

# Fast Methods for Biomolecule Charge Optimization

J. P. Bardhan\*, J. H. Lee\*, S. S. Kuo\*, M. D. Altman<sup>†</sup>, B. Tidor<sup>‡,\*</sup>, J.K.White\*

\* Department of Electrical Engineering and Computer Science

<sup>†</sup> Department of Chemistry

<sup>‡</sup> Biological Engineering Division

77 Massachusetts Avenue, Cambridge, MA 02139

## ABSTRACT

Charge optimization is an essential element of rational drug design; given ligand and receptor proteins, one wishes to determine the ligand charge distribution—a vector of partial atomic charges—that maximizes the favorable change in electrostatic free energy on binding. Work in biophysics has shown that this problem is convex, and that it can be solved using standard quadratic programming methods. However, the use of these techniques requires an initial calculation of the Hessian matrix; the present work introduces a new method that avoids this expensive computation and also scales much more favorably with problem size. The technique couples a boundary element formulation with a primal-dual interior point algorithm; initial results suggest that the cost to solve these optimization problems can be reduced by more than an order of magnitude.

## 1 INTRODUCTION

One aspect of rational drug design is the consideration of electrostatic interactions between the receptor protein, whose activity is to be alleviated or blocked, and the designed ligand and protein [1]–[5]. The interactions are long range and have significant contributions to the overall free energy of binding; optimizing these interactions is therefore important, but unfortunately also quite challenging. Two opposing effects must be carefully balanced: favorable electrostatic interactions between the ligand and receptor in the bound state, and an unfavorable ligand desolvation penalty. Prior work [5] based on continuum electrostatics models has shown that the problem of optimizing the electrostatic free energy change in binding is convex and quadratic with respect to the ligand partial atomic charges.

Traditionally, the solution of the charge optimization problem has relied on optimization schemes that depend on an explicit form for the Hessian matrix. Determination of the Hessian is computationally expensive, however, and no optimization can be performed until the entire Hessian matrix has been calculated. This work demonstrates an optimization scheme that avoids the large preliminary computation, and in addition seems capable of accelerating the optimization process by more than an order of magnitude.

The new Hessian-implicit scheme is based on a primal-dual interior point method [6]. The use of a boundary element method [7], [8] to perform the electrostatic calculations

allows the boundary element simulation to be coupled into the interior-point framework; the two systems of equations are then solved simultaneously using iterative methods.

Section 2 briefly reviews the energy optimization problem, the primal-dual interior point framework, and the boundary element method used for energy calculations. Section 3 presents the Hessian-implicit method, and Section 4 illustrates the computational advantage of the new method. Section 5 concludes the paper and closes with ideas for future work.

## 2 BACKGROUND

### 2.1 Modeling Electrostatics

The problem of interest is to take a molecular surface, either that of the ligand alone, or that of the ligand-receptor complex, and compute the reaction potential in the ligand due to a charge distribution in the ligand. Boundary element methods are used to calculate this potential both for the ligand-receptor complex in solution (called the bound system) and for the ligand alone in solution (the unbound system). The following derivation summarizes the formulation presented in [7], [8], which is based on a mixed continuum-discrete model. We derive the relation between charge distribution and reaction potential for the bound system; the derivation is identical for the unbound system except that the ligand surface is taken as the interface between solvent and protein, rather than the surface of the ligand-receptor complex.

The solvent is treated as a homogeneous medium with high permittivity, which models the polarization effect of the solvent, and Debye screening length  $\kappa$ , which models the effect of mobile ions in solution [9], [10]. The protein interior is treated as a homogeneous region with low permittivity. The partial atomic charges are treated as discrete point charges at the atom centers. An electrostatic simulation takes as input a vector of point charge values, and produces as output the reaction potential at the atom centers due to the dielectric properties of the media and the mobile ions in solution.

Green’s second theorem is applied to the solvent region and the protein interior, producing two integral equations for the potential in the two regions. Applying the boundary conditions at the surface generates a coupled pair of integral

equations:

$$\begin{aligned} \frac{1}{2}\varphi(\vec{r}_\Omega) + \int_{\Omega_{bound}} \varphi(\vec{r}') \frac{\partial G_1}{\partial n}(\vec{r}_\Omega; \vec{r}') d\vec{r}' \\ - \int_{\Omega_{bound}} G_1(\vec{r}_\Omega; \vec{r}') \frac{\partial \varphi}{\partial n}(\vec{r}') d\vec{r}' \\ = \sum_{i=1}^{n_c} \frac{q_i}{\epsilon_1} G_1(\vec{r}_\Omega; \vec{r}_i), \end{aligned} \quad (1)$$

$$\begin{aligned} \frac{1}{2}\varphi(\vec{r}_\Omega) + \int_{\Omega_{bound}} -\varphi(\vec{r}') \frac{\partial G_2}{\partial n}(\vec{r}_\Omega; \vec{r}') d\vec{r}' \\ + \int_{\Omega_{bound}} G_2(\vec{r}_\Omega; \vec{r}') \frac{1}{\epsilon_r} \frac{\partial \varphi}{\partial n}(\vec{r}') d\vec{r}' = 0. \end{aligned} \quad (2)$$

Here,  $n_c$  is the number of point charges modeled in the ligand and  $\vec{r}_\Omega$  is any point on the molecular surface. Thus,  $\varphi(\vec{r}_\Omega)$  corresponds to the potential on the surface, and  $\frac{\partial \varphi}{\partial n}(\vec{r}_\Omega)$  corresponds to the normal derivative of the potential on the surface. Also,  $\epsilon_1$  and  $\epsilon_2$  are the low dielectric constant and high dielectric constants, respectively, and  $\epsilon_r = \epsilon_2/\epsilon_1$ ;  $G_1(\vec{r}_\Omega; \vec{r}')$  and  $G_2(\vec{r}_\Omega; \vec{r}')$  represent the free space Green's functions for the Poisson equation and the linearized Poisson-Boltzmann equation, respectively. Once  $\varphi(\vec{r}_\Omega)$  and  $\frac{\partial \varphi}{\partial n}(\vec{r}_\Omega)$  are determined, the reaction potential at the charge points can be calculated from:

$$\varphi_{reac}(\vec{r}_i) = \int_{\Omega_{bound}} \left[ G_1(\vec{r}_i; \vec{r}') \frac{\partial \varphi}{\partial n}(\vec{r}') - \varphi(\vec{r}') \frac{\partial G_1}{\partial n}(\vec{r}_i; \vec{r}') \right] d\vec{r}'. \quad (3)$$

Discretizing the interface surface  $\Omega_{bound}$  into triangular panels, and then using a piecewise constant collocation scheme produces a block system of equations from (1) and (2),

$$\begin{bmatrix} \frac{1}{2}I + \int_{panel_k} \frac{\partial G_1}{\partial n} d\vec{r}' & - \int_{panel_k} G_1 d\vec{r}' \\ \frac{1}{2}I - \int_{panel_k} \frac{\partial G_2}{\partial n} d\vec{r}' & 1/\epsilon \int_{panel_k} G_2 d\vec{r}' \end{bmatrix} \begin{bmatrix} \varphi \\ \frac{\partial \varphi}{\partial n} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{n_c} \frac{q_i}{\epsilon_1} G_1 \\ 0 \end{bmatrix} \quad (4)$$

where  $\varphi$  and  $\frac{\partial \varphi}{\partial n}$  are vectors of panel potentials and normal fields, respectively. For notational simplicity, we write the block matrices as follows:

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix} \begin{bmatrix} \varphi \\ \frac{\partial \varphi}{\partial n} \end{bmatrix} = \begin{bmatrix} -A_{1,3}q \\ 0 \end{bmatrix}. \quad (5)$$

Similarly, the discretized version of (3) can be written as

$$\varphi_{reac} = A_{3,1}\varphi + A_{3,2} \frac{\partial \varphi}{\partial n}. \quad (6)$$

Each of the blocks in (5) is square with dimension  $n_p$ , the number of panels used to discretize the surface. The blocks  $A_{3,1}$  and  $A_{3,2}$  are each of dimension  $n_c$  by  $n_p$ , and  $A_{1,3}$  is  $n_p$  by  $n_c$  in size.

Because the reaction potential is a linear function of the charge distribution, it can also be written  $\varphi_{reac} = Lq$ , where  $L$

is implicitly defined by (5) and (6). The  $i^{th}$  column of  $L$  can be determined by computing  $\varphi_{reac}$  from (5) and (6) assuming  $q = e_i$ , the  $i^{th}$  unit vector. Calculating the  $n_c$  columns of  $L$  for both the unbound and bound systems therefore requires  $2n_c$  simulations.

Iterative methods such as GMRES [11] are used to solve equation (5) in a matrix-free way, and accelerated methods [12] have been developed to quickly perform the required matrix-vector multiplications. Complete derivations of the above formulation are presented in [7], [8].

## 2.2 Free Energy Optimization

The atomic structures of the ligand, receptor, and ligand-receptor complex are taken as given, as is the receptor charge distribution. The ligand charge distribution  $q$  is varied to minimize the electrostatic free energy of binding. The objective function has been shown [5] to be convex and of the form

$$\Delta\Delta G_L^0 = q^T (L_{bound} - L_{unbound})q + d^T q. \quad (7)$$

The notation  $\Delta\Delta G_L^0$  denotes the change in free energy with respect to varying the ligand charge distribution only. To simplify the derivation of the implicit-Hessian technique, the quadratic term corresponding to the unbound system will be considered only at the end of the discussion.

Linear equality constraints are imposed on the ligand charge distribution to enforce charge conservation on various functional groups, and bound constraints are imposed on each charge value to ensure that the calculated charges are physically reasonable. We thus have a linearly constrained quadratic program:

$$\begin{aligned} \text{minimize} \quad & q^T L_{bound} q + d^T q \\ \text{subject to} \quad & A_c q = b \\ & \text{and} \quad m_i \leq q_i \leq M_i, \forall i \in \{1, \dots, n_c\} \end{aligned} \quad (8)$$

where  $A_c$  is a matrix of ones and zeros used to enforce sum of charge constraints on subsets of the vector  $q$ .

## 2.3 Primal-Dual Interior Point Methods

The standard form of a quadratic program is typically written as

$$\begin{aligned} \text{minimize} \quad & y^T Q y + d^T y \\ \text{subject to} \quad & A y = b \\ & \text{and} \quad y \geq 0 \end{aligned} \quad (9)$$

where  $Q$  is symmetric and positive definite [6], [13]. The convexity of the objective and the linear constraints meet a constraint qualification; that is, to find the global minimizer, it suffices to find a primal vector  $y^*$ , Lagrange multiplier vector  $\lambda^*$ , and dual slack vector  $s^*$  that together satisfy the

Karush-Kuhn-Tucker (KKT) conditions:

$$s^* = 2Qy^* + d - A^T\lambda^* \quad (10)$$

$$Ay^* = b \quad (11)$$

$$0 = y_i^* s_i^* \quad \forall i \in \{1, \dots, n\} \quad (12)$$

$$(y^*, s^*) \geq 0 \quad (13)$$

where  $n$  is the number of primal variables in the problem. Primal-dual interior point methods find a primal-dual solution  $(y^*, \lambda^*, s^*)$  by applying a specialized Newton-Raphson method to find the zeros of the function

$$F(y, \lambda, s) = \begin{bmatrix} 2Qy + d - A^T\lambda - s \\ b - Ay \\ Ys \end{bmatrix}. \quad (14)$$

The matrices  $Y$  and  $S$  are diagonal with diagonal entries equal to the corresponding elements of  $y$  and  $s$ . The Newton-Raphson steps are scaled to enforce the condition (13) at every iteration, and biased to improve convergence [6].

### 3 HESSIAN-IMPLICIT OPTIMIZATION

#### 3.1 Problem Transformation

The transformation of (8) into the standard form (9) is accomplished by introducing slack variables  $t$  and  $r$  such that

$$m + t = q, \quad t \geq 0 \quad (15)$$

$$q + r = M, \quad r \geq 0. \quad (16)$$

The problem then can be reduced into standard form using the substitution

$$y = \begin{bmatrix} t \\ r \end{bmatrix} \quad (17)$$

along with a number of other easily derived substitutions.

#### 3.2 Incorporating an Implicit Hessian

With the substitutions above, the  $(k+1)^{th}$  Newton step can be calculated by linearizing  $F(y^k, \lambda^k, s^k)$ ,

$$\begin{bmatrix} 2Q & -A^T & I \\ A & 0 & 0 \\ S^k & 0 & Y^k \end{bmatrix} \begin{bmatrix} \Delta y^{k+1} \\ \Delta \lambda^{k+1} \\ \Delta s^{k+1} \end{bmatrix} = \begin{bmatrix} -d + s^k - 2Qy^k + A^T\lambda^k \\ b - Ay^k \\ Y^k S^k e - \frac{(y^k)^T s^k}{n} e \end{bmatrix}. \quad (18)$$

The representation of  $Lq$  in (5) and (6) can be coupled with this system equations, and thus:

$$\begin{bmatrix} 0 & 0 & -A_c^T & -I & I & 2A_{3,1} & 2A_{3,2} \\ 0 & 0 & 0 & -I & 0 & 0 & 0 \\ A_c & 0 & 0 & 0 & 0 & 0 & 0 \\ I & I & 0 & 0 & 0 & 0 & 0 \\ S^k & 0 & 0 & 0 & Y^k & 0 & 0 \\ A_{1,3} & 0 & 0 & 0 & 0 & A_{1,1} & A_{1,2} \\ 0 & 0 & 0 & 0 & 0 & A_{2,1} & A_{2,2} \end{bmatrix} \begin{bmatrix} \Delta t^{k+1} \\ \Delta r^{k+1} \\ \Delta \lambda_c^{k+1} \\ \Delta \lambda_r^{k+1} \\ \Delta s^{k+1} \\ \Delta \phi^{k+1} \\ \Delta \frac{\partial \phi}{\partial n}^{k+1} \end{bmatrix}$$

$$= \begin{bmatrix} -d - 2L(m + t^k) + s_r^k + A_c^T \lambda_c^k + \lambda_t^k \\ s_r^k + \lambda_t^k \\ b - A_c m - A_c t^k \\ M - m - t^k - r^k \\ Y^k S^k e - \frac{(y^k)^T s^k}{2n_c} e \\ -A_{1,3} t - A_{1,1} \phi^k - A_{1,2} \frac{\partial \phi}{\partial n}^k \\ -A_{2,1} \phi^k - A_{2,2} \frac{\partial \phi}{\partial n}^k \end{bmatrix}. \quad (19)$$

Preconditioned GMRES is used to iteratively solve (19), and requires approximately the same time as solving (5) by itself. The full implementation also includes the gradient from the second quadratic term  $q^T L_{unbound} q$ ; the integral operators corresponding to the unbound system are coupled to (19) in an analogous manner to those for the bound system.

## 4 COMPUTATIONAL RESULTS

Test optimization problems were generated in the following manner: two concentric spheres were considered to be the bound and unbound molecular surfaces, and  $n_c$  charge locations were randomly selected within the unbound sphere;  $n_e$  random equality constraints were generated, and random box constraints  $m$  and  $M$  were also determined. The unbound surface was discretized into 124 panels, and the bound surface was discretized into 166 panels.

To verify that the Hessian-implicit method calculates the same optimal solution as a method which first calculates an explicit Hessian, a set of problems with varying problem dimension were solved using both methods. The convergence criterion for the optimization process was set to  $y^T s \leq 10^{-8} n_c$ . Figure 1 shows the norm of the error between the Hessian-explicit and Hessian-implicit based solutions, normalized by  $1/\sqrt{n_c}$ . It should be noted that the condition number of  $L = L_{bound} - L_{unbound}$  increases by 11 orders of magnitude between the smallest example problem and the largest, but the error does not grow in this fashion.

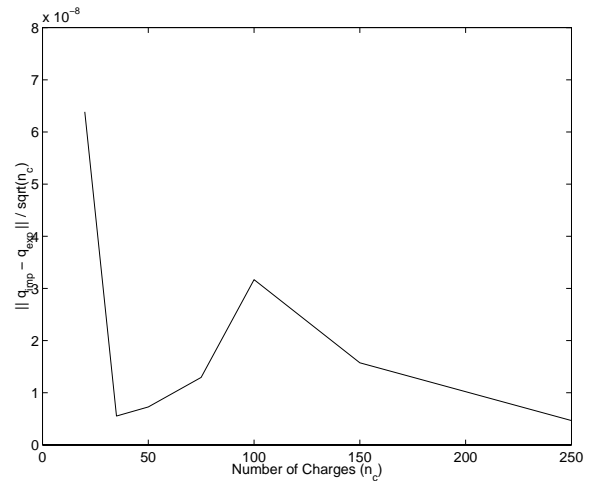


Figure 1: Verification of the Hessian-implicit method accuracy

The computational cost for the Hessian-implicit method grows very slowly with problem dimension, as Figure 2 illustrates. Because the integral operators dominate the augmented Jacobian in (19), and the primal-dual system of equations is very sparse, the cost metric used here is the number of integral operator matrix-vector products required to complete the optimization problem. Realistic problems may have as many as 3000 charges, which would make the Hessian-implicit method even more attractive as an accelerated optimization technique.

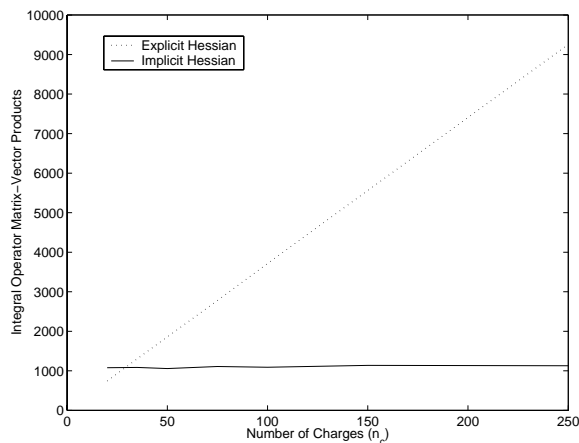


Figure 2: Computational scaling

## 5 CONCLUSION

This paper presented a Hessian-implicit primal-dual interior point optimization method that rapidly solves the charge optimization problem. The method couples two boundary element simulations (each of which computes one term of the objective function gradient) into a traditional primal-dual framework. Using the new optimization method, each Newton step requires approximately twice the computation of one electrostatic simulation. Primal-dual methods converge rapidly, and the number of iterations required scales exceptionally well with the number of charges; this results in an order of magnitude improvement in performance over the calculation of an explicit Hessian, even for moderately sized test cases. Future work will improve performance for repeated optimization (with varied constraint sets), and extend the formulation to allow nonlinear constraints.

## 6 ACKNOWLEDGMENTS

The authors are grateful for useful discussions with S. D. Senturia, J. L. Wyatt, G. Perakis, and D. Willis. This work was supported by the Singapore-MIT Alliance, the National Science Foundation, and J. Bardhan was supported by a Department of Energy Computational Science Graduate Fellowship.

## REFERENCES

- [1] L. Lee and B. Tidor. Optimization of electrostatic binding free energy. *Journal of Chemical Physics*, 106:8681–8690, 1997.
- [2] E. Kangas and B. Tidor. Electrostatic complementarity at ligand binding sites: Application to chorismate mutase. *Journal of Physical Chemistry*, 105:880–888, 2001.
- [3] E. Kangas and B. Tidor. Electrostatic specificity in molecular ligand design. *Journal of Chemical Physics*, 112:9120–9131, 2000.
- [4] L. Lee and B. Tidor. Optimization of binding electrostatics: Charge complementarity in the barnase-barstar protein complex. *Protein Science*, 10:362–377, 2001.
- [5] E. Kangas and B. Tidor. Optimizing electrostatic affinity in ligand-receptor binding: Theory, computation, and ligand properties. *Journal of Chemical Physics*, 109:7522–7545, 1998.
- [6] S. J. Wright. *Primal-Dual Interior Point Methods*. SIAM, 1997.
- [7] B. J. Yoon and A. M. Lenhoff. A boundary element method for molecular electrostatics with electrolyte effects. *Journal of Computational Chemistry*, 11:1080–1086, 1990.
- [8] S. S. Kuo, M. D. Altman, J. P. Bardhan, B. Tidor, and J. K. White. Fast methods for simulation of biomolecule electrostatics. *International Conference on Computer Aided Design (ICCAD)*, 2002.
- [9] J. G. Kirkwood. Theory of solutions of molecules containing widely separated charges with special application to zwitterions. *Journal of Chemical Physics*, 2:351, 1934.
- [10] C. Tanford and J. G. Kirkwood. Theory of protein titration curves I. general equations for impenetrable spheres. *Journal of the American Chemical Society*, 59:5333–5339, 1957.
- [11] Y. Saad and M. Schultz. GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems. *SIAM Journal of Scientific and Statistical Computing*, 7:856–869, 1986.
- [12] J. R. Phillips and J. K. White. A precorrected-FFT method for electrostatic analysis of complicated 3-D structures. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 16:1059–1072, 1997.
- [13] D. P. Bertsekas. *Nonlinear Programming: Second Edition*. Athena Scientific, 1999.