due to the delocalization error. Concave functionals also have an incorrect apparent linearity in the bulk calculation, due to the localization error and imposed symmetry. This resolves an apparent paradox and identifies the physical nature of the error to be addressed to obtain accurate band gaps from density functional theory.

DOI: 10.1103/PhysRevLett.100.146401

PACS numbers: 71.10.-w, 71.15.Mb

Accurate band-gap prediction is critical for applications in condensed matter and nanotechnology. The theory of the band gap in density functional theory (DFT) was developed in the 1980s [1–4], but practical challenges remain. Our recent work [5] shows that, in principle, it is possible to obtain the correct band gap from practical DFT calculations and it is demonstrated for finite systems.

The challenge is in the approximate exchange-correlation functionals [6–9] which, despite the success in a wide range of applications, still suffer from systematic problems in describing charge-transfer processes, excitation energies in molecules, response properties in solids, electron transport, and the band gaps of semiconductors. Previous understanding has focused on self-interaction error and the Kohn-Sham (KS) eigenvalues, but other work [3–5, 10–12] relates these problems, more usefully, to the incorrect description of systems with fractional charges.

In this Letter, we will resolve an apparent paradox on the linearity of $E(N)$ for bulk systems for any approximate functional and provide insight on the physical basis underlying the systematic errors of functionals in large or extended systems, and its implication in the calculation of the band gap and other properties.

The fundamental gap of an N -electron semiconductor can be written as energy differences from integer points

$$\begin{aligned} E_{\text{gap}}^{\text{integer}} &= \{E(N-1) - E(N)\} - \{E(N) - E(N+1)\} \\ &= I - A, \end{aligned} \quad (1)$$

or as a difference of derivatives at N

$$E_{\text{gap}}^{\text{deriv}} = \left\{ \frac{\partial E}{\partial N} \Big|_{N+\delta} - \frac{\partial E}{\partial N} \Big|_{N-\delta} \right\}, \quad (2)$$

where $E_{\text{gap}}^{\text{integer}} = E_{\text{gap}}^{\text{deriv}}$ only if the total energy is a straight line between the integers, which is the case for the exact functional as shown by Perdew *et al.* based on grand

canonical ensembles [13] and later by Yang *et al.* based on pure states [14] (related to Ref. [3]).

The expression for the derivatives on the right-hand side of Eq. (2) is different for different types of exchange-correlation functionals [5]. In the case where E_{xc} is an explicit functional of ρ , such as local density approximation (LDA) or generalized gradient approximation (GGA), then

$$E_{\text{gap}}^{\text{deriv}} = \epsilon_{\text{gap}}^{\text{KS}} = \epsilon_{\text{LUMO}}^{\text{KS}} - \epsilon_{\text{HOMO}}^{\text{KS}} \quad (3)$$

and $E_{\text{gap}}^{\text{deriv}}$ is just obtained from the KS eigenvalues of the highest occupied molecular orbital (HOMO) and lowest unoccupied molecular orbital (LUMO). The detailed expressions of $\frac{\partial E}{\partial N} \Big|_{N \pm \delta}$ for other forms of exchange-correlation functionals have been derived recently [5]. When E_{xc} has an explicit dependence on the KS orbitals, $E_{xc}[\phi_i]$, then

$$E_{\text{gap}}^{\text{deriv}} = \epsilon_{\text{gap}}^{\text{KS}} + \left\{ \frac{\partial E_{xc}}{\partial N} \Big|_{N+\delta} - \frac{\partial E_{xc}}{\partial N} \Big|_{N-\delta} \right\}. \quad (4)$$

Thus, for such an orbital functional, the KS eigenvalues from an optimized effective potential [15] calculation are no longer the derivatives of the energy expression. However, the derivatives are exactly the eigenvalues in a generalized Kohn-Sham (GKS) calculation [e.g., Hartree-Fock (HF) calculations in the case of exact exchange [2,3]]:

$$E_{\text{gap}}^{\text{deriv}} = \epsilon_{\text{gap}}^{\text{GKS}}. \quad (5)$$

We see that the second term on the right-hand side of Eq. (4), which is labeled the derivative discontinuity Δ_{xc} [1,2], is essentially the difference between an optimized effective potential and GKS calculation [5]. However, it does not contain all the error of the band gap in an LDA or GGA calculation for which $\Delta_{xc} = 0$.

The error in the band-gap calculation using approximate functionals is in the following:

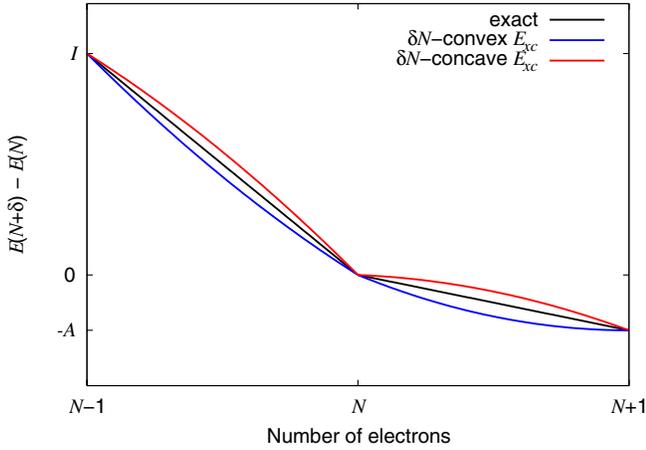


FIG. 1 (color online). Energy versus number of electron curve for a finite system with three hypothetical functionals with exact straight line behavior, δN -convex behavior, and δN -concave behavior all which gave the same I and A .

$$E_{\text{gap}}^{\text{integer}} = E_{\text{gap}}^{\text{deriv}} + \Delta_{\text{straight}}; \quad (6)$$

thus

$$E_{\text{gap}}^{\text{integer}} = \epsilon_{\text{gap}}^{\text{KS}} + \Delta, \quad \Delta = \Delta_{xc} + \Delta_{\text{straight}}. \quad (7)$$

Here Δ_{straight} , the difference between the gap from finite difference and the derivatives, accounts for the fact that an approximate functional may not have the correct straight line behavior between the integers. It is the consideration of this term, completely missing in the literature, that is the key in understanding the band gap. We have seen that Δ_{straight} has an important effect in the energy gap in molecules [5]. The focus of this Letter is to explore its implications for extended systems.

Now we present an apparent paradox for consideration of periodic or bulk systems: It has been argued in extended systems that the behavior of the total energy as a function of addition of an electron must be a straight line [1,2]; therefore, $\Delta_{\text{straight}} = 0$ for all functionals. How can it be opposite to the case of the finite systems [5,10], where the error is accounted for by the fact that $\Delta_{\text{straight}} \neq 0$? Furthermore, for LDA, $\Delta_{xc} = 0$, meaning that $\Delta = 0$. Then, what is the origin of the systematic error of LDA? It is well known [16] that the calculation of semiconductors and wide gap insulators with a local functional (such as LDA) has large systematic errors, underpredicting the band gap by up to several eV. We will resolve this apparent paradox by understanding more about the fractional charge behavior of energy functionals.

The simplest periodic system to consider is any crystal in the infinite lattice constant limit (a limit previously considered [4]). We start with a basic unit whose $E(N)$ curve for different functionals is as in Fig. 1: the exact functional with correct straight line behavior, one with incorrect convex behavior for fractional charges (δN con-

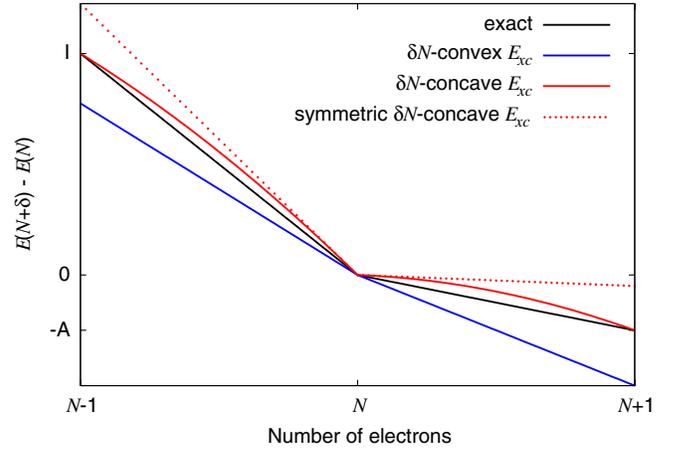


FIG. 2 (color online). The same as Fig. 1 except taken to the bulk limit. An additional curve is shown where crystal symmetry is imposed on a δN -concave functional calculation.

vex) such as LDA, and another with incorrect concave behavior (δN concave) such as HF. Next consider adding an electron to more than one unit. If we first take the case of two units, adding an electron to a δN -convex functional leads to half an electron on each unit, as it is much lower than the energy of two units, with N and $N + 1$ electrons, respectively. In this manner, it is clear that with M units

$$ME\left(N + \frac{1}{M}\right) < (M - 1)E(N) + E(N + 1). \quad (8)$$

As the number of units $M \rightarrow \infty$, the added δ electron delocalizes on to all the units such that the energy approaches the initial slope of the δN -convex curve of one unit as in Fig. 2:

$$E(MN + \delta) = ME\left(N + \frac{\delta}{M}\right) \rightarrow ME(N) + \delta \left. \frac{\partial E}{\partial N} \right|_{N+\delta}. \quad (9)$$

In this way a functional like LDA, which is δN convex for small molecules, will have an apparent linearity in large or periodic systems. It is, however, a quantitatively incorrect straight line, with the energy at the $N + 1$ integer point dictated by the fractional charge error of the functional. We can distinguish this from the correct behavior of the exact functional, which has intrinsic linearity with correct integer points for all M , whereas δN -convex functionals are only linear in the limit $M \rightarrow \infty$ with incorrect integer points. This is the delocalization error of δN -convex functionals.

The situation is reversed for δN -concave functionals, where the energy is always lower when the electron remains on one unit, even as M increases. In fact the E versus N curve is the same for all M . Delocalization actually raises the energy, as the initial slope points to above the integer points such that the inequality in Eq. (8) is reversed. For δN -concave functionals we will find $\Delta_{\text{straight}} \neq 0$ if we

carry out an energy minimized calculation at $N + 1$ in an infinite system. This is the localization error in δN -concave functionals. However, if periodic symmetry is imposed in the calculation, then the additional electron is delocalized through the entire crystal as required by the translational symmetry, a straight line will be seen that follows the initial slope of the E versus N curve. This delocalized state also has the energy of Eq. (9), which is a much higher energy, imposed by the symmetry. All the above arguments apply both to the addition of an electron and the addition of a hole (removal of electron).

In periodic band-structure calculations, we see the maximum effects of the localization and delocalization error: for both δN -convex and δN -concave functionals, we will have the apparent linearity for fractional charges, with the straight lines following the initial derivatives and leading to too low band gap for δN -convex functionals and too high band gap for δN -concave functionals. This explains the errors in the band-gap prediction.

In finite system calculations, delocalization error increases with system size until the apparent linearity appears. But localization error stabilizes at a certain system size, when the spatial extent of the added electron saturates. This also offers guidance on calculations of band gap for finite systems: The approach to calculate the band gap by explicit calculation of I and A from subtraction and addition of electron to finite neutral species, which works well with δN -convex functionals for small molecules [5], will not work for larger molecules, because delocalization error increases leading to the incorrect nature of the $N - 1$ and $N + 1$ points. However, it may work for δN -concave functionals, which do not suffer from the apparent linearity problem (without translational symmetry) and may give meaningful integer points. This issue is also slightly clouded by the fact that HF theory may not give a reasonable energy for the localized electron or even the right amount of localization. Hence it will give an additional error to the integer points, but this has a different physical basis.

The idea of delocalization error introduced here is connected to many-electron self-interaction error [10,12]. The poor performance of DFT calculations on the band gap was previously related to the self-interaction error and now we clearly relate this to the localization or delocalization of electrons. Thus, we believe that the terms localization and delocalization error capture the physics of the problem in a more useful manner than self-interaction error.

The discussion until now has focused on the energy differences and derivatives associated with the band gap as it shows very clearly the basic errors of approximate functionals. However, these errors of the functionals have much wider implications. We can see the differing behavior of the two types of functionals: δN -convex (or LDA-type) functionals tend to delocalize electrons and δN -concave (or HF-type) functionals tend to localize elec-

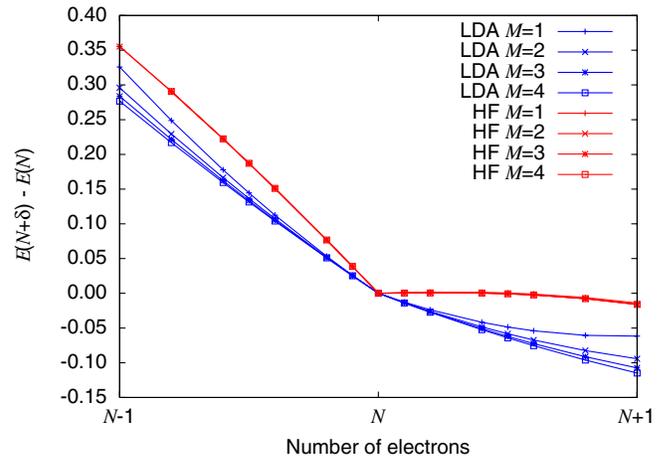


FIG. 3 (color online). Energy (in E_h) versus numbers of electrons for $(\text{H}_{16})_M$.

trons. So the nature of the electron density distribution, and the description of the HOMO and LUMO, are affected by the functional rather than being solely determined by the underlying physics of the material. These are the delocalization and localization errors of approximate functionals which are responsible for many of the errors of calculated properties in DFT which involve a change in localization, such as polarizabilities, dissociation of molecules, and barriers of chemical reactions [10]. Hybrid functionals [9] have been shown to describe band gaps accurately for certain systems, when the gaps are calculated from the band structure in the GKS scheme of Eq. (5). They contain both GGA and HF components and thus localization and delocalization errors can cancel, however, not completely [10].

The argument above is carried out in the infinite lattice constant limit but it is clear that some physical inconsistencies will be found at finite lattice constant in a normal periodic calculation. To illustrate the delocalization and localization errors we carry out a simple calculation on a one-dimensional system based on previous work [17] of H_2 polymer polarizability. We take a set of H_2 molecules that clearly shows the characteristic behavior of Fig. 1, a chain made up of 16 atoms, and repeat it with a 15 a.u. distance between the units. The results in Fig. 3 have the same behavior as in Fig. 2, showing that the energy of a δN -concave functional (HF) remains the same, independent of the number of units M . With LDA the convex behavior disappears as the delocalization error increases with M , and the $E(N)$ curve becomes linear following the initial slope of the basic unit. We can see that the point at the integer corresponding to the addition or subtraction is qualitatively wrong with a much too low energy. Even changing the distance between the units from 15 to 5 a.u. has no effect on Fig. 3.

In Fig. 4 we show a plot of the difference density $[\rho(N + 1) - \rho(N)]$ for different sized units. We observe a clear

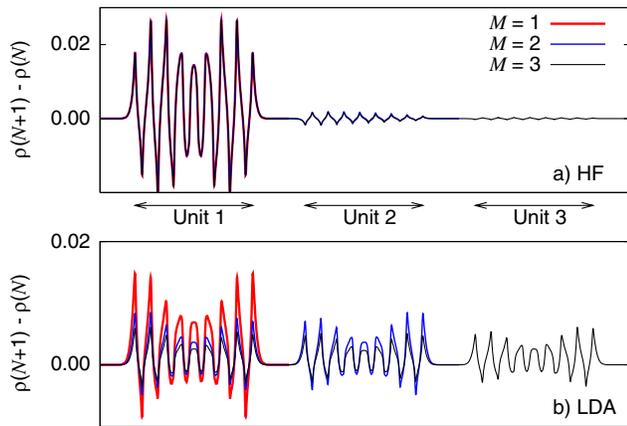


FIG. 4 (color online). Density difference for the $(\text{H}_{16})_M$ system with HF and LDA.

difference between HF and LDA. HF localizes the extra electron to just one of the units, whereas LDA delocalizes the extra electron over all the units with a corresponding drop in energy. This clearly shows the delocalizing bias of LDA and the systematic error it causes.

All of these ideas tie in with the understanding of the band gap originally from Kohn [18], who related the band gap to localization [19], and Perdew [4] who investigated a lattice of hydrogen atoms, and related the error of LDA to the creation of a delocalized hole. A systematic error in the band gap implies a systematic error in describing localization, which is just what we find. This error has widespread relevance in the calculation of many electronic properties of solids and large systems, as the delocalization error in δN -convex functionals affects large molecules and solids much more than small molecules.

The understanding given here also explains why time dependent LDA can be more successful for the calculation of small molecules and metals than for nonmetallic infinite solids and polymers [16], where there is a basic problem in the description of the response of the density due to the delocalization error. It can also explain why HF and similar methods have problems with metallic systems such as jellium due to the fact that the band-gap calculation suffers from the localization error and opens a gap when the true nature should be to be delocalized with no gap.

We should note that if it is possible to counteract the delocalizing bias of LDA by somehow localizing the added electron or hole, then the E versus N curve would resemble that of a small molecule with δN -convex behavior. The difference of the energy at the integer points may then give a reasonable estimate of the band gap, explaining some of the results in the literature (e.g., [20]). However the localization size would now be a key issue.

To conclude, we have shown that the errors in the energy for fractional charges in finite systems lead to systematic errors in larger systems. The addition of an electron (or hole) is poorly described by δN -convex functionals, such

as LDA, as they delocalize the added electron (or hole). This leads to errors in the initial slope of the E versus N curve and therefore the eigenvalues, as is clearly seen in a band-gap calculation. This also means that the explicit calculation of I and A , for example, in large cluster calculations, will suffer from the same error. Functionals which have δN -concave behavior, such as HF, have the opposite tendency and localize electrons. These errors are pervasive and lead to systematic errors in calculations of large molecules and the solid state.

The understanding offered in this work explains the physical nature of the error in the band gap from commonly used approximate functionals, and shows the implications of this error to the calculation of many other properties of solids, from optics to electron transport. A path forward is shown: by constructing functionals free from localization or delocalization error, one would be able to overcome most of the problems.

This work was supported by the National Science Foundation. Discussions with Dr. Xiangqian Hu and Dr. San-Huang Ke have been helpful.

-
- [1] J.P. Perdew and M. Levy, Phys. Rev. Lett. **51**, 1884 (1983).
 - [2] L.J. Sham and M. Schlüter, Phys. Rev. Lett. **51**, 1888 (1983).
 - [3] J.P. Perdew, in *Density Functional Methods in Physics*, edited by R. Dreizler and J. da Providencia (Plenum, New York, 1985), pp. 265–308.
 - [4] J.P. Perdew, Int. J. Quantum Chem. Symp. **19**, 497 (1986).
 - [5] A.J. Cohen, P. Mori-Sánchez, and W. Yang (to be published).
 - [6] A.D. Becke, Phys. Rev. A **38**, 3098 (1988).
 - [7] C. Lee, W. Yang, and R.G. Parr, Phys. Rev. B **37**, 785 (1988).
 - [8] J.P. Perdew, K. Burke, and M. Ernzerhof, Phys. Rev. Lett. **77**, 3865 (1996).
 - [9] A.D. Becke, J. Chem. Phys. **98**, 5648 (1993).
 - [10] P. Mori-Sánchez, A.J. Cohen, and W. Yang, J. Chem. Phys. **125**, 201102 (2006).
 - [11] A.J. Cohen, P. Mori-Sánchez, and W. Yang, J. Chem. Phys. **126**, 191109 (2007).
 - [12] A. Ruzsinszky, J.P. Perdew, G.I. Csonka, O.A. Vydrov, and G.E. Scuseria, J. Chem. Phys. **126**, 104102 (2007).
 - [13] J.P. Perdew, R.G. Parr, M. Levy, and J.L. Balduz, Jr., Phys. Rev. Lett. **49**, 1691 (1982).
 - [14] W. Yang, Y. Zhang, and P.W. Ayers, Phys. Rev. Lett. **84**, 5172 (2000).
 - [15] J.D. Talman and W.F. Shadwick, Phys. Rev. A **14**, 36 (1976).
 - [16] G. Onida, L. Reining, and A. Rubio, Rev. Mod. Phys. **74**, 601 (2002).
 - [17] P. Mori-Sánchez, Q. Wu, and W. Yang, J. Chem. Phys. **119**, 11 001 (2003).
 - [18] W. Kohn, Phys. Rev. **133**, A171 (1964).
 - [19] R. Resta and S. Sorella, Phys. Rev. Lett. **82**, 370 (1999).
 - [20] P.A. Schultz, Phys. Rev. Lett. **96**, 246401 (2006).