## Quasiparticle Self-Consistent GW Theory

M. van Schilfgaarde,<sup>1</sup> Takao Kotani,<sup>1</sup> and S. Faleev<sup>2</sup>

<sup>1</sup>Arizona State University, Tempe, AZ, 85287 <sup>2</sup>Sandia National Laboratories, Livermore, CA 94551

## Abstract

In past decades the scientific community has been looking for a reliable first-principles method to predict the electronic structure of solids with high accuracy. Here we present an approach which we call the quasiparticle self-consistent GW approximation (QPscGW). It is based on a kind of self-consistent perturbation theory, where the self-consistency is constructed to minimize the perturbation. We apply it to selections from different classes of materials, including alkali metals, semiconductors, wide band gap insulators, transition metals, transition metal oxides, magnetic insulators, and rare earth compounds. Apart some mild exceptions, the properties are very well described, particularly in weakly correlated cases. Self-consistency dramatically improves agreement with experiment, and is sometimes essential. Discrepancies with experiment are systematic, and can be explained in terms of approximations made. The Schrödinger equation is the fundamental equation of condensed matter, and the importance of being able to solve it reliably can hardly be overestimated. The most widely used theory in solids, and now in quantum chemistry, is the celebrated local density approximation (LDA)[1]. In spite of its successes, it is well known that the LDA suffers from many deficiencies, even in weakly correlated materials (see Figs. 1 and 2). This has stimulated the development of flavors of extensions to the LDA to redress one or another of its failures, such as the LDA+U method. Each of these methods improves one failing or another in the LDA, but they often bear a semi-empirical character, and none can be considered universal and parameter-free. Thus we are far from a precise and universally applicable theory for solids, with attendant limits their ability to predict materials properties.

The random phase approximation (RPA) or GW approximation (GWA, G=Green's function, W=screened coulomb interaction) of Hedin[2] is almost as old as the LDA. A major advance was put forward by Hybertsen and Louie[3] when they employed LDA eigenfunctions to generate the GW self-energy  $\Sigma = iGW$ , and showed that fundamental gaps in  $sp^3$ bonded materials were considerably improved over the LDA. Since that seminal work, many papers and some reviews[4, 5, 6] have been published on GW theory and extensions to it. One problem that has plagued the GW community has been that calculated results of the same quantities tend to vary between different groups, much as what occurred in the early days of the LDA. This is because *further* approximations are usually employed which significantly affect results. Almost ubiquitous is the 1-shot approximation where (following Hybertsen)  $\Sigma \approx i G^{\text{LDA}}W^{\text{LDA}}$ : *i.e.*,  $\Sigma$  is computed from LDA eigenfunctions. However, there is an emerging consensus[7, 8, 9, 10, 11, 12] that, when cores are treated adequately[11],  $G^{\text{LDA}}W^{\text{LDA}}$  bandgaps are underestimated even in (weakly correlated) semiconductors. (see top panel of Fig. 1; note especially CuBr).

In general, one-shot GW approaches are rather unsatisfactory. The quality of the  $G^{\text{LDA}}W^{\text{LDA}}$  approximation is closely tied to the quality of LDA starting point[11], and is adequate to construct G and W only under limited circumstances[11]. It can fail even qualitatively in transition-metal and rare earth compounds such as CoO and ErAs[11]. Some kind of self-consistency is essential: the QP levels should not be an artifact of the starting conditions. The full self-consistent GW method (full scGW) determines G self-consistently from  $\Sigma = iGW$ , which in turn generates G. Here  $W = v(1 - vP)^{-1}$  where  $P = -iG \times G$  and v are respectively the irreducible polarization function and (bare) Coulomb interaction.



FIG. 1: Fundamental gaps of sp compounds from LDA (squares) and  $G^{\text{LDA}}W^{\text{LDA}}$  (circles) in top panel, and from QPscGW, Eqn. (2), in bottom panel. The spin-orbit coupling was subtracted by hand from the calculations. The  $G^{\text{LDA}}W^{\text{LDA}}$  gaps improve on the LDA, but are still systematically underestimated. For QPscGW data, zincblende compounds with direct  $\Gamma - \Gamma$  transitions are shown as green circles; All other gaps are shown as blue squares. Errors are small and highly systematic, and would be smaller than the figure shows if the electron-phonon renormalization were included,

In the few cases where it has been applied, some difficulties were found: in particular the valence bandwidth of the homogeneous electron gas[15] is  $\sim 15\%$  wider than the noninteracting case, whereas the  $G^{\text{LDA}}W^{\text{LDA}}$  width is  $\sim 15\%$  narrower, in agreement with experiment for Na (see Fig. 2). A recent (nearly) full scGW study of Ge and Si also overestimates the



FIG. 2: Comparison of LDA (blue dashes),  $G^{\text{LDA}}W^{\text{LDA}}$  (red dots) and QPscGW (green lines) energy bands in GaAs (left) and Na (right). Circles are experimental data, with spin-orbit coupling subtracted by hand. The QPscGW fundamental gap and conduction-band effective mass  $(E_g = 1.77 \text{ eV} \text{ and } m_c^* = 0.077 m_0)$  are slightly overestimated, and the optical dielectric constant underestimated ( $\epsilon_{\infty} = 8.4$ ):  $E_g^{\text{expt}}(0\text{K}) = 1.52 \text{ eV}, m_c^{\text{*,expt}} = 0.065 m_0$ , and  $\epsilon_{\infty}^{\text{expt}} = 10.8$ . For comparison,  $E_g^{\text{G}^{\text{LDA}}W^{\text{LDA}} = 1.29 \text{ eV}$  and  $m_c^{\text{*,G}^{\text{LDA}}W^{\text{LDA}} \approx 0.059 m_0$ , while  $E_g^{\text{LDA}} = 0.21 \text{ eV}$  and  $m_c^{\text{*,LDA}} = 0.020 m_0$ . The correspondence between QPscGW and experiment at other known levels at  $\Gamma$ , L, and X, the Ga 3d level near -18 eV is representative of nearly all available data for spsystems. For Na, the QPscGW occupied bandwidth is 15% smaller than the LDA. Circles taken from photoemission data [13]; square from momentum electron spectroscopy [14].

valence bandwidth[7], though the fundamental gaps are well described.

In Ref. [16] we proposed an ansatz for a different kind of scGW, and demonstrated that it radically improves the quasiparticle (QP) levels in the oxides MnO and NiO. In this Letter, we ground the idea on an underlying principle—namely optimization of the effective one-body hamiltonian  $H^0$  by minimizing the perturbation to it—and propose it as a universal approach to the reliable prediction of the electronic structure. We show that this approach, which we call the quasiparticle self-consistent GW (QPscGW) method, results in accurate predictions of excited-state properties for a large number of weakly and moderately correlated materials. QP levels are uniformly good for all materials studied: not just fundamental gaps in semiconductors but for nearly all levels where reliable experimental data are available. Even in strongly correlated d and f electron systems we studied, errors are somewhat larger but still systematic.

The GWA is usually formulated as a perturbation theory starting from a non-interacting Green's function  $G^0$  for given one-body hamiltonian  $H^0 = \frac{-\nabla^2}{2m} + V^{\text{eff}}$ .  $H^0$  is noninteracting, so  $V^{\text{eff}}$  is static and hermitian but it can be nonlocal. Because the GWA is an approximation to the exact theory, the one-body effective hamiltonian  $H(\omega) = \frac{-\nabla^2}{2m} + V^{\text{ext}} + V^{\text{H}} + \Sigma(\omega)$ depends on  $V^{\text{eff}}$  and is a functional of it: the Hartree potential  $V^{\text{H}}$  is generated through  $G^0 = 1/(\omega - H^0 \pm i\epsilon)$ , and the GWA generates  $\Sigma(\omega)$ .  $H(\omega)$  determines the time-evolution of the one-body amplitude for the many-body system.

QPscGW is a prescription to determine the optimum  $H^0$ : we choose  $V^{\text{eff}}$  based on a self-consistent perturbation theory so that the time-evolution determined by  $H^0$  is as close as possible to that determined by  $H(\omega)$ , within the RPA. This idea means that we have to introduce a norm M to measure the difference  $\Delta V(\omega) = H(\omega) - H^0$ ; the optimum  $V^{\text{eff}}$  is then that potential which minimizes M. A physically sensible choice of norm is

$$M[V^{\text{eff}}] = \text{Tr} \left[ \Delta V \delta(\omega - H^0) \{ \Delta V \}^{\dagger} \right]$$
  
+ Tr  $\left[ \{ \Delta V \}^{\dagger} \delta(\omega - H^0) \Delta V \right]$  (1)

where the trace is taken over  $\mathbf{r}$  and  $\omega$ . Exact minimization M is apparently not tractable, but an approximate solution can be found. Note that M is positive definite. If we neglect the second term and ignore the restriction that  $V^{\text{eff}}$  is hermitian, we have the trivial minimum  $M[V^{\text{eff}}] = 0$  at  $V^{\text{eff}} = V^{\text{ext}} + V^{\text{H}} + V^{\text{xc}}$  where  $V^{\text{xc}} = \sum_{ij} |\psi_i\rangle \Sigma(\varepsilon_j)_{ij} \langle \psi_j|$ . Here  $\Sigma(\varepsilon_i)_{ij} =$  $\langle \psi_i | \Sigma(\varepsilon_i) | \psi_j \rangle$ , and  $\{\psi_i, \epsilon_i\}$  are eigenfunctions and eigenvalues of  $H^0$ . The second term is similarly minimum with  $\Sigma(\varepsilon_i) \to \Sigma(\varepsilon_j)$ . An average of the hermitian parts of these solutions results in

$$V^{\rm xc} = \frac{1}{2} \sum_{ij} |\psi_i\rangle \left\{ \operatorname{Re}[\Sigma(\varepsilon_i)]_{ij} + \operatorname{Re}[\Sigma(\varepsilon_j)]_{ij} \right\} \langle \psi_j|.$$
(2)

Re signifies the hermitian part. This result is the same as Eq. (2) in Ref. 16.

We identify solutions to  $H^0$  as "bare QP", which interact via the (bare) v. The dressed QP consists of the central bare QP plus induced polarized clouds of the other bare QPs' this is nothing but the physical picture in RPA to calculate poles of G from  $G^0$ . In the charged Fermi liquid theory [17] of Landau and Silin, the QP interact via v in addition to the short-range Landau interaction ( $f_{pp'}$  in Ref. 17; see Eq. (3.41)). We can virtually construct the Landau-Silin QP from  $G^0$  by the calculation of the non-RPA contributions to  $\Sigma$ . Or we can identify our bare QP as the Landau-Silin QP if we assume they are minimally affected by such contributions.

Our QPscGW is conceptually very different from the full scGW. In the latter case the electron-hole mediated state making  $P = -iG \times G$  is suppressed by the square of the renormalization factor  $Z \times Z$ , and includes physically unclear contributions such as: QP × (incoherent parts) [16, 18]. P loses its physical meaning as the density response function,  $P = \delta n/\delta V$ , but is merely an intermediate construction in the self-consistency cycle. Such a construction does not give reasonable W even in the electron gas [15, 19], resulting in a poor G.

We now turn to QPscGW results, focusing on the QP energies given by  $H^0$ . Fig 2 shows that the QPscGW valence band in Na properly narrows relative to the LDA by 15%. Indeed, for nearly all the sp semiconductors studied, calculated QP levels generally agree very closely with available experimental data. The best known are the fundamental gaps, shown in Fig. 1. QPscGW data is divided into circles for materials whose gap is a  $\Gamma - \Gamma$ transition and squares for all other kinds. Roughly,  $\Gamma - \Gamma$  transitions are overestimated by 0.2 eV, while the remaining gaps are overestimated by 0.1 eV. Errors appear to be larger for wide-gap, light-mass compounds (bearing elements C, N, and especially O); however, the calculations omit reduction in the gaps by the nuclear zero-point motion. This effect has been studied through varying isotopic mass in some tetrahedral semiconductors[20]. It is largest for light compounds: T=0 the gap is reduced by ~0.3 eV in diamond and ~0.2 eV in AlN, but  $\leq 0.1$  eV for heavier compounds. Because the renormalization been measured only for a few cases, we do not include it here.

Apart from some mild exceptions, QPscGW generates a consistently precise description of the electronic structure in sp systems, including other known excitations. This is illustrated in Fig. 2, where GaAs was chosen because of the abundance of available experimental data. It is notable that the errors are not only small, but unlike the  $G^{\text{LDA}}W^{\text{LDA}}$  or LDA, they are highly systematic: compare, for example, the fundamental gaps (Fig. 1). We may expect that the bandgaps should be overestimated, because the RPA dielectric function omits electronhole correlation effects. Thus  $\epsilon^{RPA}$  should be too small and under-screen W. Indeed, the optical dielectric constant  $\epsilon_{\infty}$  is systematically underestimated slightly (Fig. 2). Similar consistencies are found in the effective masses. The conduction-band mass at  $\Gamma$ ,  $m_{c\Gamma}^{*,\text{QPscGW}}$  consistently falls within a few percent of experimental data for wide-gap materials, but as the gap becomes smaller (induced by, e.g. scaling  $\Sigma$ ),  $m_{c\Gamma}^{*,\text{QPsc}GW}/m_{c\Gamma}^{*,\text{expt}}$  scales essentially as the ratio of the QPscGW gap to the experimental one, as expected when the gap become small. Taking data for GaAs from Fig. 2 for example, we obtain  $E_g^{\text{QPsc}GW}/E_g^{\text{expt}} = 1.16$ , and  $m_{c\Gamma}^{*,\text{QPsc}GW}/m_{c\Gamma}^{*,\text{expt}} = 1.18$ .

TABLE I: Valence d bandwidths  $W_d$  (calculated at  $\Gamma$  for Ti,Cr, and Co, and at N for Fe, and at X for Ni), relative position of s and d band bottoms  $\epsilon_{sd}$ , splittings  $\Delta E_x$  between majority and minority d (or f) states, and magnetic moments in 3d compounds and Gd.

	$W_d \ (eV)$			$\epsilon_{sd} \ (eV)$		
	LDA	QPscGW	Expt	LDA	QPscGW	<sup>7</sup> Expt
Ti	6.0	5.7		3.5	4.3	
$\operatorname{Cr}$	6.6	6.2		3.5	4.3	
Fe	5.2	4.6	4.6	3.6	4.4	4.6
Co	4.1	3.8	3.7	4.6	5.3	$4.9{\pm}1$
Ni	4.4	4.0	4.0	4.4	5.0	5.5
	n	noment ( $\mu_I$	3)	$\Delta E_{\rm x} \ ({\rm eV})$		
	LDA	QPscGW	Expt	LDA	QPscGW	<sup>7</sup> Expt
Fe	2.2	2.2	2.2	1.95	1.67	1.75
Co	1.6	1.7	1.6	1.70	1.21	1.08
Ni	0.6	0.7	0.6	0.6	0.5	0.3
MnO	4.5	4.8	4.6			
NiO	1.3	1.7	1.9			
MnAs	3.0	3.5	3.4			
Gd	7.7	7.8	7.6	4.9	16.1	$\sim 12.1$

Table I shows that the 3*d* bandwidth, the relative position of *s* band, exchange splittings  $\Delta E_x$ , are systematically improved relative to the LDA in elemental 3*d* metals, and Gd. QPscGW magnetic moments are systematically overestimated slightly.  $\Delta E_x$  is overestimated in Ni, presumably owing to the neglect of spin fluctuations[21]. QPscGW also predicts with reasonable accuracy the QP levels of all magnetic 3*d* compounds studied, in particular correlated oxides such as MnO and NiO where the LDA fails dramatically. As might be expected, the accuracy deteriorates somewhat relative to sp systems. For example, the QPscGW optical gap in NiO (4.8 eV) was found to be larger than experiment (~4.3 eV). Table I compares the magnetic moments, and Ref. [16] shows in detail the QP levels are consistently well described. However, for Gd, (and for GdP and GdAs) QPscGW overestimates the position of the (empty) minority Gd f shell by ~4 eV, and hence the exchange splitting  $\Delta E_x$ .



FIG. 3: DOS in CeO<sub>2</sub>. Black dots are PES+BIS data[[22]]. Calculated DOS were broadened with Gaussian of width 0.35 eV.  $G^{\text{LDA}}W^{\text{LDA}}$ +eigenvalue-only self-consistency (dotted red line) severely overestimates the position of the Ce f level, while it is slightly overestimated by QPscGW(green line). Broadening of the valence bands relative to LDA is found in all oxides studied, e.g. MgO and TiO<sub>2</sub>, and has important consequences, e.g. in determining valence-band offsets.

Nonmagnetic oxides SrTiO<sub>3</sub>, TiO<sub>2</sub>, and CeO<sub>2</sub>, with conduction bands of d or f character, overestimate fundamental gaps slightly more than their sp counterparts. The QPscGW gaps were found to be 4.19 eV and 3.78 eV in SrTiO<sub>3</sub> and TiO<sub>2</sub>, ~0.8 eV larger than the experimental gaps (~3.3 eV and ~3.1 eV). Fig. 3 compares the QPscGW DOS of CeO<sub>2</sub> with spectroscopic data: the Ce f band is similarly overestimated by QPscGW. This is reasonable, because electron-hole correlation effects are stronger in the narrow d (f) conduction bands. Fig. 3 also shows DOS computed by  $G^{\text{LDA}}W^{\text{LDA}}$ , but with eigenvalueonly self-consistency, where only the diagonal part in Eq. (2) is kept. This constrains the eigenfunctions to the starting (LDA) eigenfunctions; thus the charge and spin densities do not change. While the off-diagonal parts of  $\Sigma$  add a small effect in, e.g. GaAs, their contribution is essential in CeO<sub>2</sub>, even though the occupied states contain only a small amount of Ce f character.

To summarize, the QPscGW theory (apart from some mild exceptions) appears to be an excellent predictor of QP levels for a variety materials selected from the entire periodic table. Self-consistency is an essential part of the theory. From the results obtained so far, this approach shows promise to be universally applicable scheme, sufficient in its own right for QP levels in many materials. In contrast to the LDA or any other popular theory of electronic structure of solids in the literature today, the method is truly *ab initio* with errors that are generally small and highly systematic across many different materials classes. The errors can be attributed missing electron-hole correlation contributions to  $\epsilon$ . When better calculations are necessary (usually where the physics lies completely outside the domain of a one-particle picture, such as the description of excitons, multiplets, or Mott transitions), QPscGW can be taken as an optimum starting point where the relevant many-body contributions to the hamiltonian are (nearly) as small as possible. The systematic character of the error suggests that the dominant terms left out can be described by a few diagrams, in particular the ladder diagrams coupling electrons and holes; the smallness of the error suggests that the additional terms can be added as a perturbation around the QPscGW  $H^0$ , without the need for further self-consistency.

This work was supported by ONR contract N00014-02-1-1025 and BES Contract No. DE-AC04-94AL85000.

- [1] W. Kohn and L. J. Sham, Phys. Rev. 140, A1133 (1965).
- [2] L. Hedin, Phys. Rev. **139**, A796 (1965).
- [3] M. S. Hybertsen and S. Louie, Phys. Rev. B **34**, 5390 (1986).
- [4] F. Aryasetiawan and O. Gunnarsson, Rep. Prog. Phys 61, 237 (1998).
- [5] W. G. Aulbur, L. Jönsson, and J. Wilkins, in *Solid State Physics, vol. 54*, edited by H. Ehrenreich and F. Saepen (2000), p. 1.
- [6] G. Onida, L. Reining, and A. Rubio, Rev. Mod. Phys 74, 601 (2002).
- [7] W. Ku and A. G. Eguiluz, Phys. Rev. Lett. 89, 126401 (2002).
- [8] S. Lebegue et al., Phys. Rev. B 67, 155208 (2003).
- [9] M. L. Tiago, S. Ismail-Beigi, and S. G. Louie, Phys. Rev. B 69, 125212 (2004).
- [10] A. Fleszar and W. Hanke, Phys. Rev. B **71**, 045207 (2005).
- M. van Schilfgaarde, T. Kotani, and S. Faleev, submitted. Preprint http://arxiv.org/abs/cond-mat/0508295.

- [12] T. Kotani and M. van Schilfgaarde, Sol. State Commun. 121, 461 (2002).
- [13] I.-W. Lyo and E. W. Plummer, Phys. Rev. Lett. 60, 1558 (1988).
- [14] M. Vos et. al., Phys. Rev. B 66, 155414 (2002).
- [15] B. Holm and U. von Barth, Phys. Rev. B 57, 2108 (1998).
- [16] S. V. Faleev, M. van Schilfgaarde, and T. Kotani, Phys. Rev. Lett. 93, 126406 (2004).
- [17] D. Pines and P. Nozieres, The Theory of Quantum Liquid. Vol I (W.A.Benjamin Inc., New York, 1966).
- [18] F. Bechstedt, K. Tenelsen, B. Adolph, and R. D. Sole, Phys. Rev. Lett. 78, 1528 (1997).
- [19] D. Tamme, R. Schepe, and K. Henneberger, Phys. Rev. Lett. 83, 241 (1999).
- [20] F. J. Manjon et al., Eur. Phys. J. B **00211**, 1 (2004).
- [21] F. Aryasetiawan, Phys. Rev. B 46, 13501 (1992).
- [22] E. Wuilloud, B. Delley, W. D. Schneider, and Y. Baer, Phys. Rev. Lett. 53, 202 (1984).