



## Toward Automatic Rhodopsin Modeling as a Tool for High-Throughput Computational Photobiology

This is the peer reviewed version of the following article:

*Original:*

Melaccio, F., Del Carmen Marín, M., Valentini, A., Montisci, F., Rinaldi, S., Cherubini, M., et al. (2016). Toward Automatic Rhodopsin Modeling as a Tool for High-Throughput Computational Photobiology. JOURNAL OF CHEMICAL THEORY AND COMPUTATION, 12(12), 6020-6034 [10.1021/acs.jctc.6b00367].

*Availability:*

This version is available <http://hdl.handle.net/11365/1008179> since 2017-09-15T21:59:22Z

*Published:*

DOI:10.1021/acs.jctc.6b00367

*Terms of use:*

Open Access

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. Works made available under a Creative Commons license can be used according to the terms and conditions of said license.

For all terms of use and more information see the publisher's website.

(Article begins on next page)

# Towards Automatic Rhodopsin Modeling as a Tool for High-throughput Computational Photobiology

Federico Melaccio<sup>†,\*</sup>, María del Carmen Marín<sup>†,¶</sup>, Alessio Valentini<sup>†</sup>, Fabio Montisci<sup>†</sup>, Silvia Rinaldi<sup>†</sup>, Marco Cherubini<sup>†</sup>, Xuchun Yang,<sup>¶</sup> Yoshitaka Kato<sup>¶¶</sup>, Michael Stenrup,<sup>‡</sup> Yoelvis Orozco-Gonzalez,<sup>#,¶</sup> Nicolas Ferré,<sup>‡</sup> Hoi Ling Luk<sup>¶</sup>, Hideki Kandori<sup>¶¶</sup> and Massimo Olivucci<sup>†,¶,##,\*\*\*</sup>

<sup>†</sup>Department of Biotechnology, Chemistry e Pharmacy, Università di Siena, via A. Moro 2, I-53100 Siena, Italy, <sup>¶</sup>Department of Chemistry, Bowling Green State University, Bowling Green, OH 43403, USA, <sup>¶¶</sup>Department of Frontier Materials, Graduate School of Engineering, Nagoya Institute of Technology, Gokiso-cho, Showa-ku, Nagoya, Japan, <sup>‡</sup>Aix Marseille Univ, CNRS, ICR, Marseille, France, <sup>#</sup>Institut de Physique et Chimie des Matériaux de Strasbourg, UMR 7504 Université de Strasbourg-CNRS, F-67034 Strasbourg, France, <sup>##</sup>USIAS Institut d'Études Avancées, Université de Strasbourg, 5 allée du Général Rouvillois, F-67083 Strasbourg, France.

## Abstract

We report on a prototype protocol for the automatic and fast construction of congruous sets of QM/MM models of rhodopsin-like photoreceptors and of their mutants. In the present implementation the information required for the construction of each model is essentially a crystallographic structure or a comparative model complemented with information on the protonation state of ionizable side-chains and distributions of external counterions. Starting with such information a model formed by a fixed environment system, a flexible cavity system and a chromophore system is automatically generated. The results of the predicted vertical excitation energy for 27 different rhodopsins including vertebrate, invertebrate and microbial pigments indicate that such basic models could be employed for predicting trends in spectral changes and/or correlate the spectral changes with structural variations in large sets of proteins.

## Introduction

Cheap and fast tools for gene sequencing, protein expression and analysis are routinely used for high-throughput genomics.<sup>1,2</sup> However, due to experimental difficulties and longer timescales (e.g. for crystallization), protein structure determination cannot presently be performed at the same fast pace: a fact that is slowing down the discovery of proteins with new features as well as their *ex novo* design.<sup>3</sup> These difficulties affect heavily the field of photobiology where the number of known photoreceptor structures is limited. For instance the crystal structure of bovine rhodopsin, the retina dim-light visual photoreceptor, is the only structure available for studying monochromatic and color vision in vertebrates. In principle, the gap could be filled by the construction of sufficiently accurate atomistic computer models of the set of photoreceptors of interest, which would allow *in silico* functional characterization, property screening and *ex-novo* design. However, to be useful, such models should be produced using a fast and standardized (in the sense of a well defined and replicable) protocol with known error bars.

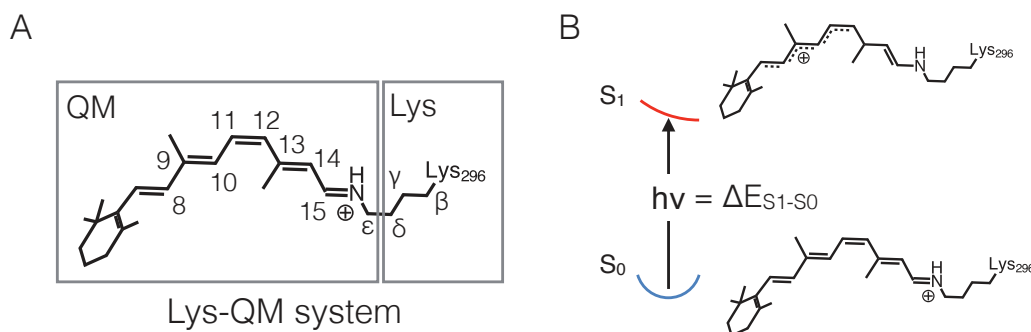
Computer models of biological photoreceptors for studying properties like spectral tuning and photoisomerization timescales are computationally complex. As originally proposed by Warshel<sup>4</sup>, these models have to incorporate a quantum mechanical (QM) description of the chromophore (i.e. the protein spectroscopically and photochemically active prosthetic group) and account for the interaction with the apoprotein using a computationally less expensive molecular mechanics (MM) force field. In the past, our group has contributed to show that hybrid quantum-mechanics / molecular-mechanics (QM/MM) approaches combining an *ab initio* multiconfigurational QM description of the light absorbing moiety (i.e. the chromophore) with a MM description of its molecular environment can be applied to different photoresponsive proteins including bovine rhodopsin and other homologous proteins.<sup>5-10</sup> Similar or distinct QM/MM models of such systems have been reported by other groups and mainly used for understanding their absorption spectra.<sup>11-14</sup> Indeed, the rhodopsin protein family represents an impressive case of regulation of the wavelength of the absorption maximum ( $\lambda_{\text{max}}^{\text{a}}$ ) playing a role in fundamental biological processes such as vision, chromatic adaptation, ion-gating and ion-pumping.<sup>15,16</sup> In these trans-membrane proteins, the chromophore  $\lambda_{\text{max}}^{\text{a}}$  is modulated by the environment from 420 nm (human short wave-sensitive pigment, hSWS) to 587 nm (sensory rhodopsin I).<sup>16</sup> As displayed in Scheme 1A, the chromophore is formed by a retinylidene isomer bound to a lysine residue via an iminium function (also known as a protonated Schiff base, PSB). The detailed knowledge of the molecular factors controlling the  $\lambda_{\text{max}}^{\text{a}}$  of rhodopsins is important not only for understanding the functions of these photoreceptors, but also for the rational design of novel, genetically encodable tools. In fact, there is a growing interest in

1  
2  
3 employing rhodopsins to switch “on” and “off” metabolic pathways, gene expression and ion  
4 channels.<sup>17-20</sup> Furthermore, rhodopsins are currently central to the field of optogenetics<sup>21</sup> where the  
5 gene expressing a light-gated ion-channel channelrhodopsin is used to engineer light-sensitivity into  
6 neurons.  
7  
8  
9

10 As mentioned above QM/MM models of rhodopsins are well established, but a quick literature  
11 survey demonstrates that even models of the same protein can be hardly compared. Taking bovine  
12 rhodopsin (Rh) as a reference, being the most representative and studied system, at least 6 different  
13 QM/MM setups have been used<sup>6,12,13,22-24</sup> (see Supporting Information for more details) resulting in a  
14 36 nm range for the predicted Rh visible absorption maximum wavelength ( $\lambda^a_{\text{max}}$ ). The substantial  
15 differences displayed by similar approaches (e.g. Andruniow *et al.*<sup>6</sup> and Tomasello *et al.*<sup>12</sup>), must be  
16 due to the distinct manipulations used for the preparation of the initial protein structure (i.e. the  
17 structure equilibration, the external environment/solvation, the protonation states of ionizable  
18 residues, the internal waters, etc.) and the multiple variables to be selected during the construction  
19 of the final QM/MM model (e.g. the choice of the QM method or methods, the treatment of the  
20 QM/MM frontier and of QM/MM non-bonding interactions, the employed MM force field).  
21 Moreover, these manipulations are prone to operator errors and, in practice, are cumbersome to  
22 replicate. Thus, having a well defined QM/MM models which can be readily replicated in other labs  
23 and allow for consistent evaluation of the error associated to the computation of specific properties,  
24 is highly desirable, if not mandatory, for the correct application and further development of the  
25 methodology.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

39 Here we report on the design and testing of an Automatic Rhodopsin Modeling protocol (from now  
40 on simply called ARM) addressing the above issue. Such protocol has the primary target of  
41 generating congruous QM/MM models of rhodopsins that may facilitate systematic rhodopsin  
42 studies.<sup>25,26</sup> More specifically, ARM automates and speeds up the construction of models based on  
43 an additive, H-link-atom scheme<sup>27</sup> for both wild type and mutant rhodopsins. ARM is *not* designed to  
44 produce the most accurate models possible (see for instance the models of ref. 28-32 targeting  
45 accurate spectroscopic studies), but basic, gas-phase and computationally fast models aimed to the  
46 *rationalization and prediction of trends between sequence variability and function*. Therefore, ARM  
47 models aim to satisfy the following desirable features for a high-throughput approach: (a)  
48 automation, to reduce accidental errors and avoid biased modeling; (b) speed, to deal with large sets  
49 of protein of the same family; (c) documented accuracy, to translate results into experimentally  
50 assessable hypothesis; (d) transferability, to treat rhodopsins with large differences in sequence (i.e.  
51 from different domains and kingdoms).  
52  
53  
54  
55  
56  
57  
58  
59  
60

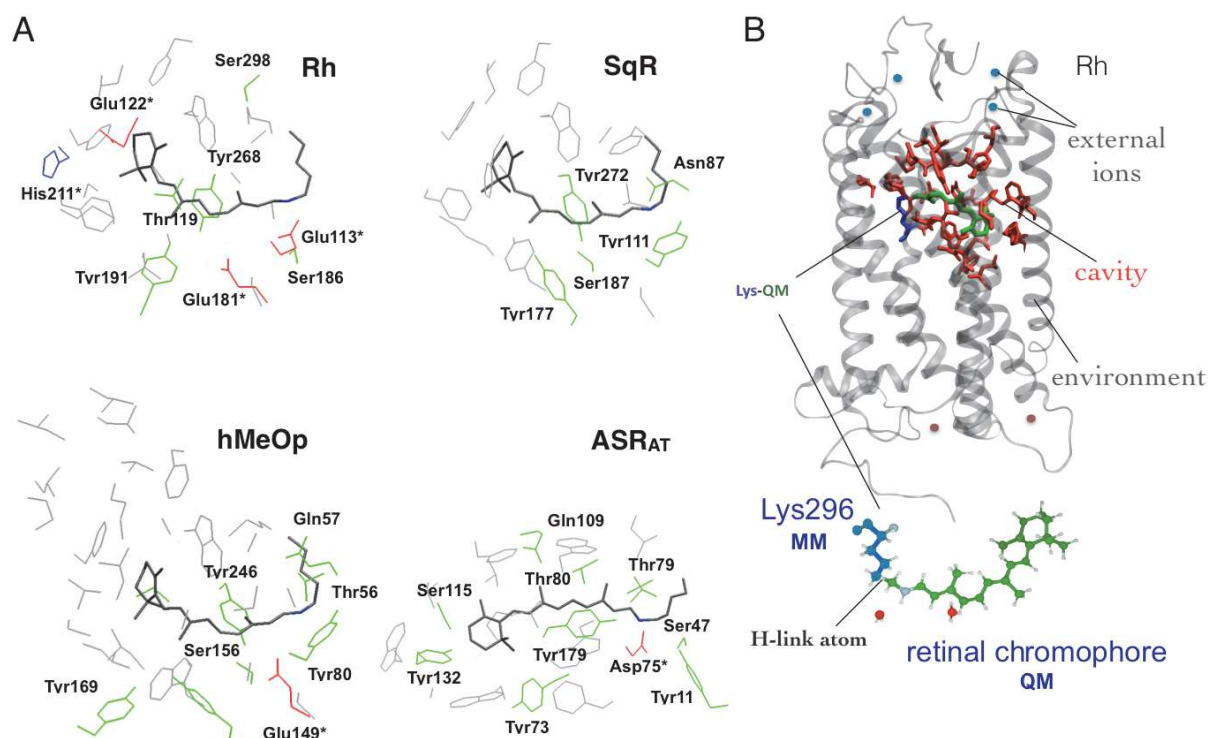
1  
2  
3 In the present work, a minimal chromophore cavity, a minimal QM subsystem and short molecular  
4 dynamics (MD) conformational sampling are considered characteristics of the investigated model,  
5 thus making its construction fast and avoiding unbiased construction steps (see features a-b above).  
6  
7 On the other hand, a benchmarking set comprising vertebrate, invertebrate and microbial  
8 rhodopsins is employed to explore features c-d through  $\lambda_{\text{max}}^{\text{a}}$  computations (i.e. the corresponding  
9 vertical excitation process of Scheme 1B). More specifically we use: the all-*trans* and 13-*cis* isomers  
10 of *Anabaena* Sensory Rhodopsin (ASR, microbial, sensorial) and 5 of their mutants for which we also  
11 measure the  $\lambda_{\text{max}}^{\text{a}}$ ,<sup>33,34</sup> bovine Rhodopsin (Rh, vertebrate, visual pigment)<sup>35</sup> and 7 of its mutants for  
12 which the measured  $\lambda_{\text{max}}^{\text{a}}$  is reported in the literature; squid Rhodopsin (SqR, invertebrate, visual  
13 pigment);<sup>36</sup> human melanopsin (hMeOp, vertebrate, non-visual pigment);<sup>9,37</sup> the light- and dark-  
14 adapted forms of bacteriorhodopsin (bR<sub>LA</sub> and bR<sub>DA</sub>, an Archaea proton pump)<sup>38,39</sup>; blue  
15 proteorhodopsin (PR, an eubacterial proton pump)<sup>40</sup> and the a chimera channelrhodopsin (ChR<sub>C1C2</sub>, a  
16 lab construct of microbial eukaryotic light-gated ion channels from the green alga *Chlamydomonas*  
17 *reinhardtii*).<sup>31</sup> The selected single-site mutants are S86D, S214D, L83Q, V112N, W76F for ASR and  
18 G90S, T94S, T118A, F261Y, W265F, W265Y and A292S for Rh. Moreover, to further increase the  
19 scope of our approach, a model of the Rh primary photocycle intermediate bathorhodopsin  
20 (bathoRh, experimental  $\lambda_{\text{max}}^{\text{a}}$  of 543 nm at low temperature<sup>41</sup>) has been built from its crystal  
21 structure<sup>42</sup>. In total the set spans 26 observed  $\lambda_{\text{max}}^{\text{a}}$  values comprised between 473 and 568 nm  
22 corresponding to vertical excitation energies (i.e., as illustrated in Scheme 1B, the  $\Delta E_{S_1-S_0}$  potential  
23 energy gap computed at the  $S_0$  equilibrium structure which is also assumed to correspond to the  
24  $\lambda_{\text{max}}^{\text{a}}$  value. Of course, in doing so, we neglect the effect of the homogeneous and inhomogeneous  
25 broadening which we assume to be similar in different rhodopsin species) in the 52 to 61 kcal/mol  
26 range. The set has been designed in such a way to include evolutionary distant sensory rhodopsins  
27 from different kingdoms (Animalia and Bacteria) and phyla (Chordata, Mollusca and Cyanobacteria),  
28 two proton pumps from different kingdoms (Archaea and Bacteria) and one light-gated ion channel  
29 from the kingdom Plantae. The “evolutionary distance” between representative members of the set  
30 is documented in the Supporting Information in terms of percentage of homology and in Fig. 1A by  
31 comparing the structure of the cavities of four sensory rhodopsins.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



**Scheme 1** A. Structure of the retinal chromophore of bovine rhodopsin and of the Lys-QM system of the ARM model. B. Schematic illustration of the vertical excitation of the chromophore. Notice the different charge distribution in the ground ( $S_0$ ) and excited ( $S_1$ ) states conjugated carbon framework.<sup>43,44</sup>

Photobiological studies focusing on both spectroscopy and chemical reactivity require the use of accurate, correlated QM methods.<sup>26</sup> ARM employs the multiconfigurational Complete Active Space Self Consistent Field (CASSCF)<sup>45</sup> method to obtain ground state geometries. Excitation energies are then computed using multiconfigurational second-order perturbation theory (CASPT2)<sup>46</sup> to recover the missing dynamical electron correlation associated with the CASSCF description. Such a CASPT2//CASSCF treatment has been extensively investigated and its limitations are well understood. Here, we stress that, as previously documented, the rather limited ca. 3-4 kcal/mol error<sup>5-8,14,25,47-50</sup> in excitation energy reported in several studies<sup>51-53</sup> and comprising CASPT2//CASSCF/6-31G\*/MM calculations are partly due to error cancellations associated to the limited quality of CASSCF/6-31G\* equilibrium geometries.<sup>54</sup> Nevertheless we decided to develop ARM on the basis of such QM treatment because CASSCF allows the use of state-averaged energies for mapping thermal and photochemical reaction paths with different electronic characters and for propagating semi-classical excited state trajectories. Therefore CASSCF/6-31G\*/MM equilibrium models of the selected rhodopsin benchmark set (from now on called ARM models) are suitable for consistent and comprehensive photochemical studies.

Below we show that ARM represents a first-order answer to the need for high-throughput QM/MM model generation for rhodopsins, yielding blue-shifted  $\Delta E_{S_1-S_0}$  values with a standard deviation of few kcal/mol and absolute deviations suggesting the presence of a systematic error. Additionally, using the same set, we document the sensitivity and limitations of ARM models with respect to: (i) the uncertainty in the assignment of the ionization state of the protein residues, (ii) the external (i.e. protein surface) counterion placement and distribution, and (iii) the stability of the side-chain conformation affecting the HBN near the chromophore.



**Figure 1** Rhodopsin models. A. Examples of cavity and chromophore variability in the selected benchmark set. The displayed structures are based on x-ray crystallographic data except for hMeOp that has been produced via comparative modelling (see Methods section). Ionizable side-chains are marked with a \*: acidic residues are shown in red and alkaline in blue. Non-ionizable polar residues are shown in green. B. General structure of the QM/MM model constructed using ARM displaying the environment (in grey with external counterions in blue and red), cavity (residues in red) and Lys-QM systems (in blue and green).

## Methods

### Overview of ARM

The QM/MM model constructed by ARM is a gas-phase and globally uncharged monomer model whose structure is illustrated in Fig. 1B. The model is divided in the three systems called environment, cavity and Lys-QM. The environment has backbone and side-chain atoms fixed at the crystallographic (or comparative model) structure, and incorporates suitably placed external counterions also kept fixed at computed positions (see the details below). The cavity, which includes the Lys-QM system detailed in the left part of Scheme 1, also features fixed backbone atoms but its side-chain atoms are free to relax. The Lys-QM system comprises the atoms of the Lys side-chain in contact (through C<sub>δ</sub>) with the QM/MM frontier and the entire QM subsystem that corresponds to an N-methylated retinal chromophore. All Lys-QM atoms are free to relax. The environment, cavity and Lys side-chain form the MM subsystem. Notice that while ARM models feature a fixed environment system (see Fig. 1B), other bare QM/MM rhodopsin models reported in the literature are based on geometrically optimized apoprotein structures.<sup>13,14,55,56</sup> In this respect ARM models, are

1  
2  
3 “conservative” in the sense that they ensure that the information from the x-ray crystallographic  
4 structure is retained in the final QM/MM model.  
5

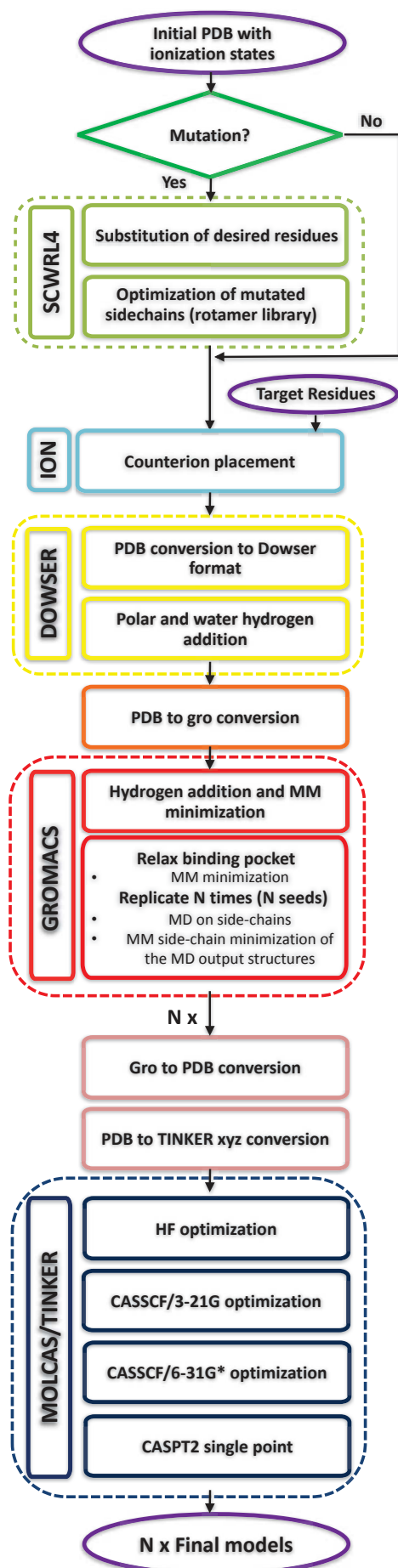
6  
7 Technically, ARM is a Linux command-line tool that interfaces various publicly available or  
8 specifically designed (e.g. the ION module) computational chemistry programs (see the SI for  
9 detailed information) carrying out the different tasks required to construct QM/MM models. The  
10 only mandatory input includes a Protein Data Bank (PDB) file containing an initial (guess) rhodopsin  
11 structure with assigned ionization states, and a list of target residues required to drive the placement  
12 of external counterions. However, the original PDB may be pre-processed by completing residues  
13 with missing atoms (missing residues are not added) using the specific utility (Automatic PSF Builder)  
14 in the VMD program.<sup>57</sup> The sequence of steps from the input to the final results is pre-defined, and  
15 either automatic or automatable in a future version of the script (see the *Present automation limits*  
16 subsection in the Results and Discussion section).  
17

18  
19 The currently established overall ARM workflow is depicted in Fig. 2. The workflow starts by reading  
20 the PDB file mentioned above. Then a series of automatic manipulations are performed including  
21 (sequentially): external counterion placement, hydrogen atom placement, MM energy minimization,  
22 multiple MD relaxations and lower-level QM/MM calculations, which are carried out to prepare the  
23 input for multiple QM/MM geometry optimizations. The ION module (see below and the SI) adds in  
24 specific positions the external Na<sup>+</sup> and Cl<sup>-</sup> ions required to neutralize the model, by independently  
25 neutralizing the inner and outer (i.e. intra-cellular and extra-cellular) protein surfaces. This would  
26 better mimic the situation of the protein in a micelle rather than in a biological membrane where  
27 one has a different ion concentration in the inner and outer protein surfaces. These two surfaces are  
28 defined by a list of target residues provided by the user (see Fig. 2). The list include a sets of (solvent  
29 exposed) charged residues located above and below the chromophore, which may potentially harbor  
30 one or more counterions. Instead, the exact position of the counterions is determined by an energy  
31 minimization procedure. More specifically, for each surface, the list of target residues usually  
32 includes all residues whose charges have the same sign as the corresponding surface net charge. An  
33 exception is the retinal-bonded lysine and its counterion, which may be excluded in order to prevent  
34 the placement of ions inside the protein. Below, we refer to the chosen neutralization of the inner  
35 and outer protein surface scheme as the No Surface Charge (NSC) scheme (see also section 2.4 in the  
36 SI). Thus, as also discussed below the NSC scheme defines the *distribution* of counterions between  
37 the two surfaces while the ION module determines in which exact *position* they have to be placed.  
38 Throughout the MM energy minimization and MD calculations, the AMBER94<sup>58</sup> force field is used.  
39 Firstly, the hydrogen atoms are added and their positions optimized at the MM level. Consistently  
40 with the workflow in Fig. 2, the hydrogen atoms bound to water and polar atoms are added by  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 DOWSER, a program for hydrogen addition according to energy criteria,<sup>59</sup> while TINKER<sup>60</sup> is used  
4 elsewhere. Since water and hydrogen bond networks (HBN) affect side chain conformations and  
5 long-range electrostatics (thus ultimately modifying spectral and photochemical properties<sup>6,61-63</sup>), it is  
6 important to use DOWSER quantitative criteria for hydrogen and water relocation. After the MM  
7 cavity atoms have been automatically selected using CASTp<sup>64</sup> online server (with default settings), a  
8 single MM geometry optimization and N independent (i.e. starting with N different seeds which  
9 provide N independent sets of initial velocities. See also the SI) MD room-temperature relaxations  
10 are performed at the MM level using GROMACS<sup>65</sup> on the cavity and the Lys-QM systems. The  
11 chromophore non-bonding and bonding interactions are modeled according to AMBER94 rules: van  
12 der Waals parameters for retinal are taken from our custom AMBER94 parameters set<sup>50,66</sup>; partial  
13 charges were calculated as AMBER-like RESP charges.<sup>58</sup> Distinct sets of parametrized charges are  
14 used for different chromophore isomers. The MD relaxations consist of a 50 ps heating followed by  
15 150 ps room-temperature equilibration and 800 ps production for a total of 1 ns. The MD output  
16 structures constitute the guess for N corresponding QM/MM calculations<sup>67</sup> performed according to  
17 the scheme described in ref. 50. Briefly, the full N-methyl retinal chromophore (53 atoms)  
18 corresponding to the QM subsystem is connected with the designated lysine side-chain of the MM  
19 subsystem (taken together these form the aforementioned Lys-QM system). The Lys residue linked to  
20 the retinal chromophore is automatically identified on the basis of the distance between the  
21 chromophore C15 atom (see Scheme 1) and the Lys residue N atoms in the input structure.

22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36 The QM/MM frontier is treated within a link atom approach,<sup>68</sup> whose position is restrained according  
37 to the Morokuma scheme,<sup>69</sup> and it is placed across the lysine C $\delta$ -C $\epsilon$  bond (where C $\epsilon$  is a QM atom).  
38 The lysine charges are modified by setting the C $\delta$  charge to zero, to avoiding hyperpolarization and  
39 redistributing the residual fractional charge on the most electronegative atoms of the lysine, thus  
40 ensuring a +1 integer charge of the Lys-QM system. A QM/MM geometry optimization carried out at  
41 the HF/3-21G/MM level completes the preparatory phase. In this and subsequent QM/MM  
42 computations the AMBER94<sup>58</sup> MM force field is used.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



**Figure 2** Automatic Rhodopsin Modeling workflow. Vertical boxes refer to programs and their associated actions are circled with dashed lines and colored in light green (SCWRL4), light blue (ION), yellow (DOWSER), red (GROMACS) and blue (MOLCAS/TINKER). Intermediate steps corresponding to file format conversions are shown in orange (PDB to GROMACS type) and pink (GROMACS to PDB, PDB to TINKER XYZ format) boxes.

The obtained  $N$  independent HF/3-21G/MM models are processed via two sequential QM/MM optimizations at the single-root CASSCF(12,12)/3-21G/MM level and CASSCF(12,12)/6-31G(d)/MM level to get  $N$  final structures. These sequential optimizations have the scope to achieve a more rapid convergence of both the molecular orbital and geometry. Subsequently a CASPT2(12,12)/6-31G(d)/MM 3-single root single point calculation is performed for each model, taking a 3-roots State Average (SA)<sup>70</sup> CASSCF(12,12)/6-31G(d) wavefunction as a reference. Finally,  $N$  vertical excitation energy ( $\Delta E_{S1-S0}$ ) values are computed as the difference between the first two roots. The average  $\Delta E_{S1-S0}$  value is then compared with the corresponding experimentally observed value (i.e. we assume that this value corresponds to the energy carried by a photon with the observed  $\lambda_{max}^a$ ). Oscillator strengths ( $f$ ) are computed at the SA-CASSCF level to check if  $\Delta E_{S1-S0}$  corresponds to a spectroscopically allowed transition. All QM/MM calculations are performed by the distributed MOLCAS/TINKER interface<sup>67</sup>, as implemented in MOLCAS version 7.4 or higher (ESPF module),<sup>71</sup> interfaced with TINKER version 4.2 or higher.<sup>60</sup> ARM can optionally perform amino acid substitutions on the starting PDB structure to generate ARM models of rhodopsin mutants. The list of desired mutations must be provided as an additional input file (see the SI for further details). Mutations are carried out by SCWRL4,<sup>72</sup> a

1  
2  
3 program for predicting side-chain conformations from a given protein backbone, using a backbone-  
4 dependent rotamer library. Once the mutations have been made, the ARM workflow proceeds  
5 exactly as described above.  
6  
7

8 While the results presented below are all consistently produced with *the same* standardized ARM  
9 workflow using pre-assigned *default* parameter values (see points a-f below), the following  
10 customization is possible: (a) DOWSER can retain the crystallographic water molecules in the starting  
11 guess structure (default) or neglect them. In the latter case, it generates and places water molecules  
12 in all cavities of suitable sizes, (b) the default number (N=10) of independent MD relaxations as well  
13 as the default heating (50 ps), room temperature (298 K) pre-equilibration (150 ps) and production  
14 (800 ps) stages may be changed, (c) the default chromophore cavity is obtained automatically using  
15 the CASTp<sup>64</sup> online server which identifies cavities and pockets in proteins by geometrical methods  
16 (we used a probe radius set at 1.4 Å). This cavity may be replaced by cavities constructed  
17 automatically specifying the maximum distance of any atom of the protein residues from any atom of  
18 the chromophore using the program VMD<sup>57</sup>. Alternatively, a list of residues whose side-chains will be  
19 relaxed can be provided using an additional input file. (d) The chromophore can be either allowed to  
20 relax (default) or kept frozen during the MD relaxations and (e) the final MD structures can be  
21 obtained by choosing either the last snapshot of the run, or the snapshot featuring the most similar  
22 structure to the unphysical average structure obtained from the MD production stage (nearest-to-  
23 average structure). Such similarity is evaluated as the Root Mean Square Deviation (RMSD) value of  
24 the snapshot against the average structure, so that the most similar snapshot has the lowest RMSD.  
25 The nearest-to-average structure is the default. (f) The QM/MM calculations are not customizable,  
26 except for the use of microiterations<sup>73</sup> to improve the convergence: during each geometry  
27 optimization step the cavity side-chains are relaxed after the optimization of the Lys-QM system by  
28 default, unless they are disabled by the user.  
29  
30

31 The above procedure is automatic once the requested input is provided by the user, who has to type  
32 values and choices and run scripts following on-screen instructions. TINKER 5.1,<sup>60</sup> GROMACS 4.5.4<sup>65</sup>  
33 and MOLCAS 7.8<sup>67</sup> were used throughout the testing. The final outcome is a set of N TINKER XYZ files  
34 containing the protein geometries, which can be converted to PDB by running a provided script. For  
35 each TINKER file, two files with information such as the CASPT2 absolute energies in a.u. of the  
36 computed electronic states, the associated vertical excitation energies expressed in both kcal/mol  
37 and nm and the oscillator strengths for the same transitions, are produced.  
38  
39  
40  
41  
42  
43  
44

45 *Input data for the benchmark set.*  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

58 As mentioned above the construction of an ARM model requires an initial geometry and the relative  
59  
60

1  
2  
3 residue protonation states. The geometry is provided by PDB crystallographic structures for Rh  
4 (1U19<sup>35</sup>), bathoRh (2G87<sup>42</sup>), ASR (1XIO<sup>34</sup>), SqR (2Z73<sup>36</sup>), bR<sub>LA</sub> (1C3W<sup>38</sup>), bR<sub>DA</sub>(1XOS<sup>39</sup>), PR (4JQ6<sup>40</sup>) and  
5 C1C2 (3UG9<sup>74</sup>), while for hMeOp the comparative model of ref. 9, is employed. The residue  
6 protonation states are instead determined, in all cases, using the program PROPKA version 3.0<sup>75</sup> for  
7 all members of the benchmark set including the ASR<sub>AT</sub>, ASR<sub>13C</sub> and Rh mutants. The target residues  
8 required by the ION module for placing the chloride and sodium counterions on the two rhodopsin  
9 surfaces (as requested by the NSC scheme. See also the description above) are unambiguously  
10 selected according to energy criteria.

11  
12 The list of residues defining the retinal chromophore cavity, the position of the water oxygens and  
13 the set up of the MD relaxation, are considered pre-assigned parameters. The positions of the  
14 crystallographic water oxygens were used as a starting guess for the DOWSER-3<sup>59</sup> calculation  
15 assigning the initial orientation of all water hydrogen atoms (in the initial hMeOp geometry obtained  
16 via comparative modeling we included the water oxygens present in its SqR template for  
17 consistency).

#### 28 29 *Sensitivity of ARM models to the protonation states.*

30  
31 As already mentioned above the assignment of the protonation states of the ionizable residues is  
32 part of the ARM input. When a different ionization state is assigned to a cavity (or near cavity) side-  
33 chain, significant changes in the  $\Delta E_{S1-S0}$  values are expected. In the Result and Discussion section we  
34 investigate this effect in two cases: the E181 side-chain of Rh (Fig. 1A), whose ionization state is still  
35 debated;<sup>76-78</sup> and the D85 and D212 residues in bR, whose complex hydrogen bond network allows  
36 for different combinations of protonation states.<sup>79</sup> In the case of Rh, E181 is predicted to be neutral  
37 by PROPKA, which gives an estimated pKa>7.6, corresponding to >80% neutral molecules in a sample  
38 at pH 7. Therefore we constructed a model with an ionized E181 side-chain (labeled Rh<sub>181</sub>) by using  
39 the same initial geometry as for the reference model (Rh) so that the effect of the charged E181  
40 could be singled out. For bR we used the PROPKA prediction giving a neutral D85 and a charged  
41 D212. Nevertheless, PROPKA labels them as “correlated” residues, i.e. the choice of ionization state  
42 on one residue determines the other, mostly because they are only about 6 Å apart. Also, D85 and  
43 D212 are connected via a water molecule, so protons can, in principle, shuffle between them.  
44 However, since D212 makes hydrogen bonds with Y57 and Y185, a protonated state seems unlikely.  
45 On the other hand, D85 can accommodate a proton but the presence of many polar residues and  
46 water molecules would stabilize a charged form too. For this reason we constructed an additional  
47 model with a charged D85 side chain (bR<sub>85</sub>, see the subsection *Alternative assignments of ionization*  
48 *states: Rh and bR* below), while keeping D212 charged as in the reference ARM model.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5 *Sensitivity of ARM models to the neutralization scheme.*  
6

7 The sensitivity of the ARM models to the neutralization scheme has also been assessed because the  
8 standard NSC neutralization scheme does not necessarily reflect the physiological conditions (i.e. in  
9 general, membrane proteins are exposed to trans-membrane electrostatic fields in the range of few  
10 tens of mV originating from an asymmetrical distribution of the surface ions). As we report below  
11 (see the *Assessing the protein environment (counterion distribution and position) approximation*  
12 subsection below), we look at the  $\Delta E_{S1-S0}$  variations obtained when, instead of neutral inner and  
13 outer surfaces, external  $\text{Cl}^-$  and  $\text{Na}^+$  counterions are placed in such a way that oppositely charged  
14 surfaces are obtained, while still achieving an overall neutral model. In the latter case, an  
15 electrostatic field across the chromophore is generated. This effect is tested for Rh, Rh<sub>181</sub> and the  
16 bacterial rhodopsin isomer ASR<sub>AT</sub>. Below we also documented the  $\Delta E_{S1-S0}$  variations between a  
17 manual (i.e. biased) external counterion placement scheme (see the SI for details) and the automatic  
18 placement obtained with the ION module for Rh, ASR, SqR, hMeOp, bR, PR, bathoRh.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29 *Sensitivity of ARM models to the residue side-chain conformation.*  
30

31 Residue side chain conformations affecting the HBN in the vicinity of the chromophore could  
32 significantly affect the computed  $\Delta E_{S1-S0}$  value in rhodopsins. Thus the effect of the flexibility of the  
33 side-chain conformation and HBN stability has been assessed (see the subsection *Stability of the*  
34 *computed side-chain conformations and hydrogen bond network* below) as a function of: (i) the  
35 number N of replicated 1 ns MD relaxations for Rh and ASR and (ii) the length (up to 30 ns) of a  
36 single MD relaxation for Rh and Rh<sub>181</sub>. Such significantly longer MD relaxation is used to search for  
37 instabilities or alternate conformations that are unreachable during single 1 ns relaxation employed  
38 when building ARM models and to estimate their effects on computed  $\Delta E_{S1-S0}$  values.  
39  
40  
41  
42  
43  
44  
45

46 *Preparation and  $\lambda_{max}^a$  measurements of ASR mutants.*  
47

48 In order to expand the benchmarking of mutants beyond the set of Rh mutants found in the  
49 literature, we have prepared 5 ASR mutants. The corresponding  $\lambda_{max}^a$  values have been determined  
50 via HPLC analysis and spectral measurements for both the ASR<sub>AT</sub> and ASR<sub>13C</sub> isomers providing 10  
51 experimentally observed values to be used in the benchmark study. Full-length ASRs having six  
52 histidine residues at the C terminus were expressed in E. coli. BL21, C41(DE3) or UT5600 strains. The  
53 protein was solubilized in n-dodecyl- $\beta$ -D-maltoside (DDM), and purified via  $\text{Ni}^{2+}$ -affinity column  
54 chromatography in manner previously described.<sup>33,80,81</sup> The  $\lambda_{max}^a$  of ASR<sub>AT</sub> and ASR<sub>13C</sub> forms were  
55 determined from the UV-vis spectroscopy and the HPLC analysis on dark- and light-adapted sample  
56  
57  
58  
59  
60

1  
2  
3 at pH 7, according to the method described previously.<sup>82</sup> For these measurements, DDM solubilized  
4 sample were used. The detailed results are reported in the SI.  
5  
6

### 7 *Present automation limits.*

8  
9 The automation of ARM is presently not complete. First of all the assignment of the amino acid  
10 ionization states, including the generation of models of the same rhodopsin with different possible  
11 ionizations and ionization equilibria (e.g. a possible equilibrium between Rh and Rh<sub>181</sub>), are manually  
12 handled after inspection of the PROPKA output or, when available, on the basis of experimental  
13 data. Secondly, as described above, the target residues of the ION module are manually provided for  
14 the outer and inner protein surfaces even if such a list could be automatically generated on the basis  
15 of the ionization state information and the geometrical parameters of the residue. Thirdly, the list of  
16 the CASTp cavity residues is presently generated using the CASTp<sup>64</sup> online server whose output is  
17 pasted in a command line. While the previous limitations require only suitable coding to be removed  
18 a fourth limitation requires more serious efforts to be overcome. This is related to the selection of  
19 the CASSCF active space. In all benchmark models illustrated here the HF/3-21G optimization (see  
20 Fig. 1) yields the 12  $\pi$  molecular orbitals in the HOMO-5 to LUMO+5 position. While this  
21 automatically leads to a CASSCF correct active space selection (i.e. no orbital permutation is  
22 necessary and the HF molecular orbitals can be directly used as molecular orbital guess for the  
23 successive CASSCF calculation), such favorable ordering must be considered an exception. In other  
24 words, in different rhodopsin or mutant models, the HOMO-5 to LUMO+5 positions may contain  
25 orbitals not belonging to the skeletal  $\pi$ -system. An automatic way of detecting this situation, based  
26 on CASSCF orbital occupations, is implemented in ARM. A wrong orbital is likely to have either 2 or 0  
27 (larger than 1.999 or lower than 0.001) orbital occupation rather than a value in between. When this  
28 happens, it is necessary to visually inspect the orbitals and permute their order before the CASSCF  
29 calculation is restarted, thus making the workflow not automatic. However, a research line on the  
30 automatic selection of the CASSCF active space is presently carried on<sup>83</sup>, making the possibility of a  
31 robust automation of ARM realistic.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49

## 50 **Results and Discussion**

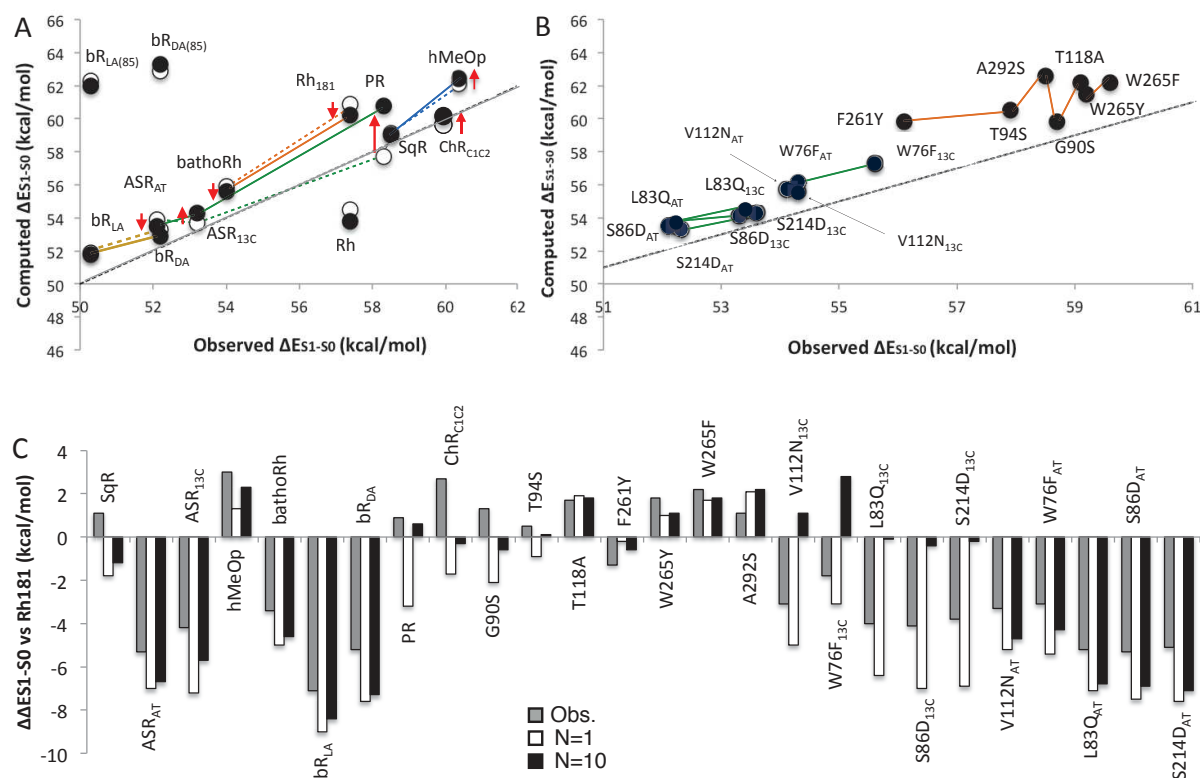
51  
52 The average (N=10)  $\Delta E_{S1-S0}$  values calculated with ARM using the default parameters defined above  
53 are reported in Table 1 (wild-type rhodopsins) and Table 2 (mutants) and plotted as full circles in Fig.  
54 3A and 3B where they are compared to the experimental values. Since the observed  $\lambda_{max}^a$  is not well  
55 defined for hMeOp, the 473 nm value (i.e. a  $\Delta E_{S1-S0}$  value of 60.4 kcal/mol) corresponding to the  
56 average of the observed range (467-480 nm)<sup>9,84,85</sup> was taken as a reference. All input structures (i.e.  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

crystallographic structures and the hMeOp comparative model with assigned ionization states and external counterions) are provided as SI. As part of the SI we also provide one selected output structure for each wild-type rhodopsin and rhodopsin mutant, which corresponds to the model with the  $\Delta E_{S1-S0}$  closest to the average  $\Delta E_{S1-S0}$ .

The overall standard deviation of the average  $\Delta E_{S1-S0}$  values with respect to observed values (connected by a full line in Fig. 3) and when excluding the Rh, bR<sub>DA(85)</sub> and bR<sub>LA(85)</sub> models is 3.6 kcal/mol with the absolute deviations (the difference between the calculated and observed  $\Delta E_{S1-S0}$  values) ranging from -0.2 to +4.1 kcal/mol. When considering the wild-type proteins only, the standard deviation increases to 3.8 kcal/mol with the largest value found for the Rh<sub>181</sub> (+2.8 kcal/mol). The trend, which parallels the observed values, and blue-shifted deviation of the computed values indicate the presence of a systematic error of few kcal/mol whose exact origin is difficult to characterize. Notice that the MS-CASPT2 level of theory yield the same trend (see Table S6 in the SI) as the CASPT2 level even if the excitation energies appear further blue-shifted. Such behaviour is attributed to an oversplitting of the electronic energies in the presence of state mixing (i.e. two large off-diagonal elements in the MS-CASPT2 hamiltonian especially when using small atomic basis) already described in the literature.<sup>86,87</sup>

As it will be discussed later, using the average of ten  $\Delta E_{S1-S0}$  values generated from uncorrelated MD relaxations rather than using a single  $\Delta E_{S1-S0}$  value, decreases the  $\Delta E_{S1-S0}$  deviations. Indeed, the comparison of full (N=10) and open (N=1) circles in Fig. 3A shows a general improvement in the first set with respect to the second. One extreme case is that of PR which pass from a negative deviation (-0.6) to a positively deviation (+2.5) when considering the  $\Delta E_{S1-S0}$  average value. A similar but less dramatic effect is seen for ASR<sub>13C</sub>. Furthermore, in the case of bovine rhodopsin, the N=10 averaging moves the Rh<sub>181</sub> model in better agreement with the general trend with respect to the Rh model. This is why in this work the Rh<sub>181</sub> model is selected as the reference model of bovine rhodopsin. While we cannot presently exclude the presence of an equilibrium between Rh and Rh<sub>181</sub> or that our results are a consequence of the specific ARM model setup, a further investigation of these issues is beyond the scope of the present work mainly focusing on the accuracy of computed  $\Delta E_{S1-S0}$  trends and not specific  $\Delta E_{S1-S0}$  absolute values.



**Figure 3** Benchmark results. (A) Observed vs. computed values for vertical excitation energies of wild-type rhodopsins. The diagram includes the results of sensitivity tests for input parameters, such as ionization state (for  $Rh_{181}$ ,  $bR_{LA}$  and  $bR_{DA}$ ), HBN variations and side-chain conformational changes. The effects of HBN variations and conformational changes are assessed by comparison between average and single  $\Delta E_{51-50}$  values (see text). Full colored lines connect the average  $\Delta E_{51-50}$  values obtained from N=10 ARM models generated using independent MD runs. The dashed colored lines connect the single  $\Delta E_{51-50}$  values obtained with a single N=1 ARM model. The colors refer to vertebrate rhodopsin (orange), invertebrate and non-visual rhodopsins (blue), bacterial rhodopsins (green) and archaea rhodopsin (brownish). Channelrhodopsin  $ChR_{C1C2}$  is the only rhodopsin originating from microbial eukaryotic rhodopsins. The red arrows show the  $\Delta E_{51-50}$  changes when passing from single (open circles) to average (full circles) values. (B) Comparison between computed and observed average  $\Delta E_{51-50}$  values (full circles) from N=10 ARM models of a set of  $ASR_{AT}$ ,  $ASR_{13C}$  (ten values on the left part of the diagram) and  $Rh_{181}$  (seven values on the right part of the diagram) mutants. (C) Relative observed and computed  $\Delta E_{51-50}$  changes ( $\Delta\Delta E_{51-50}$ ) with respect to  $Rh_{181}$ . Grey bars are observed values, white bars correspond to  $\Delta\Delta E_{51-50}$  values computed using single  $\Delta E_{51-50}$  values while black bars correspond to  $\Delta\Delta E_{51-50}$  values computed using average  $\Delta E_{51-50}$  values.

**Table 1.** Comparison between computed and observed vertical excitation energies  $\Delta E_{51-50}$  in kcal/mol (maximum absorption wavelengths  $\lambda_{max}^a$  in nm) for our benchmark set of 13 wild-type ARM models computing single and multiple calculations. Computed oscillator strengths  $f_{Osc}$  are also shown. The corresponding MS-CASPT2 values of the relevant models are given in the SI.

	Obs. $\Delta E$ ( $\lambda_{max}^a$ )	Calc. $\Delta E$ ( $\lambda_{max}^a$ ) Single	Calc. $\Delta E$ ( $\lambda_{max}^a$ ) Average (N=10)	Error Single	Error Average (N=10)	$f_{Osc}$ Average (N=10)
<b>Rh</b>	57.4 (498) <sup>81</sup>	54.5 (524)	53.8 (531)	-2.9	-3.6	1.08
<b>Rh<sub>181</sub></b>	57.4 (498) <sup>81</sup>	60.9 (469)	60.2 (474)	3.5	2.8	0.99



<b>bathoRh</b>	54.0 (529) <sup>82</sup>	55.9 (511)	55.6 (513)	1.9	1.6	1.04
<b>SqR</b>	58.5 (489) <sup>83</sup>	59.1 (484)	59.0 (484)	0.6	0.5	0.81
<b>hMeOp</b>	60.4 (473) <sup>a</sup>	62.2 (459)	62.5 (457)	1.8	2.1	0.78
<b>ASR<sub>AT</sub></b>	52.1 (549) <sup>34</sup>	53.9 (529)	53.5 (534)	1.8	1.4	1.14
<b>ASR<sub>13C</sub></b>	53.2 (537) <sup>34</sup>	53.7 (532)	54.5 (526)	0.5	1.3	1.04
<b>bR<sub>LA</sub></b>	50.3 (568) <sup>84</sup>	51.9 (550)	51.8 (551)	1.6	1.5	1.30
<b>bR<sub>LA(85)</sub></b>	50.3 (568) <sup>84</sup>	62.3 (458)	62.0 (461)	12.0	11.7	0.6
<b>bR<sub>DA</sub></b>	52.2 (548) <sup>84</sup>	53.3 (535)	52.9 (540)	1.1	0.7	0.90
<b>bR<sub>DA(85)</sub></b>	52.2 (548) <sup>84</sup>	62.9 (454)	63.3 (451)	10.7	11.1	0.74
<b>PR</b>	58.3 (490) <sup>85</sup>	57.7 (495)	60.8 (470)	-0.6	2.5	0.88
<b>ChR<sub>C1C2</sub></b>	60.1 (476) <sup>b</sup>	59.2 (483)	59.9 (477)	-0.9	-0.2	0.95

<sup>a</sup> Average of available values from refs. 84, 85

<sup>b</sup> Available value is provided in ref. 74

**Table 2.** Comparison between computed and observed vertical excitation energies  $\Delta E_{S1-S0}$  in kcal/mol (maximum absorption wavelengths  $\lambda_{max}^a$  in nm) for our benchmark set of ARM models for 17 ASR<sub>13C</sub>, ASR<sub>AT</sub> and Rh<sub>181</sub> mutants, computing single and multiple calculations. Computed oscillator strengths  $f_{Osc}$  are also shown.

	Obs. $\Delta E$ ( $\lambda_{max}^a$ )	Calc. $\Delta E$ ( $\lambda_{max}^a$ ) Single	Calc. $\Delta E$ ( $\lambda_{max}^a$ ) Average (N=10)	Error Single	Error Average (N=10)	$f_{Osc}$ Average (N=10)
<b>G90S</b>	58.7 (487) <sup>86</sup>	58.8 (485)	59.8 (478)	0.1	1.1	0.95
<b>T94S</b>	57.9 (494) <sup>87</sup>	60.0 (475)	60.5 (472)	2.1	2.6	0.86
<b>T118A</b>	59.1 (484) <sup>86</sup>	62.8 (454)	62.2 (459)	3.7	3.1	0.78
<b>F261Y</b>	56.1 (510) <sup>88</sup>	60.7 (470)	59.8 (477)	4.6	3.7	0.91
<b>W265F</b>	59.6 (480) <sup>89</sup>	62.6 (456)	62.2 (460)	3.0	2.6	0.89
<b>W265Y</b>	59.2 (483) <sup>89</sup>	61.9 (461)	61.5 (464)	2.7	2.3	0.87
<b>A292S</b>	58.5 (489) <sup>86</sup>	63.0 (453)	62.6 (457)	4.5	4.1	0.83
<b>V112N<sub>13C</sub></b>	54.3 (526) <sup>a</sup>	54.5 (524)	55.6 (514)	0.2	1.3	1.07
<b>W76F<sub>13C</sub></b>	55.6 (514) <sup>a</sup>	57.0 (501)	57.3 (498)	1.4	1.7	0.96
<b>L83Q<sub>13C</sub></b>	53.4 (535) <sup>a</sup>	54.1 (528)	54.4 (524)	0.7	1.0	0.86
<b>S86D<sub>13C</sub></b>	53.3 (536) <sup>a</sup>	54.0 (529)	54.1 (528)	0.7	0.8	1.01
<b>S214D<sub>13C</sub></b>	53.3 (536) <sup>a</sup>	54.1 (528)	54.3 (526)	0.5	0.7	0.99
<b>V112N<sub>AT</sub></b>	54.1 (529) <sup>a</sup>	55.7 (512)	55.7 (512)	1.6	1.6	1.18
<b>W76F<sub>AT</sub></b>	54.3 (527) <sup>a</sup>	56.6 (504)	56.1 (509)	2.3	1.8	0.98
<b>L83Q<sub>AT</sub></b>	52.2 (548) <sup>a</sup>	53.2 (536)	53.6 (533)	1.0	1.4	1.17
<b>S86D<sub>AT</sub></b>	52.1 (549) <sup>a</sup>	54.0 (529)	53.5 (533)	1.9	1.4	1.15
<b>S214D<sub>AT</sub></b>	52.0 (550) <sup>a</sup>	53.6 (532)	53.3 (536)	1.3	1.0	1.16

<sup>a</sup> Present work. See the SI for details.

In principle, the set of mutants would be more challenging with respect to wild-type models due to the smaller changes in the observed  $\Delta E_{S1-S0}$  values and the lack of crystallographic structures. Indeed, the deviations of the computed average  $\Delta E_{S1-S0}$  values for ASR and Rh<sub>181</sub> mutant set are 0.7-4.1 kcal/mol and therefore higher than those of the -0.2-2.8 kcal/mol of the wild-type set. In fact, upon replacement of a side chain, the local steric hindrance and polarity are perturbed upon substitution

1  
2  
3 and this induced a re-organization of the nearby residues that may not be accurately described by  
4 the SCWRL4 module and limited CASTp cavity relaxation. On the other hand, as shown in Fig. 3C, the  
5 sign of the experimentally observed  $\Delta\Delta E_{S1-50}$  are always qualitatively reproduced, with the only  
6 exception of G90S, when using  $\Delta\Delta E_{S1-50}$  computed on the basis of average  $\Delta E_{S1-50}$  values. With single  
7  $\Delta E_{S1-50}$  values both G90S and T94S Rh<sub>181</sub> mutants do not provide the correct shift direction further  
8 indicating the importance of the sampling on the HBN and side-chain conformations. However, these  
9 two “less accurate” mutant models are different from the other models since their CASTp cavity  
10 does not incorporate the mutated side-chains. Accordingly, these side-chains are kept frozen during  
11 the MD and QM/MM optimization runs and not relaxed as for the mutated side-chains of the other  
12 members of the set. However, inclusion of the 90 and 94 side chains in the chromophore cavities  
13 (which therefore does not correspond to the default CASTp cavity) does not change the trend when  
14 compared with a wild-type structure with the same expanded CASTp cavity.

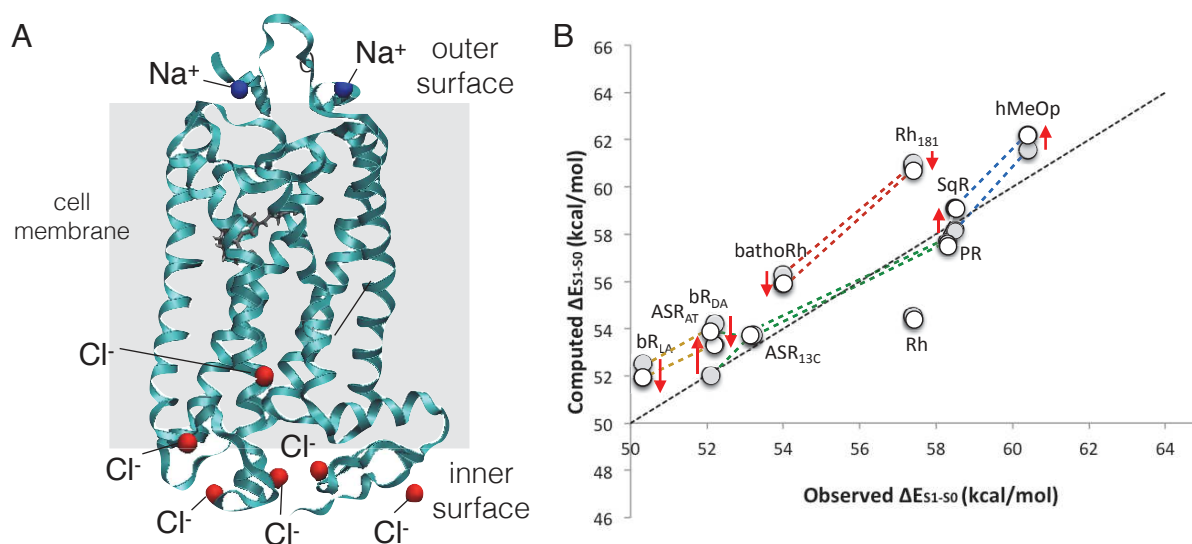
15  
16  
17  
18  
19  
20  
21  
22  
23  
24 The ability of ARM to reproduce/predict the observed  $\Delta E_{S1-50}$  values may appear surprising when  
25 considering the relatively small group of 27 values (i.e. excluding the Rh, bR<sub>LA(85)</sub> and bR<sub>DA(85)</sub>  
26 discussed below), the general uncertainty associated with the generation of comparative (i.e.  
27 hMeOp) and mutant models and, of course, the basic models generated by ARM. Furthermore, we  
28 expect the computed values to be affected by systematic error cancellations. An error of this type  
29 has been already mentioned above and occurs when the multiconfigurational CASPT2//CASSCF/6-  
30 31G\* strategy is used and it is exclusively limited to the QM part of the computation.<sup>87</sup> However,  
31 what is most important here is the ability to reproduce the observed  $\Delta E_{S1-50}$  trends (i.e. the sign and  
32 magnitude of the  $\Delta\Delta E_{S1-50}$  values). Such ability appears to be satisfactory (see Fig. 3C) especially  
33 when comparing rhodopsin originating in similar organisms and therefore with higher homology (see  
34 orange, blue, green and brownish full lines in Fig. 3A and 3B).

#### 45 *Assessing the protein environment (counterion distribution and position) approximation.*

46 While rhodopsins are trans-membrane proteins, ARM models do not incorporate the environment  
47 provided by the complex hydrated membrane bilayer and/or account for the existence of protein  
48 dimers or trimers. The effect of the protein environment is indirectly incorporated in the model at  
49 the input level by using the experimentally derived structure to build the fixed part of the protein  
50 body (see Fig. 1B), by assuming a total charge of zero and, according to our NSC reference  
51 neutralization scheme (see Method section), assuming that the solvent molecule distribution on  
52 both the inner and outer surfaces screens the surface charged residues and counterions. For  
53 instance, since Rh has a +4 total charge after the assignment of the residue protonation states (i.e. at  
54 pH=7), six Cl<sup>-</sup> are distributed on the intracellular surface, where there are 16 positively and 10  
55  
56  
57  
58  
59  
60

negatively charged residues, and two  $\text{Na}^+$  on the extracellular side, where there are 6 positively and 8 negatively charged residues (see Rh structure in Fig. 4A). This procedure results in a neutral rhodopsin model with neutral (0 charge) inner and outer surfaces.

The NSC scheme defines how the counterions shall be distributed between the two protein surfaces. However, the exact position of each counterion is decided automatically by the ION module, which employs energy criteria (see the SI for details). The sensitivity to the counterion position is assessed in Fig. 4B, where the single ( $N=1$ )  $\Delta E_{S1-50}$  values computed using ION are compared to the  $\Delta E_{S1-50}$  values computed by placing the counterions manually according to operator defined rules (i.e. as detailed in the SI,  $\text{Cl}^-$  or  $\text{Na}^+$  are placed at qualitatively predefined distances and orientations close to the positive and negative residues respectively). The comparison shows an improvement in the trend of computed  $\Delta E_{S1-50}$  values for the set of wild-type rhodopsins of Fig. 3A when using the automatic placement, with the exception of the red-shifted Rh and PR values, whose deviations are discussed in the following sections dealing with ionization states and thermal sampling.



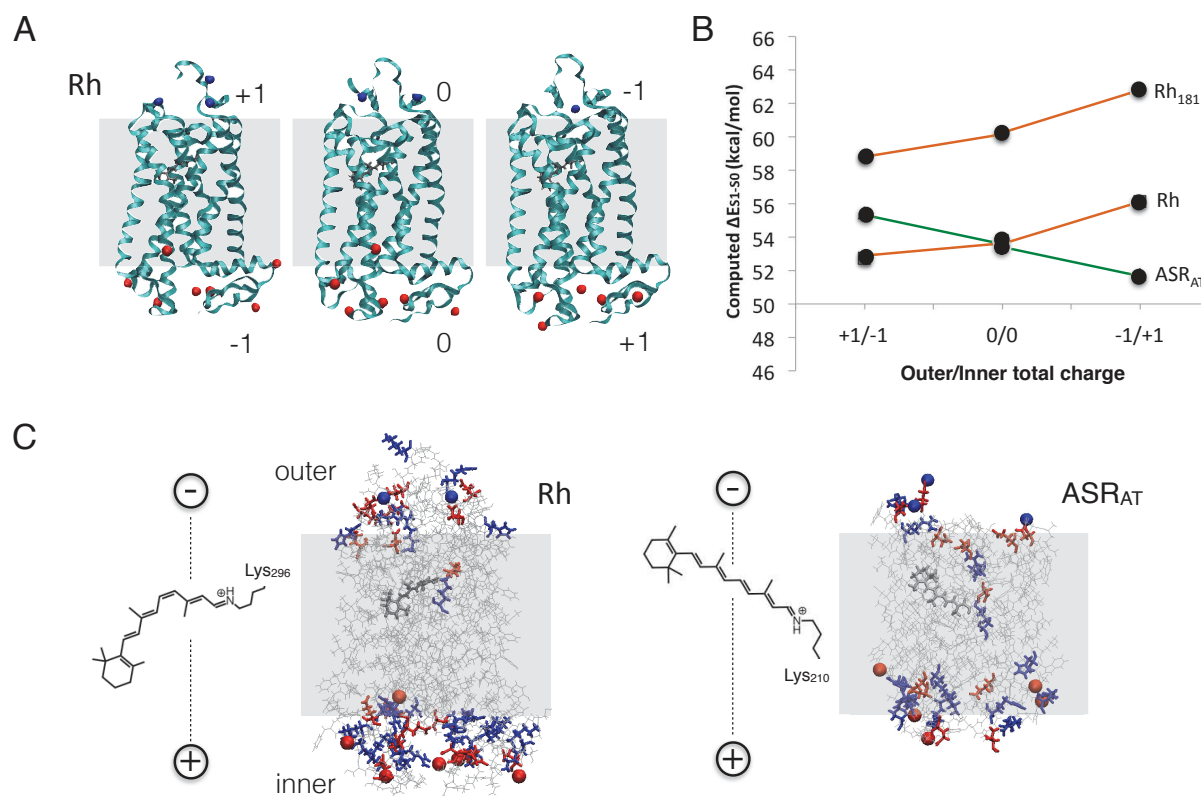
**Figure 4** Features affecting the vertical excitation energy in ARM rhodopsin models. (A) ION placement of the  $\text{Na}^+$  and  $\text{Cl}^-$  counterions in Rh. (B) Differences between the  $\Delta E_{S1-50}$  values obtained via a manual placement (see the SI) and the ION placement. The red arrows show the  $\Delta E_{S1-50}$  changes when passing from a manual (dashed circles) to an automatic (open circles) placement.

According to the NSC scheme the counterions are distributed in such a way to yield uncharged surfaces. However, in all rhodopsins the sum of the ionizable residue charges is more positive in the inner surface and more negative (less positive) in the outer surface (due to the physiological excess of extra-cellular positive ions of cell membranes compared to their intra-cellular one). This unbalanced situation creates an electrostatic field (a voltage) across the membrane, which ultimately results from differently charged inner and outer protein surfaces. For this reason, and without targeting a simulation of the physiological state, we test the sensitivity of ARM models with respect to an increase (decrease) of one unit of surface charge on the  $\text{ASR}_{AT}$ , Rh and  $\text{Rh}_{181}$  models

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

(see Fig. 5A for Rh). From inspection of the results displayed in Fig. 5B, it is apparent that the models have opposite sensitivities to electrostatic fields approximately aligned along the protein axis. Rh and Rh<sub>181</sub> display a  $\Delta E_{S_1-S_0}$  increase when increasing the charge on the outer (extracellular) surface and decreasing the inner (intracellular) surface. Instead, ASR<sub>AT</sub> displays the opposite trend. In both cases the  $\Delta E_{S_1-S_0}$  changes are limited to ca. 2 kcal/mol for unit charge.

The opposite sensitivities displayed by ASR<sub>AT</sub> and Rh (or Rh<sub>181</sub>) call for a molecular-level explanation. As previously reported,<sup>88</sup> the effects of the surface charges on the  $\Delta E_{S_1-S_0}$  value can be understood using the mechanistic model illustrated in Scheme 1B. An electrostatic field stabilizing the positive charge in its S<sub>0</sub> location near the protonated Schiff base will increase the  $\Delta E_{S_1-S_0}$  value, while a destabilizing field would decrease  $\Delta E_{S_1-S_0}$ . In contrast, an electrostatic field stabilizing the positive charge in its S<sub>1</sub> location near the  $\beta$ -ionone ring will decrease the  $\Delta E_{S_1-S_0}$  value, while a destabilizing field would increase  $\Delta E_{S_1-S_0}$ . As shown in Fig. 5C left, the chromophore in Rh and Rh<sub>181</sub> is placed far from the two protein sides and oriented in such a way to expose its Schiff base to the negative surface and its  $\beta$ -ionone region to the more positive surface. The same arrangement is found for SqR, hMeOp and the bathoRh intermediate even if bathoRh features a distorted all-*trans* chromophore. In these situations an extra negative charge placed on the outer surface together with a positive charge placed on the inner surface would stabilize the S<sub>0</sub> and destabilize the S<sub>1</sub> state, thus resulting into a  $\Delta E_{S_1-S_0}$  value increase. Therefore the computed ca. 2 kcal/mol increase of the  $\Delta E_{S_1-S_0}$  values is consistent with the creation of an electrostatic field by the counterion placement scheme (see open circles in Fig. 3B) with respect to the reference NSC placement.



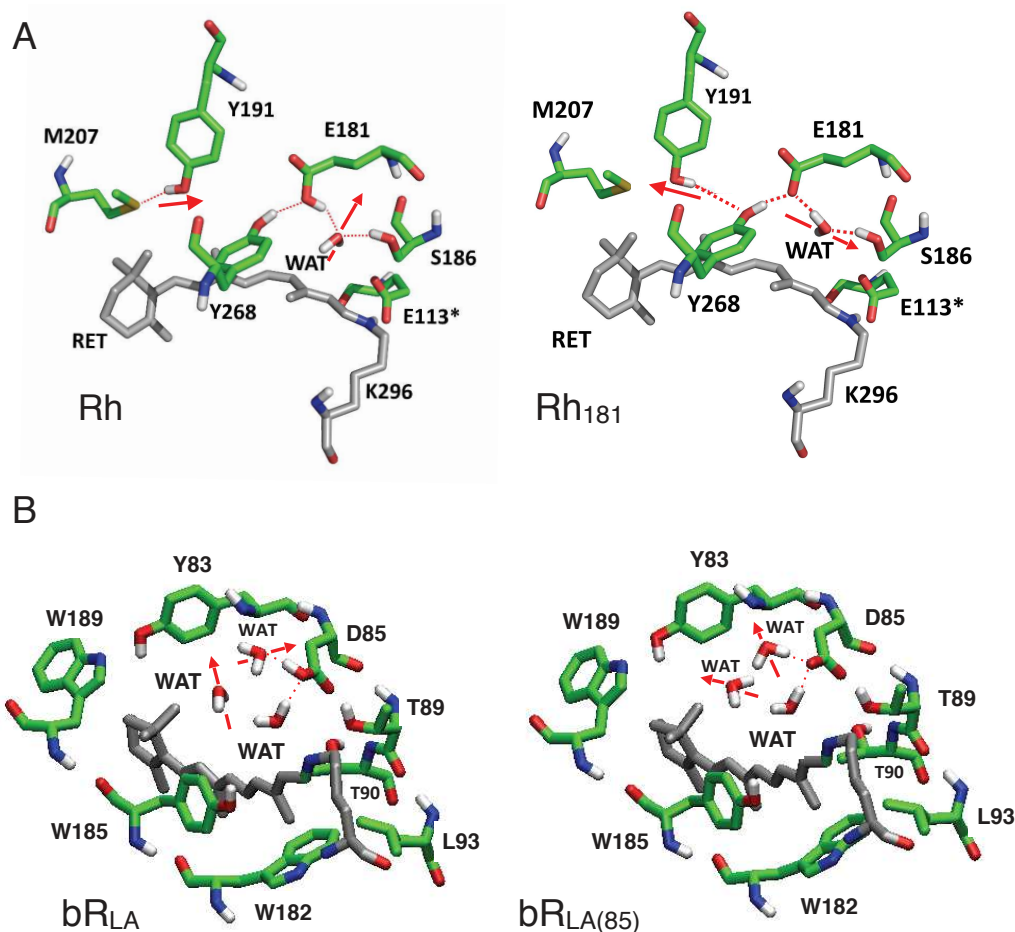
**Figure 5** Features affecting the vertical excitation energy in the Rh model. (A) Different external counterion distribution. Red spheres and blue spheres correspond to  $\text{Cl}^-$  and  $\text{Na}^+$  ions respectively. The Rh models correspond to the standard ARM model using the NSC scheme (center) and to two alternative model where the surface total charge is increased or decreased of a one unit. (B) Variations in average  $\Delta E_{S_1-S_0}$  values as a function of excess surface charges. (C) Orientation of the electrostatic field imposed by the reference external counterion distribution in Rh (left) and ASR<sub>AT</sub> (right). Positively charged side-chains are given in blue while negatively charged side-chains are given in red. The charge symbols at the top and bottom of the schematic 11-*cis* and all-*trans* chromophores indicate the direction of the protein field in the case of the reference counterion distribution and not the position of specific point charges.

The same reasoning explains the opposite effect of ASR<sub>AT</sub>. As shown in Fig. 5C in ASR<sub>AT</sub> the chromophore is placed closer to the negative protein surface and oriented in such a way to expose its  $\beta$ -ionone region to the negative charge. The consequence would be a decrease of the  $\Delta E_{S_1-S_0}$  value (stabilization of the  $S_1$  state with respect to the  $S_0$  state), which explains why the value of  $\Delta E_{S_1-S_0}$  for both ASR forms decreases upon creation of an electrostatic field by placing an extra negative counterion in the outer surface and a positive counterion in the inner surface. These qualitative explanations of the behaviour reported in Fig. 5B on the basis of the charge translocation effect (see Scheme 1) and chromophore orientation (see Fig. 5C) are supported by electrostatic potential calculations whose results are given in Section 10 of the SI.

#### Alternative assignments of ionization states: Rh and bR

As reported in Table 1 and Fig. 3A the effect of the ionization on the average  $\Delta E_{S_1-S_0}$  value for rhodopsin is large. Such value increases of ca. 6 kcal/mol leading to a strong  $\lambda_{\text{max}}^a$  blue-shift when

going from the Rh model featuring a protonated E181 residue, to the deprotonated Rh<sub>181</sub> model featuring a negatively charge E181 residue. The Rh<sub>181</sub> model yields a blue-shifted  $\Delta E_{S_1-S_0}$  value with respect to the experimental data and, in contrast to the Rh model, appears in line with the type of positive deviation observed for the other ARM models.



**Figure 6** Features affecting the vertical excitation energy in vertebrate and archaea rhodopsin ARM models. (A) Ionization state effects of E181. The two structures illustrate the most frequent HBN and water orientation in the Rh (E181 protonated) and Rh<sub>181</sub> (E181 deprotonated) models. (B) Ionization state effects of D85. The two structures illustrate the most frequent HBN and water orientation in bR<sub>LA</sub> (D85 protonated) and bR<sub>LA</sub>(85) (D85 deprotonated) models. The red arrows indicate the changes in dipole direction.

As revealed by a comparative analysis of the chromophore cavities in Rh<sub>181</sub> and Rh models, the computed  $\Delta E_{S_1-S_0}$  change represents the response to the large variation of electrostatics and induced HBN changes in the chromophore surroundings. More specifically, as shown in Fig. 6A, Rh<sub>181</sub> features a HBN where a water molecule (WAT) has, with respect to Rh (see Fig. 6A), an O-H group flipped and, rather obviously, forms a strong hydrogen bond with one of the E181 carboxylate oxygen. As a consequence the new WAT orientation modifies the local dipole moment direction leading to a further stabilization of the chromophore positive charge in its S<sub>0</sub> location (i.e. on the C=N bond). Another change regards the flipping of the OH group of the Y191 side-chain, which further stabilize the Y268 conformation and extends the local HBN. In conclusion, while the direct stabilization of S<sub>0</sub>

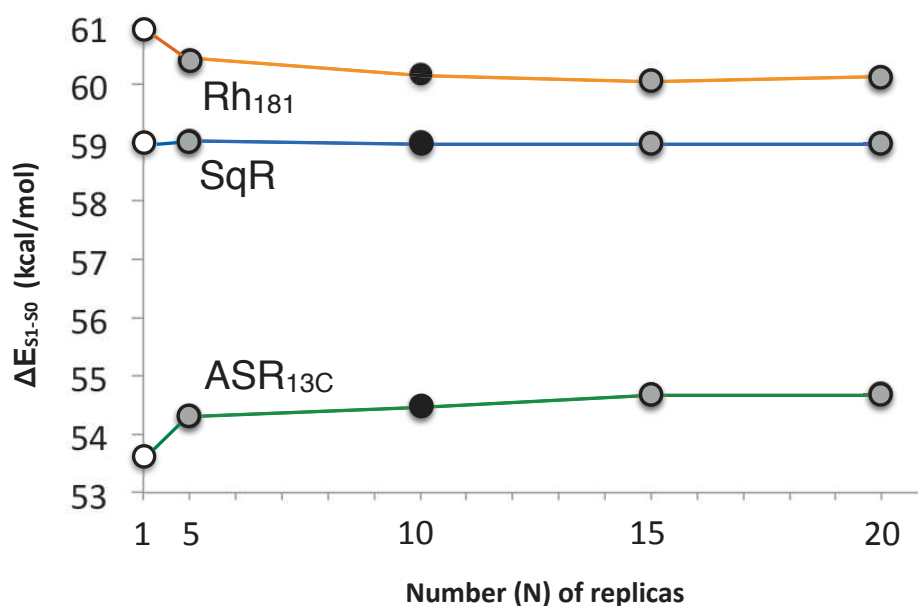
1  
2  
3 due to the E181 negative charge may not be large in Rh (i.e. such charge is located half-way along  
4 the chromophore backbone), the Rh<sub>181</sub> model indicates that the same negative charge can stabilize  
5 S<sub>0</sub> indirectly through a modification of the local HBN (compare left and right structures in Fig. 6A).  
6 This is confirmed by the reported HBN resulting from extensive MD simulations of bovine rhodopsin  
7 model with a ionized E181 side-chain.<sup>28</sup>  
8

9  
10 An even larger effect is observed for the bR<sub>LA(85)</sub> and bR<sub>DA(85)</sub> models (10 kcal/mol increase in  $\Delta E_{S_1-S_0}$ ),  
11 as the negatively charged D85 residue would further stabilize the Schiff base, together with the  
12 negatively charged D212 residue. The additional charge causes a further stabilization of the S<sub>0</sub> state  
13 relative to the S<sub>1</sub>, predicting a further blue-shifted value relatively to the experiment. The  
14 corresponding HBN modification is reported in Fig. 6B where we compare part of the chromophore  
15 cavity of light-adapted bacteriorhodopsin with protonated and deprotonated D85 residue (bR<sub>LA</sub> and  
16 bR<sub>LA(85)</sub> respectively). In contrast, to the bovine rhodopsin case, these results indicate that the bR<sub>LA</sub>,  
17 and not bR<sub>LA(85)</sub>, is the model best accommodating the trend of  $\Delta E_{S_1-S_0}$  values. A similar results is  
18 obtained for bR<sub>DA</sub>.  
19  
20  
21  
22  
23  
24  
25  
26  
27

### 28 29 *Stability of the computed side-chain conformations and hydrogen bond network*

30  
31 The ten 1 ns MD relaxations incorporated, for each rhodopsin model, in the ARM workflow (see  
32 Methods section) allow the sampling of different cavity conformations and HBN variations. The  
33 sensitivity of the resulting average  $\Delta E_{S_1-S_0}$  values as a function of the number of replicas (N=1, N=5,  
34 N=10, N=15 and N=20) has been directly assessed for the very different Rh<sub>181</sub>, SqR and ASR<sub>13C</sub>  
35 models. The results reported in Fig. 7 show that the average values become substantially stable, with  
36 variations less than 1 kcal/mol beyond ten replicas. N=10 is therefore assumed to be a suitable  
37 default value.  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60





**Figure 7** Features affecting the average vertical excitation energy in vertebrate, invertebrate and bacterial rhodopsin models. Effect of the sampling via  $N$  uncorrelated replicas (using different random seeds) 1 ns MD. The open circles and full circles correspond to the values of Fig. 3A and Table 1.

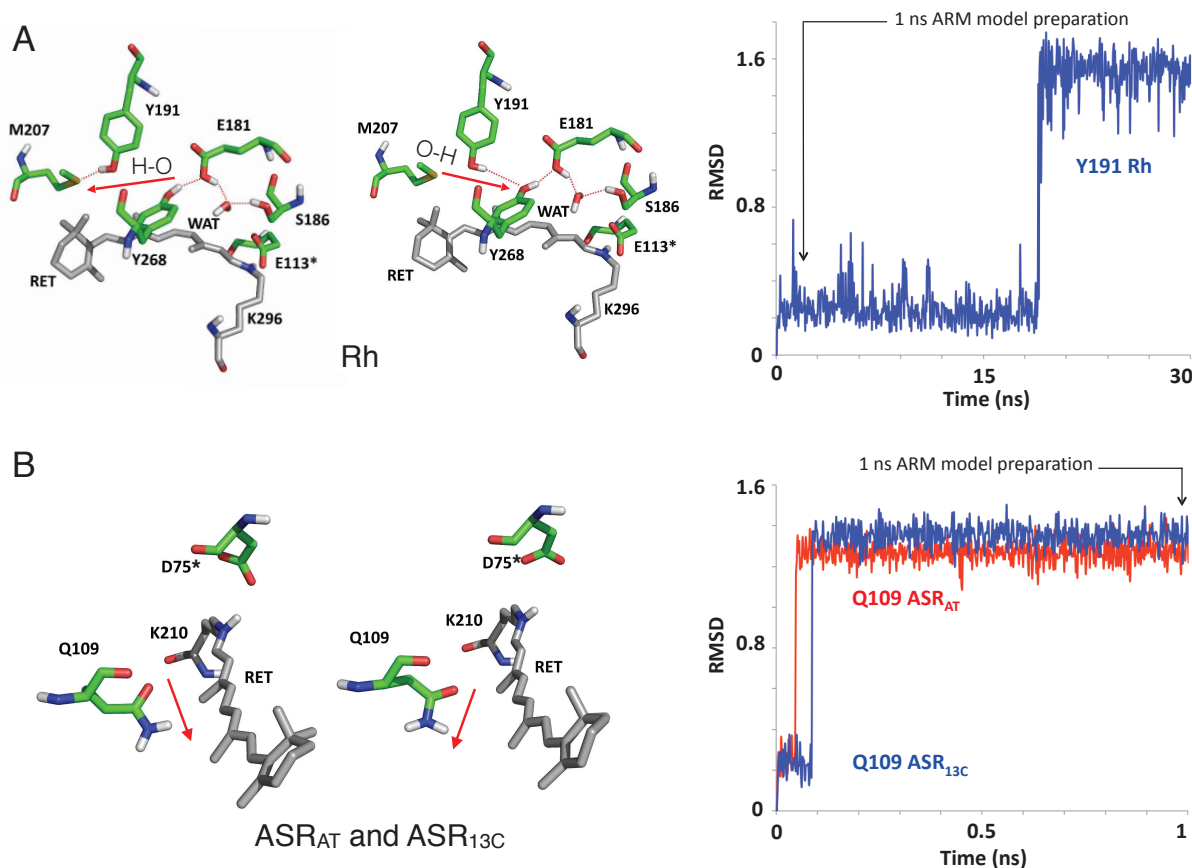
In Fig. 8A we instead assess the ARM model stability along a single, room-temperature 30 ns MD of the Rh model which, due to its protonated E181 residue and, therefore, weaker HBN, is more likely to undergo conformational changes with respect to Rh<sub>181</sub>. It can be seen that during the first 20 ns the HBN is stable and consistent with that of the Rh model (compare Fig. 8A left structure with the Rh structure of Fig. 6A). On the other hand, at times longer than ca. 20 ns the HBN changes and displays a flipped Y191 side-chain as shown in the right structure of Fig. 8A. An ARM model (Rh<sub>Y191</sub>), generated by a snapshot taken at the end of the 30 ns MD and constructed via geometry optimization produces a 1.3 kcal/mol lower  $\Delta E_{S1-S0}$  value than the Rh model (see Table 3). Such difference does not change the predicted trend for wild-type rhodopsin set, but can be used to assess the model conformational flexibility. In the Rh case, the variation is attributed to a re-orientation of the dipole moment associated with the Y191 terminal O-H bond (compare the two structures in Fig. 8A). In Rh<sub>Y191</sub> the O-H is hydrogen bound to the Y268 residue rather than to the M207 residue and therefore projects a positive rather than negative potential in the C=N region leading to a destabilization of the  $S_0$  state with respect to the  $S_1$  state. However, Rh<sub>Y191</sub> has a potential energy ca. 9 kcal/mol higher than the Rh model (see the SI). Both the higher stability and the better  $\lambda_{max}^a$  prediction would favor the initial conformation, in spite of the fact that the observed HBN does not revert back during the last ten ns of the 30 ns MD. In fact, this is still a too short simulation to conclude that this is a stable HBN. As discussed above this orientation becomes stabilized when the E181 side-chain is deprotonated such as in Rh<sub>181</sub>. This is because the negative E181 charge strongly stabilizes the HBN formed by Y191, Y268, E181, WAT, S186. A similar 30 ns MD

test carried out on SqR never displays an HBN change, further pointing to system-dependent results of these tests. A comprehensive investigation of the HBN dynamics would require longer and computationally demanding MD relaxations and are beyond the scope of the present work.

**Table 3.** Computed vertical excitation energies  $\Delta E_{S_1-S_0}$  in kcal/mol and maximum absorption wavelengths  $\lambda_{max}^a$  (nm) for alternative ARM models of Rh and ASR. The deviation from the experimental value and the computed oscillator strengths  $f_{osc}$  are also given.

	Calc. $\Delta E$ ( $\lambda_{max}^a$ ) Average (N=10)	Error Average (N=10)	$f_{osc}$
<b>Rh<sub>Y191</sub></b>	52.5 (544)	-4.9	0.90
<b>ASR<sub>AT/Q109</sub></b>	54.2 (527)	2.1	1.08
<b>ASR<sub>13C/Q109</sub></b>	55.0 (519)	1.8	1.13

An investigation of the side chain stability and conformational effects has also been carried out for the two ASR models. Inspection of the conformations of side chains involved in the HBN of both ASR<sub>AT</sub> and ASR<sub>13C</sub> reveals that these are similar in the two models, so the HBN effect on the computed  $\Delta E_{S_1-S_0}$  value is thought to be approximately the same. However, as shown in Fig. 8B, such conformations are different from those initially assigned to the crystallographic structure templates (left structure in Fig. 8B). In fact a conformational change occurs already during the initial 150 ps MD equilibration, which modifies the orientation of the Q109 side-chain (yielding the structure of Fig. 8B right). To estimate the effect of such change we have constructed two models (ASR<sub>AT/Q109</sub> and ASR<sub>13C/Q109</sub>) starting from average MD structures where the Q109 side-chain orientation was modified to resemble the crystal structure conformation. As reported in Table 3 a slightly larger  $\Delta E_{S_1-S_0}$  value is obtained for these models (e.g. for ASR<sub>AT/Q109</sub> has an excitation energy 0.7 kcal/mol higher than ASR<sub>AT</sub>), consistently with the change in orientation of the dipole associated to the amino acid side chain (see Fig. 8B), which stabilizes  $S_0$  with respect to  $S_1$  and increases the energy gap. However the ASR<sub>AT/Q109</sub>  $S_0$  energy is ca. 10 kcal/mol higher than the default ARM model. So, effectively the ARM protocol seems to have “improved” on the ASR crystal structure by generating a more stable structure. Nevertheless, these results show again the importance of side chain orientation and the potentially large effect they can have on the computed vertical excitation energy.



**Figure 8** Changes in the side-chain conformation and related HBN. (A) Change induced in the reference ARM Rh model when a 30 ns room-temperature MD run is employed for the rhodopsin cavity sampling. The red arrow illustrates the change in orientation of the dipole moment of the Y191 residue characterizing the conformational change occurring after 20 ns MD as shown in the RMSD vs time plot. The ARM model generated starting from the 30 ns snapshot is called Rh<sub>Y191</sub> and has a different HBN with respect to Rh. (B) Change in the HBN occurring during the protocol 1 ns MD run of ASR. The left structure corresponds to the original assigned X-ray crystal structure and it is obtained by geometry optimization of the initial snapshot. The right structure is the one of the final ARM model. RMSD vs time plot for the MD runs showing the corresponding conformational changes of ASR<sub>AT</sub> and ASR<sub>13C</sub> Q109 side-chain.

## Conclusions and Outlook

Accurate QM/MM models of rhodopsins require a realistic treatment of the protein environment. This includes a membrane patch hosting the protein; inner and outer membrane water molecule layers with chloride and sodium ions (i.e. the external counterions) to counterbalance the solvent exposed ionized amino acids and obtain a globally neutral model; and the right number of protein monomers according to experimental data. A set of models of this type, targeting accurate spectroscopic studies,<sup>28,30-32</sup> were indeed reported including an hydrated bilayer preparation, relaxed and equilibrated with several nanoseconds of molecular dynamics (MD), followed by multiple QM/MM treatments for further sampling the protein interior and chromophore before computing the  $\lambda_{\max}^a$  values. Even more resource-intensive methodologies were applied to artificial soluble retinal proteins.<sup>30</sup> Such computationally demanding methodologies are, presently, not adequate for

1  
2  
3 a fast *in silico* screening because they are time-consuming and their complex model building makes  
4 the workflow control difficult to automate. In order to overcome such difficulties, we have  
5 developed and benchmarked ARM as a prototype protocol for high throughput computational  
6 screening of photoreceptors and a reference for further model improvement. Accordingly, ARM  
7 generates approximated rhodopsin models useful for quickly screening experimental or computer  
8 generated sequences with the target of tracking trends in spectral or reactivity properties. Such  
9 analysis would provide candidates worth of more demanding and accurate modeling or direct  
10 experimental assessment. Thus, possible ARM applications would be the screening of the sequences  
11 associated to an evolutionary tree for tracking a specific ancestor or generated via *in silico*  
12 mutagenesis looking for sequences where a certain property is optimized.  
13  
14  
15  
16  
17  
18  
19  
20

21 We have shown that the information required for the construction of each ARM model is essentially  
22 a crystallographic structure or a comparative model provided in a PDB file format. However, in the  
23 current version of ARM such primary input needs to be completed by assigning the state of all  
24 ionizable side-chains and the list of the target amino acid required for counterion placement. In  
25 contrast, features such as the HBN structure (including the internal water molecule position), the  
26 Na<sup>+</sup> and Cl<sup>-</sup> counterion position and the side-chain conformation of the amino acid forming the  
27 chromophore cavity are generated automatically. Thus, presently, ARM models would be more  
28 correctly defined as semi-automatically generated models.  
29  
30  
31  
32  
33  
34  
35

36 The ARM models (see the SI for the corresponding PDB files) generated for 6 different sensory  
37 rhodopsins, 3 proton-pumping rhodopsins and one ion channel show that it is possible to reproduce  
38 the observed general trends in vertical excitation energies with blue-shifted values not significantly  
39 displaced from the observed data. This is encouraging when considering that the benchmarked  
40 rhodopsins come from evolutionarily distant organisms with different physiological functions.  
41 Remarkably, a similar level of accuracy is preserved when predicting the more challenging trend of  
42 17 bacterial and vertebrate rhodopsin mutants. While we have shown that different factors affect  
43 the computed  $\Delta E_{S1-S0}$  value, the changes due to different external counterion distributions and side-  
44 chain conformations are usually small (e.g. ca. 3-2 kcal mol<sup>-1</sup>). Larger changes have been shown to be  
45 associated with different ionization state assignments of the internal residues that would also affect  
46 the stability of the HBN. For cases when these situations occur during the preliminary MD run, our  
47 results suggest that distinct models of the same rhodopsin with different ionization state  
48 assignments should be included in the set in order to examine the possible variations in predicted  
49 trends.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Future work will address the present uncertainties, by introducing a more accurate placement and optimization of all hydrogen atoms, a better treatment of the external environment, an improved treatment of the chromophore cavity. Furthermore, in the case of hMeOp, we used a comparative model as starting structure (there is no X-ray crystal structure for this protein) derived according to a procedure discussed in ref. <sup>9</sup> and taking as a template SqR (see the Supporting Information). The sensitivity of the predicted  $\Delta E_{S1-S0}$  values to the specific procedure used to generate the starting hMeOp model and its impact on the error relative to SqR remain to be investigated.

In conclusion, ARM constitutes the first step towards a fast and cheap computational tool filling the gap between slow experimental structure determination and fast protein expression and characterization achieved by genetic engineering techniques. It also represents an attempt to provide a basis for the standardization and reproducibility of rhodopsin QM/MM models. If such research line will be proven to be fruitful, our ultimate goal would be the extension of the same methodologies to other photoresponsive protein families, so that ARM could eventually become a general tool in computational photobiology.

### **ASSOCIATED CONTENT**

Supporting text and coordinates of the input wild-type structures (11 models) and optimized (28 models) benchmark wild-type and mutant models in PDB format. This material is available free of charge via the Internet at <http://pubs.acs.org>.

### **ACKNOWLEDGEMENT**

This work was supported by the National Science Foundation under grant no. CHE-1152070 and CHE-1551416 and the Human Frontier Science Program Organization under grant RGP0049/2012CHE09-56776. M.O. is grateful to the Center for Photochemical Sciences and School of Arts and Sciences of the Bowling Green State University. The authors are indebted to NSF-XSEDE and Ohio Supercomputer Center for granted computer time. M.O. and Y.O. are grateful to the Foundation of the University of Strasbourg for an USIAS fellowship supporting a visiting professorship at the IPCMS. M.S. and N.F. thank the French Agence Nationale de la Recherche for funded project FEMTO-ASR, n° ANR-14-CE35-0015-02.

### **References**

- (1) Kircher, M.; Kelso, J. *BioEssays* **2010**, *32*, 524-536.
- (2) Doyle, S.; Koehn, J.; Hunt, I. in *High Throughput Protein Expression and Purification*; Humana Press: Totowa, NJ, 2009; 498, pp 1-18.
- (3) Berry, S. M.; Lu, Y. Protein Structure Design and Engineering in *eLS*; John Wiley & Sons, Ltd: published on-line, 2001.
- (4) Warshel, A.; Levitt, M. *J. Mol. Biol.* **1976**, *103*, 227-249.

- 1  
2  
3 (5) Altoè, P.; Cembran, A.; Olivucci, M.; Garavelli, M. *Proc. Natl. Acad. Sci. U. S. A.* **2010**, *107*, 20172-  
4 20177.  
5  
6  
7 (6) Andruniow, T.; Ferré, N.; Olivucci, M. *Proc. Natl. Acad. Sci. U. S. A.* **2004**, *101*, 17908-17913.  
8  
9  
10 (7) Coto, P. B.; Strambi, A.; Ferré, N.; Olivucci, M. *Proc. Natl. Acad. Sci. U. S. A.* **2006**, *103*, 17154-  
11 17159.  
12  
13 (8) Strambi, A.; Durbeij, B.; Ferré, N.; Olivucci, M. *Proc. Natl. Acad. Sci. U. S. A.* **2010**, *107*, 21322-  
14 21326.  
15  
16 (9) Rinaldi, S.; Melaccio, F.; Gozem, S.; Fanelli, F.; Olivucci, M. *Proc. Natl. Acad. Sci. U. S. A.* **2014**, *111*,  
17 1714-1719.  
18  
19 (10) Huntress, M. M.; Gozem, S.; Malley, K. R.; Jailaubekov, A. E.; Vasileiou, C.; Vengris, M.; Geiger, J.  
20 H.; Borhan, B.; Schapiro, I.; Larsen, D. S.; Olivucci, M. *J. Phys. Chem. B.* **2013**, *117*, 10053-10070.  
21  
22 (11) Schreiber, M.; Sugihara, M.; Okada, T.; Buss, V. *Angew. Chem. Int. Ed.* **2006**, *45*, 4274-4277.  
23  
24 (12) Tomasello, G.; Olaso-González, G.; Altoè, P.; Stenta, M.; Serrano-Andrés, L.; Merchà, M.;  
25 Orlandi, G.; Bottoni, A.; Garavelli, M. *J. Am. Chem. Soc.* **2009**, *131*, 5172-5186.  
26  
27 (13) Altun, A.; Yokoyama, S.; Morokuma, K. *J. Phys. Chem. B.* **2008**, *112*, 16883-16890.  
28  
29 (14) Altun, A.; Yokoyama, S.; Morokuma, K. *J. Phys. Chem. B.* **2008**, *112*, 6814-6827.  
30  
31 (15) Ernst, O. P.; Lodowski, D. T.; Elstner, M.; Hegemann, P.; Brown, L. S.; Kandori, H. *Chem. Rev.*  
32 **2014**, *114*, 126-163.  
33  
34 (16) Spudich, J. L.; Yang, C. S.; Jung, K. H.; Spudich, E. N. *Annu. Rev. Cell. Dev. Bi.* **2000**, *16*, 365-392.  
35  
36 (17) Gorostiza, P.; Isacoff, E. Y. *Science* **2008**, *322*, 395-399.  
37  
38 (18) Boyden, E. S.; Zhang, F.; Bamberg, E.; Nagel, G.; Deisseroth, K. *Nat. Neurosci.* **2005**, *8*, 1263-  
39 1268.  
40  
41 (19) Han, X.; Boyden, E. S. *PLoS One* **2007**, *2*, e299.  
42  
43 (20) Schoenenberger, P.; Gerosa, D.; Oertner, T. G. *PLoS One* **2009**, *4*, e8185.  
44  
45 (21) Deisseroth, K. *Nat. Methods.* **2011**, *8*, 26-29.  
46  
47 (22) Sekharan, S.; Sugihara, M.; Buss, V. *Angew. Chem. Int. Ed.* **2007**, *46*, 269-271.  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4 (23) Bravaya, K.; Bochenkova, A.; Granovsky, A.; Nemukhin, A. *J. Am. Chem. Soc.* **2007**, *129*, 13035-  
5 13042.  
6  
7  
8 (24) Fujimoto, K.; Hayashi, S.; Hasegawa, J.; Nakatsuji, H. *J. Chem. Theory Comput.* **2007**, *3*, 605-618.  
9  
10 (25) Schapiro, I.; Ryazantsev, M. N.; Ding, W. J.; Huntress, M. M.; Melaccio, F.; Andruniow, T.;  
11 Olivucci, M. *Aust. J. Chem.* **2010**, *63*, 413-429.  
12  
13  
14 (26) El-Khoury, P. Z.; Schapiro, I.; Huntress, M. M.; Melaccio, F.; Gozem, S.; Frutos, L. M.; Olivucci, M.  
15 in *CRC Handbook of Organic Photochemistry and Photobiology*; Griesbeck, A., Oelgemöller, M., and  
16 Ghetti, F., Eds.; CRC Press: Boca Raton, FL, 2012; 1029-1056.  
17  
18  
19  
20 (27) Singh, U. C.; Kollman, P. A. *J. Comput. Chem.* **1986**, *7*, 718-730.  
21  
22  
23 (28) Valsson, O.; Campomanes, P.; Tavernelli, I.; Rothlisberger, U.; Filippi, C. *J. Chem. Theory Comput.*  
24 **2013**, *9*, 2441-2454.  
25  
26  
27 (29) Campomanes, P.; Neri, M.; Horta, B. A. C.; Röhrig, U. F.; Vanni, S.; Tavernelli, I.; Rothlisberger, U.  
28 *J. Am. Chem. Soc.* **2014**, *136*, 3842-3851.  
29  
30  
31 (30) Cheng, C.; Kamiya, M.; Uchida, Y.; Hayashi, S. *J. Am. Chem. Soc.* **2015**, *137*, 13362-13370.  
32  
33  
34 (31) Kato, H. E.; Kamiya, M.; Sugo, S.; Ito, J.; Taniguchi, R.; Orito, A.; Hirata, K.; Inutsuka, A.;  
35 Yamanaka, A.; Maturana, A. D.; Ishitani, R.; Sudo, Y.; Hayashi, S.; Nureki, O. *Nat. Commun.* **2015**, *6*,  
36 7177.  
37  
38  
39 (32) Guo, Y.; Beyle, F. E.; Bold, B. M.; Watanabe, H. C.; Koslowski, A.; Thiel, W.; Hegemann, P.;  
40 Marazzi, M.; Elstner, M. *Chem. Sci.* **2016**, *7*, 3879-3891.  
41  
42  
43 (33) Jung, K. H.; Trivedi, V. D.; Spudich, J. L. *Mol. Microbiol.* **2003**, *47*, 1513-1522.  
44  
45  
46 (34) Vogeley, L.; Sineshchekov, O. A.; Trivedi, V. D.; Sasaki, J.; Spudich, J. L.; Luecke, H. *Science* **2004**,  
47 *306*, 1390-1393.  
48  
49  
50 (35) Okada, T.; Sugihara, M.; Bondar, A. -N.; Elstner, M.; Entel, P.; Buss, V. *J. Mol. Biol.* **2004**, *342*,  
51 571-583.  
52  
53  
54 (36) Murakami, M.; Kouyama, T. *Nature* **2008**, *453*, 363-367.  
55  
56  
57 (37) Provencio, I.; Jiang, G.; De Grip, W. J.; Hayes, W. P.; Rollag, M. D. *Proc. Natl. Acad. Sci. U. S. A.*  
58 **1998**, *95*, 340-345.  
59  
60

- 1  
2  
3 (38) Luecke, H.; Schobert, B.; Richter, H. -T.; Cartailier, J. -P.; Lanyi, J. K. *J. Mol. Biol.* **1999**, *291*, 899-  
4 911.  
5  
6  
7  
8 (39) Nishikawa, T.; Murakami, M.; Kouyama, T. *J. Mol. Biol.* **2005**, *352*, 319-328.  
9  
10 (40) Ran, Tingting; Ozorowski, Gabriel; Gao, Yanyan; Sineshchekov, Oleg A.; Wang, Weiwu; Spudich,  
11 John L.; Luecke, Hartmut *Acta Crystallogr. D* **2013**, *69*, 1965-1980.  
12  
13 (41) Yoshizawa, T.; Wald, G. *Nature* **1963**, *197*, 1279-1286.  
14  
15  
16 (42) Nakamichi, H.; Okada, T. *Angew. Chem. Int. Ed.* **2006**, *45*, 4270-4273.  
17  
18  
19 (43) Mathies, R.; Stryer, L. *Proc. Natl. Acad. Sci. U. S. A.* **1976**, *73*, 2169-2173.  
20  
21  
22 (44) Bonacic-Koutecky, V.; Köhler, K.; Michl, J. *Chem. Phys. Lett.* **1984**, *104*, 440-443.  
23  
24  
25 (45) Roos, B. O. in *Ab. initio. Methods. in. Quantum. Chemistry.*; Lawley, K. P., Ed. Wiley: New York,  
26 1987; pp 399-446.  
27  
28  
29 (46) Andersson, K.; Malmqvist, P. -; Roos, B. O.; Sadlej, A. J.; Wolinski, K. J. *J. Phys. Chem.* **1990**, *94*,  
30 5483-5488.  
31  
32  
33 (47) Coto, P. B.; Martí, S.; Oliva, M.; Olivucci, M.; Merchán, M.; Andrés, J. *J. Phys. Chem. B* **2008**, *112*,  
34 7153-7156.  
35  
36  
37 (48) Navizet, I.; Liu, Y. -J.; Ferré, N.; Xiao, H. -Y.; Fang, W. -H.; Lindh, R. *J. Am. Chem. Soc.* **2009**, *132*,  
38 706-712.  
39  
40  
41 (49) Altoè, P.; Climent, T.; De Fusco, G. C.; Stenta, M.; Bottoni, A.; Serrano-Andrès, L.; Merchàn, M.;  
42 Orlandi, G.; Garavelli, M. *J. Phys. Chem. B* **2009**, *113*, 15067-15073.  
43  
44  
45 (50) Ferré, N.; Olivucci, M. *J. Am. Chem. Soc.* **2003**, *125*, 6868-6869.  
46  
47  
48 (51) Schapiro, I.; Ryazantsev, M. N.; Frutos, L. M.; Ferré, N.; Lindh, R.; Olivucci, M. *J. Am. Chem. Soc.*  
49 **2011**, *133*, 3354-3364.  
50  
51  
52 (52) Polli, D.; Altoe, P.; Weingart, O.; Spillane, K. M.; Manzoni, C.; Brida, D.; Tomasello, G.; Orlandi,  
53 G.; Kukura, P.; Mathies, R. A.; Garavelli, M.; Cerullo, G. *Nature* **2010**, *467*, 440-443.  
54  
55  
56 (53) Frutos, L. M.; Andruniow, T.; Santoro, F.; Ferré, N.; Olivucci, M. *Proc. Natl. Acad. Sci. U. S. A.*  
57 **2007**, *104*, 7764-7769.  
58  
59  
60



- 1  
2  
3  
4 (54) Gozem, S.; Melaccio, F.; Lindh, R.; Krylov, A. I.; Granovsky, A. A.; Angeli, C.; Olivucci, M. *J. Chem.*  
5 *Theory Comput.* **2013**, *9*, 4495-4506.  
6  
7  
8 (55) Sekharan, S.; Altun, A.; Morokuma, K. *Chem-Eur. J.* **2010**, *16*, 1744-1749.  
9  
10 (56) Sekharan, S.; Wei, J. N.; Batista, V. S. *J. Am. Chem. Soc.* **2012**, *134*, 19536-19539.  
11  
12 (57) Humphrey, W.; Dalke, A.; Schulten, K. *J. Mol. Graphics* **1996**, *14*, 33-38.  
13  
14 (58) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Gould, I. R.; Merz, K. M.; Ferguson, D. M.; Spellmeyer, D.  
15 C.; Fox, T.; Caldwell, J. W.; Kollman, P. A. *J. Am. Chem. Soc.* **1995**, *117*, 5179-5197.  
16  
17 (59) Zhang, L.; Hermans, J. *Proteins* **1996**, *24*, 433-438.  
18  
19 (60) Ponder, J. W.; Richards, F. M. *J. Comput. Chem.* **1987**, *8*, 1016-1024.  
20  
21 (61) Strambi, A.; Coto, P. B.; Frutos, L. M.; Ferré, N.; Olivucci, M. *J. Am. Chem. Soc.* **2007**, *130*, 3382-  
22 3388.  
23  
24 (62) Jardón-Valadez, E.; Bondar, A. -N.; Tobias, D. J. *Biophys. J.* **2009**, *96*, 2572-2576.  
25  
26 (63) Jardón-Valadez, E.; Bondar, A. -N.; Tobias, D. J. *Biophys. J.* **2010**, *99*, 2200-2207.  
27  
28 (64) Dundas, J.; Ouyang, Z.; Tseng, J.; Binkowski, A.; Turpaz, Y.; Liang, J. *Nucleic. Acids. Res.* **2006**, *34*,  
29 W116-W118.  
30  
31 (65) Pronk, S.; Páll, S.; Schulz, R.; Larsson, P.; Bjelkmar, P.; Apostolov, R.; Shirts, M. R.; Smith, J. C.;  
32 Kasson, P. M.; van der Spoel, D.; Hess, B.; Lindahl, E. *Bioinformatics* **2013**, *29*, 845-854.  
33  
34 (66) Ferré, N.; Cembran, A.; Garavelli, M.; Olivucci, M. *Theor. Chem. Acc.* **2004**, *112*, 335-341.  
35  
36 (67) Aquilante, F.; De Vico, L.; Ferré, N.; Ghigo, G.; Malmqvist, P. A.; Neogady, P.; Pedersen, T. B.;  
37 Pitonak, M.; Reiher, M.; Roos, B. O.; Serrano-Andres, L.; Urban, M.; Veryazov, V.; Lindh, R. *J. Comput.*  
38 *Chem.* **2010**, *31*, 224-247.  
39  
40 (68) Singh, U. C.; Kollman, P. A. *J. Comput. Chem.* **1986**, *7*, 718-730.  
41  
42 (69) Humbel, S.; Sieber, S.; Morokuma, K. *J. Chem. Phys.* **1996**, *105*, 1959-1967.  
43  
44 (70) Stålring, J.; Bernhardsson, A.; Lindh, R. *Mol. Phys.* **2001**, *99*, 103-114.  
45  
46 (71) Ferré, N.; Angyan, J. G. *Chem. Phys. Lett.* **2002**, *356*, 331-339.  
47  
48 (72) Krivov, G. G.; Shapovalov, M. V.; Dunbrack, R. L. *Proteins* **2009**, *77*, 778-795.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 (73) Melaccio, F.; Olivucci, M.; Lindh, R.; Ferré, N. *Int. J. Quant. Chem.* **2011**, *111*, 3339-3346.  
4  
5  
6 (74) Kato, H. E.; Zhang, F.; Yizhar, O.; Ramakrishnan, C.; Nishizawa, T.; Hirata, K.; Ito, J.; Aita, Y.;  
7 Tsukazaki, T.; Hayashi, S. *Nature* **2012**, *482*, 369-374.  
8  
9  
10 (75) Olsson, M. H. M.; Sondergaard, C. R.; Rostkowski, M.; Jensen, J. H. *J. Chem. Theory Comput.*  
11 **2011**, *7*, 525-537.  
12  
13  
14 (76) Främcke, J. S.; Wanko, M.; Phatak, P.; Mroginski, M. A.; Elstner, M. *J. Phys. Chem. B.* **2010**, *114*,  
15 11338-11352.  
16  
17  
18 (77) Sandberg, M. N.; Amora, T. L.; Ramos, L. S.; Chen, M. -H.; Knox, B. E.; Birge, R. R. *J. Am. Chem.*  
19 *Soc.* **2011**, *133*, 2808-2811.  
20  
21  
22  
23 (78) Hall, K. F.; Vreven, T.; Frisch, M. J.; Bearpark, M. J. *J. Mol. Biol.* **2008**, *383*, 106-121.  
24  
25  
26 (79) Lanyi, J. K. *BBA-Bioenergetics* **2006**, *1757*, 1012-1018.  
27  
28 (80) Kandori, H.; Shimono, K.; Sudo, Y.; Iwamoto, M.; Shichida, Y.; Kamo, N. *Biochemistry* **2001**, *40*,  
29 9238-9246.  
30  
31  
32 (81) Furutani, Y.; Kawanabe, A.; Jung, K. H.; Kandori, H. *Biochemistry* **2005**, *44*, 12287-12296.  
33  
34  
35 (82) Kawanabe, A.; Furutani, Y.; Jung, K. H.; Kandori, H. *Biochemistry* **2006**, *45*, 4362-4370.  
36  
37  
38 (83) Stein, C. J.; Reiher, M. *J. Chem. Theory Comput.* **2016**, *12*, 1760-1771.  
39  
40  
41 (84) Matsuyama, T.; Yamashita, T.; Imamoto, Y.; Shichida, Y. *Biochemistry* **2012**, *51*, 5454-5462.  
42  
43  
44 (85) Bailes, H. J.; Lucas, R. J. *P. Roy. Soc. B-Biol. Sci.* **2013**, *280*, 20122987.  
45  
46  
47 (86) Granovsky, A. A. *J. Chem. Phys.* **2011**, *134*, 214113-214127.  
48  
49 (87) Gozem, S.; Huntress, M.; Schapiro, I.; Lindh, R.; Granovsky, A. A.; Angeli, C.; Olivucci, M. *J. Chem.*  
50 *Theory Comput.* **2012**, *8*, 4069-4080.  
51  
52  
53 (88) Melaccio, F.; Ferré, N.; Olivucci, M. *Phys. Chem. Chem. Phys.* **2012**, *14*, 12485-12495.  
54  
55  
56  
57  
58  
59  
60

## TOC

