



# LUND UNIVERSITY

## **NMR structure determination of proteins supplemented by quantum chemical calculations: Detailed structure of the Ca<sup>2+</sup> sites in the EGF34 fragment of protein S**

Hsiao, Ya-Wen; Drakenberg, Torbjörn; Ryde, Ulf

*Published in:*  
Journal of Biomolecular NMR

*DOI:*  
[10.1007/s10858-004-6729-7](https://doi.org/10.1007/s10858-004-6729-7)

2005

*Document Version:*  
Peer reviewed version (aka post-print)

[Link to publication](#)

*Citation for published version (APA):*  
Hsiao, Y.-W., Drakenberg, T., & Ryde, U. (2005). NMR structure determination of proteins supplemented by quantum chemical calculations: Detailed structure of the Ca<sup>2+</sup> sites in the EGF34 fragment of protein S. *Journal of Biomolecular NMR*, 31(2), 97-114. <https://doi.org/10.1007/s10858-004-6729-7>

*Total number of authors:*  
3

*Creative Commons License:*  
Unspecified

### **General rights**

Unless other specific re-use rights are stated the following general rights apply:  
Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: <https://creativecommons.org/licenses/>

### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117  
221 00 Lund  
+46 46-222 00 00

# NMR structure determination of proteins supplemented by quantum chemical calculations: Detailed structure of the $\text{Ca}^{2+}$ sites in the EGF34 fragment of protein S

Ya-Wen Hsiao<sup>a</sup>, Torbjörn Drakenberg<sup>b</sup>, Ulf Ryde<sup>a\*</sup>

<sup>a</sup>*Department of Theoretical Chemistry and* <sup>b</sup>*Department of Biophysical Chemistry, Chemical Centre, Lund University, P. O. Box 124, S-221 00 Lund; \*Corresponding author*  
*Ulf.Ryde@teokem.lu.se, Phone: +46-46 222 45 02, Fax +46-46 222 45 43*

June 29, 2004

**Abstract.** We present and test two methods to use quantum chemical calculations to improve standard protein structure refinement by molecular dynamics simulations restrained to experimental NMR data. In the first, we replace the molecular mechanics force field (employed in standard refinement to supplement experimental data) for a site of interest by quantum chemical calculations. This way, we obtain an accurate description of the site, even if a molecular-mechanics force field does not exist for this site, or if there is little experimental information about the site. Moreover, the site may change its bonding during the refinement, which often is the case for metal sites. The second method is to extract a molecular mechanics potential for the site of interest from a quantum chemical geometry optimisation and frequency calculation. We apply both methods to the two  $\text{Ca}^{2+}$  sites in the epidermal growth factor-like domains 3 and 4 in the vitamin K-dependent protein S and compare them to various methods to treat these sites in standard refinement. We show that both methods perform well and have their advantages and disadvantages. We also show that the glutamate  $\text{Ca}^{2+}$  ligand is unlikely to bind in a bidentate mode, in contrast to the crystal structure of an EGF domain of factor IX.

**Abbreviations:** EGF: Epidermal growth factor; EGF34: epidermal growth factor-like domains 3 and 4 in the vitamin K-dependent protein S; cbEGF: calcium-binding EGF domains; MM: molecular mechanics; NOE: nuclear Overhauser effect; QM: quantum mechanics; rMD: restrained molecular dynamics; SANI: susceptibility anisotropy.

**Keywords:** NMR refinement, QM/MM methods, density functional calculations, EGF modules,  $\text{Ca}^{2+}$  sites



© 2004 Kluwer Academic Publishers. Printed in the Netherlands.

## 1. Introduction

Nuclear magnetic resonance (NMR) and X-ray crystallography are the two major sources of structural information for large biomolecules, such as proteins. Both methods have in common that they do not directly give a three-dimensional image of the structure. Instead, the structure is determined by an involved process of interpreting the experimental raw data. In crystallography, the problem is that the phases of the reflections are unknown. Approximate phases can be obtained from related crystal structures or from heavy-metal derivatives, and they are then improved by repeated cycles of model building and refinement of the structure (Kleywegt and Jones, 1995). In NMR structure determination, the raw data consist mainly of a number of estimated distances between pairs of atoms, constraints in dihedral angles, and hydrogen bonds (Cavanagh et al., 1996). These are converted to a three-dimensional structure by the use of distance geometry methods or restrained molecular dynamics (rMD).

Both methods have also in common that the experimental data is usually supplemented by empirical chemical data, typically in the form of a molecular-mechanics force field, with terms for the ideal geometry of bonds, angles, dihedrals, planar groups, chirality, and non-bonded interactions. The force field is used to ensure that the bond lengths and angles are chemically reasonable and that aromatic systems are planar.

As a consequence, the quality of the resulting structures will depend on the force field used in the structure refinement (Kleywegt and Jones, 1998; Nilsson et al., 2003). For standard amino acids and nucleic acids, accurate target values for bond lengths and angles exist (Engh and Huber, 1991). However, for more unusual molecules, such as substrates, inhibitors, coenzymes, and metal centres, i.e. *hetero-compounds*, experimental data are often incomplete or less accurate (Kleywegt and Jones, 1998). In particular, force constants are normally not available and the force field has to be constructed by the experimentalist, a complicated and error-prone procedure.

A conceivable way to solve these problems is to replace the force field for the site of interest by more accurate quantum chemical calculations: Density functional calculations with a medium-sized basis set typically reproduce experimental bond lengths within 0.02 Å for organic molecules and 0.07 Å for bonds to metal ions (Jensen, 1999; Sigfridsson et al., 2001; Olsson and Ryde, 2001; Ryde and Nilsson, 2003a), making them more accurate than standard low- and medium-resolution crystal structures. We have recently developed such a method, *quantum refinement* (Ryde et al., 2002), in which we replace the empirical force field for a small part of the protein in a standard crystallographic refinement by quantum chemical calculations. We have shown that it works properly and that it can be used to locally improve crystal struc-

tures of hetero-compounds, e.g. inhibitors and metal sites (Ryde and Nilsson, 2003a; Ryde et al., 2002).

In this paper, we show that a similar method can also be used to locally improve the results of NMR structure determinations. For such structures, this method has the additional advantage that it can be employed for sites for which the experimental data give little information about the structure, e.g. for metal sites. Therefore, we test the method for two calcium sites in the epidermal growth factor-like domains 3 and 4 in the vitamin K-dependent protein S. This is an ideal test case, because the  $\text{Ca}^{2+}$  ion is known to have flexible geometric preferences, binding to 6–8 ligands with variable Ca–ligand bond lengths (da Silva and Williams, 1991), which makes it very hard to describe by standard molecular mechanics methods. We show that the method works properly and that we can obtain a much more detailed picture of the calcium sites than with standard methods. We also test another method to automatically obtain a molecular mechanics force field for a site of interest from a theoretical frequency calculation (Nilsson et al., 2003).

## 2. Methods

### 2.1. HESS2FF AND COMQUM-N

Standard NMR refinement is performed as a restrained molecular dynamics (rMD) annealing scheme with an energy function of the type

$$E_{tot} = E_{MM} + E_{NMR}, \quad (1)$$

where  $E_{NMR}$  is the sum of all the NMR restraint energies (e.g. distance constraints based on nuclear Overhauser effect (NOE) data, dihedral constraints from the  $J$  couplings, hydrogen-bond restraints from amide proton exchange data, and susceptibility anisotropy (SANI) restraints from the residual dipolar couplings).  $E_{MM}$  is a standard MM energy function with bond, angle, dihedral angle, and non-bonded terms. Thus, the structure is obtained as a compromise between these two terms. The magnitude of the NMR restraints is arbitrary; therefore, each separate NMR term has a weight factor that determines the importance of this restraint relative to the MM restraints, which are in energy units.

QM data can be introduced into this energy in two ways. First, we can use QM calculations to construct MM parameters for the system of interest. There are many ways to perform such a parameterisation (Norrby and Liljefors, 1998). We have used a simple, fast, and automatic method, originally developed for the study of hetero-compounds in crystal structures (Hess2FF) (Nilsson et al., 2003), but it is directly applicable also to NMR systems. It extracts the ideal bond lengths, angles, and dihedrals, as well as



force constants from the Hessian matrix (i.e. the second derivative of the energy with respect to the coordinates) obtained from a QM optimised structure of a model of the interesting part of the protein. In this way, we get a more accurate description of the site of interest than a standard MM potential (if anyone exists at all), but there is still a risk that the site is not well determined at the MM level (which is likely for a  $\text{Ca}^{2+}$  site).

Alternatively and more accurately, we can replace the MM potential by a full QM calculation of the energy and the forces. Unfortunately, accurate QM methods can not yet be applied on a whole protein. Therefore, we have to restrict the QM calculations to a small but interesting part of the protein (e.g. a part that is poorly defined by the NMR restraints or not well described by the standard MM force field). This is done by partitioning the protein into two subsystems. System 1 consists of the site of interest that will be studied by QM methods, whereas system 2 contains the rest of the protein (and possibly parts of the surrounding solvent). We can then calculate the energy as:

$$E_{tot} = E_{QM1} - E_{MM1} + E_{MM} + E_{NMR}, \quad (2)$$

where  $E_{QM1}$  is the QM energy of the QM system,  $E_{MM1}$  is the MM energy of the same system, whereas  $E_{MM}$  and  $E_{NMR}$  have the same meaning as in Eqn. 1. The  $E_{MM1}$  term is needed to avoid double counting of energies in system 1 (i.e. to cancel the MM terms of system 1 in  $E_{MM}$ ).

This energy expression is similar to that used in the standard combined QM and MM method (QM/MM), which is one of the most popular ways to treat proteins with QM methods (Monard and Merz, 1999; Mulholland, 2001; Ryde, 2003; Ryde, 1996; Svensson et al., 1996):

$$E_{tot} = E_{QM1} - E_{MM1} + E_{MM}. \quad (3)$$

Thus, our approach can equivalently be seen as a QM/MM method restrained to fit the NMR data.

Special attention is needed if there is a covalent bond between the QM system and the surroundings protein (a junction). This is a well-known problem in QM/MM methods and a simple and robust solution (Nicoll et al., 2001) is to truncate the QM system with hydrogen atoms, the positions of which are linearly related to the corresponding carbon atom in the protein (Ryde et al., 2002; Ryde, 1996). Of course,  $E_{MM1}$  is also calculated with these hydrogen atoms, so that artefacts introduced by the hydrogen truncation may cancel. The forces are the negative gradient of the energy in Eqn. 2, taking into account the relation between the H and C junction atoms using the chain rule (Maseras and Morokuma, 1995).

In the quantum chemical calculations, system 1 is represented by its wavefunction and the rest of the protein is modeled by point charges that polarise the QM system in a self-consistent manner. In the MM calculations, all atoms

are described by the standard MM force field, but without any electrostatic interactions between the QM system and the surrounding protein, because they are already accounted for in the QM calculations. The electrostatic interactions within system 2 can be treated either in the QM calculations or in the MM calculations. In the former case, all interactions are considered, including interactions between bonded atoms (1–2, 1–3, and 1–4 interactions), which is normally not intended in the force field. In the latter case, the electrostatic interactions are treated in the intended way by the MM force field (typically ignoring 1–2 and 1–3 interactions and scaling down the 1–4 interactions by a constant factor). However, in order to obtain stable and reliable energies, a large cutoff distance need to be employed (ideally infinite), which is not always possible (many MM programs insist on calculating and storing a vector of all pairs of interactions to be included in the calculations, which may become too large to store in the internal memory). With an infinite cutoff, the two methods typically give very similar structures and relative energies.

We have implemented this energy expression (Eqn. 2) in a program, COMQUM–N, by constructing an interface between the QM software Turbomole 5.6 (Ahlrichs et al., 2000) and the free and widely used software Crystallography & NMR System (CNS), version 1.1 (Brunger et al., 2000). The interface is based on our QM/MM software COMQUM (Ryde, 1996; Ryde and Olsson, 2001). The philosophy behind this approach is that there should be no change in the code of the QM and MM/NMR software. Instead, COMQUM–N consists of a number of small programs which move information between the software, adding the forces and energies in a proper way.

An interface between Turbomole and CNS already exists in the quantum refinement software COMQUM–X (Ryde et al., 2002). However, this had to be slightly modified to allow for the treatment of also hydrogen atoms (these are normally not resolved in X-ray crystal structures) and for the inclusion of point charges in QM calculations (an electrostatic model without any hydrogen atoms is meaningless, because no hydrogen bonds or solvation can be described).

Moreover, procedures and input files had to be developed for the CNS calculations with NMR-based restraints. These were based on the CNS standard input file `anneal.inp` (dynamical annealing with NMR restraints, using rMD). This file was modified in a few ways (as is detailed in the web page [http://www.teokem.lu.se/~ulf/Methods/comqum\\_n.html](http://www.teokem.lu.se/~ulf/Methods/comqum_n.html)): First, an extra coordinate file has to be read and written containing the fourth to eighth decimals (standard CNS reads coordinates in PDB format, i.e. with three decimals, which is not enough for proper convergence) (Ryde et al., 2002). Second, all dynamics sections were disabled, so that only the final minimisation is performed. The reason for this is that we want to perform only a local optimisation of the site of interest at the end of the NMR refinement. It would

be a waste of computational power to perform QM calculations before the QM system was approximately assembled (i.e. before the QM groups are in proximity). In the energy and force calculations no optimisation is performed at all (the number of steps for the geometry optimisation is zeroed). Third, code was added to write out the energy terms and forces in separate files to be read by the COMQUM-N interface. Fourth, the non-bonded force field was slightly modified, as will be discussed below.

The program flow of COMQUM-N is shown in Scheme 1. It can be seen that COMQUM-N consists of five small interface routines that move information (energy, forces, charges, and coordinates) forth and back between the QM and NMR programs. In each cycle of the geometry optimisation, the geometry of system 1 is relaxed by the total forces, keeping the geometry of system 2 fixed. Then, the geometry of system 2 is relaxed by an extensive (not only a single step as for system 1) energy minimisation with NMR restraints, keeping system 1 fixed, performed by CNS (but still without any rMD annealing). This way, we take advantage of the fact that the NMR refinement is much faster than the QM calculation. This relaxation of system 2 is optional, but we see no reason not to perform it, because the NMR restraints will ensure that the structure is close to the true structure, and it will relieve strain that otherwise may build up between the QM system and the surrounding protein.

In practice, the five interface routines of COMQUM-N are divided into four separate programs: one core routine that is independent of the QM and MM programs, two input routines that construct input files to this core routine from the particular QM and MM programs in text files with a standard format, and one output routine that reads the output from the core routine and writes the data back into the specific QM or MM program (Ryde et al., 2002). In this way, the core COMQUM procedure becomes independent of the actual programs used for the QM or NMR calculations, which makes porting to other programs easier and more lucid.

## 2.2. APPLICATIONS ON EGF34

In order to test the performance of Hess2FF and COMQUM-N for NMR refinement, we have applied them to the two  $\text{Ca}^{2+}$ -binding sites in the epidermal growth factor-like domains 3 and 4 in the vitamin K-dependent protein S (EGF34). Modules homologous with epidermal growth factor (EGF) are common in extracellular proteins. They are found in wide variety of animal proteins: connective tissue fibres, complement, blood coagulation and fibrinolytic proteins, as well as proteins involved in cell morphogenesis (Appella et al., 1988; Campbell and Bork, 1993; Stenflo et al., 2000). In fact, this module is the fourth largest protein family, present in 1% of the human proteins (Henikoff et al., 1997). The EGF modules are independently folding domains that usually consist of 40–50 amino acids and three disulphide

bridges. They are often involved in protein–protein interactions that are  $\text{Ca}^{2+}$  dependent. Thus, a subset of the EGF motifs bind a  $\text{Ca}^{2+}$  ion in a conserved sequence. The Ca site is typically 6–7 coordinate, involving five residues from the protein (two back-bone amide groups and three Asp, Asn, Glu, or Gln residues) and a water molecule.

Some proteins contain many EGF modules in tandem repeats (Stenflo et al., 2000). Interestingly, the Ca affinity of such multimers is often higher than for the isolated modules. For example, in the pair of EGF modules 3 and 4 in protein S, module 3 has approximately the same Ca affinity as the isolated module 3, whereas the affinity of the fourth module is 8600 times larger in the pair than in the isolated module 4 (Stenberg et al., 1997a; Stenberg et al., 1997b).

Structures of many EGF-module protein are known, both from NMR and crystallographic studies (Stenflo et al., 2000). In particular, a 1.5 Å crystal structure of the EGF-like domains in human clotting factor IX has been presented (Rao et al., 1995). It shows two EGF domains, each binding a  $\text{Ca}^{2+}$  ion in a pentagonal bipyramidal manner (one carboxylate group binds bidentately).

The anticoagulant cofactor protein S has four EGF-like domains in tandem. Domains 2–4 are calcium-binding EGF domains (cbEGF). The  $\text{Ca}^{2+}$  binding to these domains are much stronger than to any other cbEGF studied so far. The smallest fragment with high Ca affinity is EGF34. We have used NMR to determine the three-dimensional structure of this protein fragment, hoping that it would reveal the reason to the high  $\text{Ca}^{2+}$  affinity. Standard multidimensional NMR has been used to estimate H–H distances and rMD has been used to obtain structures in agreement with the NMR distance restraints. Experimental details, as well as a discussion of the general structure, and its relation to the high  $\text{Ca}^{2+}$  affinity will be published elsewhere (Drakenberg et al., 2004). This paper is restricted to a discussion of various methods to treat the  $\text{Ca}^{2+}$  sites.

The putative Cys–Cys bridges and  $\text{Ca}^{2+}$ -binding sites in EGF34 were identified from the consensus sequence of the EGF domains (Stenflo et al., 2000) and by comparison with the EGF domains in clotting factor IX (Rao et al., 1995): It was assumed that disulphide bridges are formed by the Cys residues 164–176, 171–185, 187–200, 206–215, 211–224, and 226–241. Likewise, we assumed that the two  $\text{Ca}^{2+}$ -binding sites are formed by residues Asp–160, Val–161, Glu–163, Asn–178, and Ile–179, as well as by Asp–202, Ile–203, Glu–205, Asn–217, and Tyr–218 (the second and fifth residues of each site bind by the back-bone amide oxygen atom, whereas the others bind by the side chains). In addition, both  $\text{Ca}^{2+}$  sites were assumed to bind one water molecule.

In the COMQUM–N calculations, the protein is divided into two parts. System 1 consisted of  $\text{Ca}(\text{CH}_3\text{COO})_2(\text{CH}_3\text{CONHCH}_3)_2(\text{CH}_3\text{CONH}_2)(\text{H}_2\text{O})$

(the same for both  $\text{Ca}^{2+}$  sites) and it was treated by quantum chemistry, whereas the rest of the protein (system 2) was treated entirely with standard NMR refinement methods. Thus, there were twelve junctions between the two systems: one each for the Asp, Glu, and Asn residues, and four for the two backbone groups ( $\text{N}$  and  $\text{C}^\beta$  for the residue containing the CO group and  $\text{C}$  and  $\text{C}^\beta$  for the next residue, containing the CO group). The second backbone model in both sites contains an additional junction, because the next residue is Pro, giving a junction also for the  $\text{C}^\delta$  atom.

The parameters for the junctions were exactly the same as for the original amino acids, except for the bonds to the junction hydrogen atom. This bond length was taken from the structure of the same fragment ( $\text{CH}_3\text{COO}^-$ ,  $\text{CH}_3\text{CONHCH}_3$ , or  $\text{CH}_3\text{CONH}_2$ ) optimised with the quantum chemical method. The force constant was calculated as the force constant of the original bond times the square of the quotient of the ideal original bond length (from the force field libraries) and the ideal bond length of the junction. This way, the forces of the bond with and without the junction will be equal. In addition, a few improper dihedral angles involving both QM and junction atoms had to be removed (they will not cancel between  $E_{MM1}$  and  $E_{MM}$ , owing to the movement of the junction atoms).

### 2.3. COMPUTATIONAL DETAILS

The QM calculations were performed by density functional theory, using the Becke–Perdew-1986 exchange–correlation functional (BP86) (Becke, 1988; Perdew, 1986) and the standard medium-sized 6-31G\* basis set for all atoms (Hefre et al., 1986). Only the five pure  $d$ -type functions were used. The calculations were sped up by expansion of the Coulomb interactions in auxiliary basis sets, the resolution-of-identity approximation (Eichkorn et al., 1995; Eichkorn et al., 1997). These calculations were performed by Turbomole 5.6 (Ahlrichs et al., 2000). Such a method is known to give accurate and nearly converged geometries for metal-containing systems (Ryde and Nilsson, 2003a; Siegbahn and Blomberg, 2000; Ryde et al., 2001). The COMQUM–N optimisations were performed in two steps: First, system 2 was allowed to relax and the full geometry was optimised until the change in energy between two iterations was below  $10^{-4}$  Hartree and the maximum norm of the gradients was below  $10^{-2}$  atomic units. Then, system 2 was fixed and the structure was further optimised with stricter convergence criteria,  $10^{-6}$  Hartree (2.6 J/mole) and  $10^{-3}$  atomic units (1.4 kJ/mole/Å, i.e. the default criteria in Turbomole). In the latter calculations, the maximum allowed movement of any atom (dqmax) was reduced to 0.03 atomic units (0.016 Å). Otherwise, extensive oscillations were often seen, although the same minimum and energy was normally obtained, but after many more iterations.

In some calculations, solvation effects were estimated using the continuum conductor-like screening model (COSMO) (Klamt and Schüürmann, 1993; Schäfer et al., 2000). These calculations were performed with default values for all parameters (implying a water-like probe molecule) and a dielectric constant ( $\epsilon$ ) of 80. For the generation of the cavity, a set of atomic radii have to be defined. We used the optimised COSMO radii in Turbomole (130, 200, 183, and 172 pm for H, C, N, and O, respectively, and 200 pm for  $\text{Ca}^{2+}$ ) (Klamt et al., 1998).

The CNS calculations (Brunger et al., 2000) were performed with the standard protein-allhdg, water, and ion topology and parameter files. All protein atoms (including hydrogen atoms) were included in these calculations, but no water molecules (except the two  $\text{Ca}^{2+}$  ligands in some calculations). As mentioned above, the NMR calculations were performed with the file `anneal.inp`. In this file, we used the default values for most parameters. In particular the final weight factors for the NOE and dihedral constraints were 75 and 400 and the final force constant for the SANI restraints was 1.0. When the protein was relaxed, ten cycles of final minimization consisting of 200 steps were run. In each standard rMD calculation with CNS, 200 structures were obtained from an extended structure, using random starting velocities.

In the calculations we included four types of NMR restraints, viz. 813 NOE distance restraints, 30 hydrogen bond restraints, 89 dihedral restraints, and 43 susceptibility anisotropy restraints. NOE restraints from methyl groups, degenerate methylene groups, and ambiguous assignments were averaged using the default sum mode. A final SANI force constant of 1.0 kcal/mole was used, because it resulted in calculated residual dipolar couplings matching the experimental ones within experimental errors. The SANI coefficients were optimised through a grid search:  $a_0 = -0.0601$ ,  $a_1 = -15$ , and  $a_2 = 0.35$ .

#### 2.4. CALIBRATION OF THE QM METHOD

We began the investigation with some calibrations of the QM method. To this end, we started from the  $\text{Ca}^{2+}$  site 1 in the crystal structure of the EGF-like domain in human clotting factor IX (Rao et al., 1995). The site was truncated in the same way as in the COMQUM-N calculations (i.e. to  $\text{Ca}(\text{CH}_3\text{COO})_2(\text{CH}_3\text{CONHCH}_3)_2(\text{CH}_3\text{CONH}_2)$ ) and a water molecule was added (there is an obviously empty coordination site in the crystal structure). Then, this structure was optimised with a number of different methods and basis sets. In addition, we also optimised a number of structures starting from NMR structures of the two  $\text{Ca}^{2+}$  sites in EGF34. The most interesting results of these calculations are collected in Table S1.

It can be seen that the structure of the  $\text{Ca}^{2+}$  site is quite insensitive to the QM method. The Ca-O distances change by less than 0.01 Å when

the method is changed from BP86 to B3LYP (Hertwig and Koch, 1997). An increase of the basis set from 6-31G\* to the appreciably larger 6-311+G(2d,2p) (Hehre et al., 1986) has a somewhat larger effect on the Ca–O distance, viz. a contraction by 0.02–0.05 Å. Inclusion of a continuum solvent (COSMO model) with a dielectric constant of 80 (similar to water) changes the Ca–O distances by 0–0.03 Å in a somewhat erratic manner. Different starting structures had a similar effect on the site: The individual Ca–O distances differ by up to 0.07 Å, but the average is essentially the same, 2.41 Å.

However, if the results are compared to the crystal structure of the EGF-like domain in human clotting factor IX, it can be seen that most of the optimised structures end up in six-coordinate Ca<sup>2+</sup> sites with both carboxylate groups binding in a monodentate manner, whereas one of the carboxylate groups bind bidentately in the crystal structure. Some NMR structures also ended up in a bidentate structure and from these, it can be concluded that such a binding leads to an increase in the Ca–O distance of the carboxylate group from 2.34–2.41 Å to 2.51–2.55 Å. The average Ca–O distance also increases to 2.48–2.49 Å.

The change between mono- and bidentate binding of carboxylate groups (so called carboxylate shifts) has been studied in several other systems, e.g. for Zn<sup>2+</sup> and binuclear iron sites (Ryde, 1999; Torrent et al., 2001). From these studies, it is clear that there is only a minor energetic difference between mono- and bidentate binding. Therefore, the binding mode is mainly determined by what interactions the non-bonding carboxylate oxygen atom can form in the monodentate state. This is also observed in the present structures. In the monodentate sites, the two carboxylate groups form strong hydrogen bonds to the water ligand and to the amide hydrogen atoms of the Asn ligand. In the crystal structure, the latter hydrogen bond is retained, whereas the water molecule and the bidentate carboxylate group is exposed to solvent, where more ideal hydrogen bonds can be provided by the surrounding water molecules. Such effects can be simulated in the calculations by adding a water molecule in the second coordination sphere of the Ca<sup>2+</sup> ion. This also led to bidentate structures (Table S1).

Considering that the crystal structure is bidentate, it is somewhat alarming that it has an average Ca–O distance of 2.40 Å, which is more similar to the monodentate than to the bidentate optimised structures. This is partly an effect of the missing water ligand, which has a 0.01–0.06 Å longer Ca–O distance than the average value in all monodentate structures. It is also partly an effect of the basis set (~0.03 Å). In order to check if the remaining difference is caused by the uncertainty in the crystal structure (typically at least ~0.1 Å (Nilsson et al., 2003; Fields et al., 1994; Cruickshank, 1999) or by systematic errors in the QM method, we also optimised the structure of Ca<sup>2+</sup> in water. Experimentally, it is known that Ca<sup>2+</sup> on average has eight ligands in water with a distance of 2.48 Å (Jalilvand et al., 2001). A QM optimisation

of  $\text{Ca}(\text{H}_2\text{O})_8^{2+}$  in  $D_{4d}$  symmetry (to avoid internal hydrogen bonds between the water molecules) gave a Ca–O distance of 2.47 Å with the BP86/6-31G\* method and 2.49 Å with the 6-311+(2d,2p) basis set, i.e. both in excellent agreement with experimental data. On the basis of these results, we decided to use the BP86/6-31G\* method, which is much faster than with the larger basis set. However, it should then be kept in mind that this method overestimates the Ca–O distances by  $\sim 0.03$  Å.

### 3. Results and Discussion

There are several ways to treat a metal site in standard NMR structure determination by restrained molecular dynamics (rMD) simulations. First, the metal site can be totally ignored, using restraints only from the NMR raw data. This should give a structure as close as possible to the NMR data, but still bring the metal ligands in proximity, if the site is well-defined by the NMR data. However, it would not give any information about the binding mode of the ligands or the detailed structure of the metal site. Second, if the metal ligands are known beforehand, metal–ligand distances could be included as normal NOE distance restraints, using reasonable estimates of the bond lengths. This seems to be the most common method in standard structure determination (Downing et al., 1996; Saha et al., 2001; Wang et al., 2001; Tossavainen et al., 2003). This should give an improved structure of the metal site. Third, the metal–ligand interactions could be described by an MM potential like the surrounding protein, rather than by NOE restraints. This should give a much more accurate description of the details of the metal site. However, accurate MM potentials for metal ions are hard to construct, especially if the actual number and geometry of the ligands are not known or may change.

In the following sections, we will first test these three approaches for the two  $\text{Ca}^{2+}$  sites in EGF34 to see how the predicted structure of the  $\text{Ca}^{2+}$  site changes and how well the NMR restraints can be fulfilled. This is important to ensure that we do not enforce a site that violates the NMR raw data. It should be remembered that we do not have any direct experimental evidence of the actual  $\text{Ca}^{2+}$  ligands – the suggested ligands, mentioned above come simply from sequence alignments (Stenflo et al., 2000). We will then see if we can improve the structure of the  $\text{Ca}^{2+}$  site by the use of COMQUM–N.

#### 3.1. REFINEMENT WITHOUT ANY CALCIUM RESTRAINTS

We first performed an rMD annealing without any restraints for the  $\text{Ca}^{2+}$  ion (i.e. we employed only the standard NMR restraints and no quantum chemistry or MM potential for  $\text{Ca}^{2+}$ ). Consequently, the structure contained



neither  $\text{Ca}^{2+}$  nor any water molecules. In total, 200 different structures were obtained in this way using random starting velocities. The refinement consisted of a high-temperature dynamics, two slow-cool annealings, and a final minimisation. The dynamics and the first annealing simulations were performed in torsional space and a soft repulsion potential was used in all simulations.

The results of these calculations (Table S2) show that the two  $\text{Ca}^{2+}$  sites have poor geometries in all the structures obtained. For example, for site 1, only one structure has a maximum O–O distance (for the Ca ligands) shorter than 10 Å (8.8 Å; it is 4.8 Å in the crystal structure). The lowest maximum Ca–O distance for any structure is 6.4 Å, which is much larger than for a typical  $\text{Ca}^{2+}$  site (it is 2.6 Å in the crystal structure; we here assume that the  $\text{Ca}^{2+}$  ion resides at the midpoint between the two carbonyl ligand atoms – in the crystal structure of the EGF domain in factor IX, these two atoms are on opposite sides of the  $\text{Ca}^{2+}$  ion with an O–Ca–O angle of 172–175°). In addition, the total energy of the best Ca structures is very high. For site 2, the situation is somewhat better, with carbonyl distances almost half as long as for site 1. However, even the best structures have a maximum O–O distance of 5.6 Å, a maximum Ca–O distance of 3.1 Å, and high total energies.

The reason why site 2 is better defined by the NMR data than site 1 is that site 1 is at one end of the protein, whereas site 2 is in the middle of the protein, as can be seen in Figure 1. However, the main conclusion from this section is that the NMR data alone does not lead to any reasonable  $\text{Ca}^{2+}$  sites. In particular, it is impossible to speculate about any details of the sites, e.g. whether the carboxylic groups are bidentate or how many water molecules may coordinate to the site. Finally, we can also note that there is an appreciable spread in the energies obtained for the various structures, both the total energies and the energies of the various NMR terms. For example, the average total and NMR energies of the best 20 structures are 654 and 92 kJ/mole, whereas the corresponding energies for the best structure are 529 and 72 kJ/mole. This is important to remember when judging the effect of various  $\text{Ca}^{2+}$  restraints.

### 3.2. REFINEMENT WITH O–O RESTRAINTS

Next, we tried to define the  $\text{Ca}^{2+}$  site by including a set of ten O–O restraints between the five putative protein  $\text{Ca}^{2+}$  ligands for each site, defined in the same way as normal NOE restraints, with a flat-bottomed (between 3.0 and 5.1 Å) harmonic potential. Thus, we still did not include any  $\text{Ca}^{2+}$  ion or water molecules in the calculations.

Quite naturally, the O–O restraints ensure that all ligand atoms are relative close in space, but in all structures, the carbonyl O–O distances are around the upper limit of the restraints, 5.1–5.2 Å and the maximum O–O distance is

even longer 5.4–5.6 Å (Table S3). This is also reflected by the maximum Ca–O distance (the Ca<sup>2+</sup> was inserted in the middle of the carbonyl O–O bond, as before), which is over 3.0 Å for all structures of site 1 and 2.8 Å for site 2. Thus, all Ca<sup>2+</sup> sites still have effectively lost at least one ligand. Moreover, in most of the structures, some of the ligand oxygen atoms are not directed towards the putative centre of the Ca<sup>2+</sup> site, as can be seen in Figure 2.

The energies of these structures are slightly larger than for those without any Ca<sup>2+</sup>-related restraints, e.g. by 205 kJ/mole for the average of the total energy for the 20 best structures. However, most of this difference comes from the MM energy: The difference in the total NMR energy is only 42 kJ/mole, originating mainly from the NOE term (36 kJ/mole). This can partly be explained by the additional O–O restraints for the Ca<sup>2+</sup> sites, which are included in this term. It is a shortcoming of this method that the experimental data and empirical restraints are both mixed into the NOE term, so that it cannot be clearly determined how much the new restraints have affected the fit to the experimental data. However, it is notable that many of the best structures of the Ca<sup>2+</sup> sites are also among the structures with the lowest energy. This shows that good sites are not unnatural.

The general structure of the protein obtained with these restraints (Figure S1) is similar to that obtained without any constraints. However, the Ca<sup>2+</sup> sites, especially site 1, is somewhat better defined with these constraints. Thus, we can conclude that this seems to be a better method to obtain reasonable structures of the protein than without any Ca-restraints, but it still does not give any good geometries of the Ca<sup>2+</sup> site.

### 3.3. REFINEMENT WITH CA–O NOE RESTRAINTS

In order to improve the structure of the Ca<sup>2+</sup> site and get better starting points for the COMQUM–N calculations, we decided to introduce the two Ca<sup>2+</sup> ions and two water ligands in the refinement calculations. First, we tried to describe the Ca–O interaction with a flat-bottom potential, similar to the NOE restraints. This seems to be the most common way to treat a Ca<sup>2+</sup> ion in NMR structures (Downing et al., 1996; Saha et al., 2001; Wang et al., 2001). The potential was zero between 2.0 and 3.0 Å and harmonic outside this range. However, this did not give any satisfactory results (Table S4): In all the obtained structures, the maximum Ca–O distance was 3.0 Å (the upper limit of the flat bottom). Thus, all sites have effectively lost at least one ligand and show little variation. Of course, we could have cured this problem by making the flat bottom tighter, but we would still only get what we put in (the upper limit), without any physical relevance of the results. Moreover, this approach, like the previous one, has the shortcoming of mixing up experimental data and the empirical potential, because the Ca<sup>2+</sup> sites are described by NOE

restraints and the corresponding energies will appear in the NMR term, rather than in the MM term.

### 3.4. REFINEMENT WITH A BONDED MM CA–O POTENTIAL

Therefore, we decided to use another approach, where the Ca–O interactions are described by a standard MM potential. The potential was obtained by the program Hess2FF (Nilsson et al., 2003) from QM vacuum optimisations and frequency calculations of the two Ca<sup>2+</sup> sites in the EGF domain of factor IX. We used the structures in Table S1 and took force constants as the average of the similar interactions (i.e. for water, carboxylate, and carbonyl groups). The ideal bond lengths and force constants used are listed in Table S5. We decided to use only the bonded terms (i.e. no angle or dihedral restraints), because we did not want to bias the results towards any particular coordination number or geometry and also because the vacuum structure is somewhat distorted by interactions between the carboxylate groups and the methyl groups (an unavoidable vacuum effect).

The carboxylate groups pose a special problem because they can bind to Ca with either both or only one of the carboxylate oxygens. In the crystal structure, one of the carboxylate groups in each site binds in the bidentate mode, whereas the other binds monodentately. We decided to test both these possibilities in our calculations and therefore designed two sets of parameters, one for a monodentate site, based on the structure in the first row of Table S1, and the other bidentate, based on the sixth row in the same table (the calculation with an additional water molecule).

This approach gave excellent Ca<sup>2+</sup> sites in all structures (Table S6). For both sites, the monodentate parameters gave an average maximum Ca–O bond length of 2.6–2.7 Å for all structures. The shortest maximum Ca–O bonds were 2.45 and 2.43 Å for the two sites, i.e. similar to what is found in the crystal structure. Likewise, the carbonyl O–O distances also show a quite restricted variation (averages 4.9–5.2 Å, slightly shorter for site 2 than for site 1; almost the same values were obtained for the maximum O–O distances).

The energies are similar to those obtained with the O–O restraints: The total energy of the best structure is 1 kJ/mole lower, but the average values of the total and the NMR energies in the 20 best structures are slightly higher. In particular, it seems that the dihedral terms have increased slightly.

The calculation with parameters for a bidentate binding of Glu–163/205 gave slightly worse results (Table S7): The average maximum Ca–O distances are ~0.2 Å longer in the bidentate structure, whereas in the best structures for each site, the difference is smaller (from the force field, the difference should be 0.09 Å). Likewise, the best and average energies are also larger for the bidentate site, by 40 kJ/mole for the total energy and by 11 kJ/mole for the

NMR energy, this time originating mainly from the NOE term (average of the 20 best structures).

There is also the possibility that it is the Asp-160/202 residues that bind in a bidentate mode, although this is not observed in the crystal structure of the human clotting factor IX. Therefore, such a coordination was also tested. This gave actually slightly lower total and NMR energies than both the bidentate Glu sites and the monodentate sites (Table S8). For example, the average total and NMR energies of the 20 best structures were 880 and 121 kJ/mole for the bidentate Asp sites, whereas they were 946 and 144 kJ/mole for the monodentate sites. On the other hand, the Ca-O distances are appreciably longer in the bidentate Asp sites: The average maximum Ca-O distance among the 20 best structures is 2.78 and 2.90 Å for the two sites, whereas it is only 2.60 and 2.68 Å for the monodentate site. Once again, this is larger than what would be expected from the respectively equilibrium bond lengths in the force field (cf. Table S5).

In conclusion, bonded Ca-O terms of MM type, seems to be an excellent method to obtain reasonable structures for the Ca<sup>2+</sup> sites. The resulting structure is at least as well determined as with the O-O restraints (Figure S2). The bidentate parameters for the Asp residues seem to give a slightly better site than the other two possibilities, but the NMR energies are not very different.

### 3.5. CALIBRATION OF THE COMQUM-N METHOD

So far, we have only used the QM data to construct an MM potential of the Ca<sup>2+</sup> site. Of course, much information is lost by this conversion and there is always the risk that errors are introduced, especially if the QM calculation is performed on a structure that is different from the final NMR structure. Moreover, a standard MM potential, such as the one in CNS does not allow the coordination number to change during the refinement. Therefore, we enforce a certain structure when we set up the MM potential, and this may not change during the refinement. Thus, we cannot model the dissociation of a ligand or the transition from mono- to bidentate binding of a carboxylate group.

All these problems can be avoided by using the QM calculations directly in the refinement, as in the COMQUM-N method. However, we first have to test out the method and decide how it is optimally used. Therefore, we performed a number of test calculations, using various starting structures. We looked especially at four issues: the choice of repulsive parameters, the number of MM iterations, the weight of the NMR restraints, and the use of electrostatics in the QM and MM calculations.

In standard NMR refinements, CNS uses a soft repulsive potential, which allows atoms to go through each other. Unfortunately, this potential frequently led to failures in COMQUM-N, because QM atoms ended up very

close to MM atoms. For this reason, and also because it has been shown that such soft potentials, when employed also in the final minimisation of the refinement, may lead to poor structures in terms of the Ramachandran plots (Doreleijers et al., 1999), we decided to instead use the standard van der Waals (Lennard–Jones) potential of CNS (the default method for X-ray refinement). In addition, we used an infinite cut-off to avoid instabilities.

Second, we looked at the optimum number of iterations in the NMR-restrained minimisation (cf. Scheme 1). In CNS, default number is 2000 (10 cycles of 200 iterations). However, in the combination of crystallographic refinement and QM calculations, divergence was observed unless only one iteration was used. However, this is not the case with COMQUM–N (Table S9; it should be noted that the calculations in Tables S9–S11 used slightly different NMR restraint than in the other tables; therefore the NMR energies are larger). On the contrary, the convergence was faster with many iterations. Likewise, the total energy was lower. It turned out to be favourable to start the calculation with a normal NMR-restrained MM minimisation (without any QM) of the protein to convergence ( $\sim 10\,000$  iterations), using the standard van der Waals parameters. Still, it can be seen that the results are not fully converged until 40 000 steps of MM minimisations are used. However, if that many steps are allowed, most of the time is spent in the MM minimisations and the full optimisations will take a very long time (several weeks). Therefore, we decided to use the default 2000 MM steps, for which the Ca–O distances are converged to within 0.01 Å and the total and NMR energies within 3 and 1 kJ/mole, respectively.

Third, we tested the effect of changing the weight of the NMR restraints. In COMQUM–X, the results strongly depend on the relative weight between crystallography and the MM and QM energy functions (Ryde et al., 2002). For COMQUM–N, the effect seems to be less pronounced: The average Ca–O bond length does not change at all (Table S10), whereas the individual distances change by up to 0.05 Å. Of course, the energy terms change more when increasing the NOE weight from 75 to 750. The NOE energy decreases by almost a factor of four, whereas the other two terms increase slightly. However, the MM energy increases even more, so that the total energy increases by 446 kJ/mole. The energy of the quantum system changes by less than 10 kJ/mole (data not shown). Therefore, we see no reason to modify the NMR weights from their default values.

Finally, we also tested the treatment of electrostatics in the COMQUM–N calculations. By default, CNS ignores all electrostatic interactions in NMR refinement. However, in the QM calculations, electrostatics within the QM system is included. We can then choose to include only these electrostatic interactions, include also the polarisation of the QM system by the surrounding protein (which is the standard choice in QM/MM optimisations), or even to turn on the electrostatics also in the NMR-restrained minimisation. We tried

all three alternatives for several systems (Table S11 shows two typical cases). However, we always saw a strong increase of the NMR energies if electrostatics were included in NMR-restrained minimisation. This is in accordance with the consensus that NMR refinement should be run without electrostatics, unless the protein is explicitly solvated.

The other methods gave quite similar results with NMR energies within 10 kJ/mole. However, it was invariably observed that calculations without any point charges in the QM calculations gave a lower energy than with the point charges (by 4–7 kJ/mole). Furthermore, calculations where the QM system is dispersed into a continuum solvent (COSMO method (Klamt and Schüürmann, 1993)) gave even lower NMR energies (by 1–5 kJ/mole). We doubt that these results are general, because the energies involved are so small and because the effect of point charges and COSMO should depend on the detailed structure of the surroundings. For water-exposed sites, as the present  $\text{Ca}^{2+}$  sites, it is possible that continuum calculations with a dielectric constant of  $\sim 80$  may improve the results. However, for sites that are buried inside the protein and interacts with the surroundings with many hydrogen bonds, it is likely that a point-charge model would be the best choice.

For the general use of COMQUM–N, our best recommendation is to run without any point charges in the QM calculation, unless there is extensive hydrogen bonding to the site of interest, and without any continuum model. This would be in accordance with the treatment of the surrounding protein. However, the optimum solution would probably be to include electrostatics in all calculations (i.e. both for the QM system and the surrounding protein). Then it would also be necessary to include full solvation of the protein by explicit water molecules. CNS is not set up and calibrated for such calculations, whereas other programs, e.g. AMBER (Case et al., 2002), allow for such an approach.

### 3.6. REFINEMENTS WITH THE COMQUM–N METHOD

After this calibration of COMQUM–N, we run production calculations on EGF34. As for standard NMR refinement, our aim is to obtain an ensemble of possible structures of the  $\text{Ca}^{2+}$  sites. Therefore, we started COMQUM–N calculations for each of the two  $\text{Ca}^{2+}$  sites from the ten best structures (in terms of total energy) in Tables S6, S7, and S8, i.e. for both the mono- and bidentate sites. Of course, we could also have started from structures obtained with other methods, but for structures obtained without or with only O–O restraints, this would have been waste of computer resources, because the starting structures are too poor and there is essentially no force in QM to assemble the  $\text{Ca}^{2+}$  site if the ligands are far away.

For the same reason, we performed only the final minimisation of the  $\text{Ca}^{2+}$  site, i.e. no high-temperature dynamics was run. Such a local optimisation

of the  $\text{Ca}^{2+}$  site in an ensemble of structures obtained by standard NMR refinement exploit the computer resources in the best way; if the COMQUM-N method had been used already in the early phases of the refinement, most of the structures would have ended up with unrealistic  $\text{Ca}^{2+}$  sites (like those in Table S2) at a very high computational cost. Therefore, it is more favourable to produce reasonable starting structures for the  $\text{Ca}^{2+}$  sites with the MM-restraint methods and then perform a final minimisation of the  $\text{Ca}^{2+}$  site with COMQUM-N. Provided that the weight factors are appropriate, this will still allow for significant modifications of the sites, if necessary.

The results show that in many of the final structures, the coordination has changed (Tables S12 and S13). In five of the calculations starting from the monodentate site, the final structure is bidentate (four with Asp and one with Glu). Moreover, in seven of the calculations (mostly for site 2), one of the carbonyl groups dissociate (to Ca-O distance of 3.04–3.67 Å; 4.46 Å in one case), but in two cases this is compensated by the bidentate binding of Asp. Likewise, only four of the bidentate Glu structures keep all the seven ligands, whereas four of them become monodentate, seven become six-coordinate with the Glu ligand still bidentate, two become six-coordinate with the Asp ligand bidentate, three become five-coordinate, and two actually become seven-coordinate with both the Asp and Glu ligand bidentate, but with the water ligand dissociated. Finally, seven of the bidentate Asp sites keep all the ligands, whereas two become monodentate, nine lose one ligand, and two lose two ligands (not Asp). Thus, in total, there are 14 monodentate structures, 5 bidentate structures with Glu, 11 bidentate structures with Asp, 19 bidentate structures with one lost ligand (12 with Asp and 7 with Glu), 9 five-coordinate structures, and 2 seven-coordinate structures with both Asp and Glu bidentate. Typical examples of the monodentate and bidentate structure with Asp are shown in Figure 3.

The average Ca-O distances follow the coordination number of the site: The five-coordinate sites have an average Ca-O distance of 2.38 Å, the six-coordinate sites have average Ca-O distances of 2.42–2.44 Å, and the seven-coordinate sites have average Ca-O distances of 2.49–2.51 Å (in all these averages, the dissociated ligands have been excluded, in variance to the averages in Tables S12 and S13). The Glu and Asp residues give the shortest Ca-O bonds (average 2.36 and 2.37 Å; the shortest bond encountered was 2.23 Å for Glu in a five-coordinate site). The Asn ligand gives slightly longer distances (2.39 Å), whereas water and the two carbonyl groups give the longest bonds (averages 2.49–2.57 Å). The longest bond encountered was 2.89 Å for the first carbonyl in a monodentate site 1 (our limit for a dissociated ligand was 3.04 Å).

Looking on the energies, it can be seen that the NMR energies are moderate, 98–188 kJ/mole (average 140 kJ/mole). This is slightly higher than for the structures obtained without any Ca restraints, but similar to all the

calculations with restraints (averages of the 20 best structures of 121–155 kJ/mole). Thus, the change in van der Waals parameters and the heavy initial MM minimisation do not significantly affect the energies. It is hard to discern any clear trends among the various coordination modes, except that all sites with a bidentate Glu ligand have relatively high NMR energies (averages 151–171 kJ/mole compared to 129–140 kJ/mole for the other types of sites). The five-coordinate structures and six-coordinate structure with a bidentate Asp give a slightly lower average NMR energy (129 and 131 kJ/mole) than the monodentate and bidentate structures with Asp (137 and 140 kJ/mole). The lowest NMR energies are obtained for two six-coordinate structures with a bidentate Asp ligand, whereas the mono- and bidentate Asp sites give the lowest QM and total energies. There is a clear correlation between the NMR and total energies, whereas there are hardly any correlation between the QM energy and the other energies (possibly a slight anticorrelation for the dissociated sites). Interestingly, for all coordination modes, the QM energies are  $\sim 20$  kJ/mole lower for site 1, whereas the NMR energies are  $\sim 10$  kJ/mole lower for site 2. This most likely reflect that there are more NMR restraints for site 2 than for site 1.

In conclusion, the COMQUM–N method works properly and give a general structure of the protein similar to the other methods, as can be seen in Figure 4 (note that there three times as many structures in these than the previous three figures). Both the COMQUM–N and MM results quite clearly show that a bidentate binding of the Glu ligand is energetically unfavourable and therefore unlikely. However, we cannot unambiguously decide if the  $\text{Ca}^{2+}$  sites are monodentate or bidentate with Asp. These two structures give similar energies (NMR, QM, and total) and they also arise in calculations started from other coordination modes. Perhaps, a slight higher tendency to bidentate Asp coordination can be seen for site 2. It is even likely that the two types of sites may show a fast interchange on an NMR time-scale, because the barrier between the mono- and bidentate binding of a carboxylate group is small (e.g.  $\sim 10$  kJ/mole for  $\text{Zn}^{2+}$  complexes (Ryde, 1999)).

#### 4. Concluding remarks

In this paper, we have tested and compared a number of methods to treat a metal site in NMR protein structure refinement. In particular, we have developed and tested two new methods to employ QM data in the refinement to supplement the refinement and obtain a more accurate description of a site of interest.

We have seen that it is not enough to describe the site as simple ligand–ligand restraint: This may lead to a structure in which the putative ligands do



not have the proper orientation to bind the metal (Figure 2). Instead, explicit metal–ligand bonds seem to be necessary to yield a realistic metal site.

QM data can be introduced in the refinement either directly, as in the COMQUM–N method or via the construction of an accurate MM potential. The latter method (Hess2FF) has been developed and tested out for hetero-compounds in crystal structures (Nilsson et al., 2003). However, it is equally suited and applicable for NMR refinement. The present results (Tables S6–S8) shows that it performs quite well also for the theoretically complicated plastic  $\text{Ca}^{2+}$  sites.

We have also developed the COMQUM–N method, as an NMR variant of crystallographic quantum refinement (Ryde et al., 2002; Ryde and Nilsson, 2003b). It is an appreciably more accurate method, because it avoids the risk of introducing errors during the conversion of QM data to the MM potential. Moreover, COMQUM–N allows changes in the coordination number during the refinement, reducing the risk of biasing the results by the choice of the MM potential. This is nicely illustrated in the application to EGF34, for which the unexpected possibility of a bidentate binding of the Asp ligand was discovered by the initial COMQUM–N calculations. On the other hand, this also means that there is a risk that COMQUM–N loses ligands by chance, e.g. if the starting structure is poor (as was seen in several of the COMQUM–N calculations on EGF34). In this sense, COMQUM–N is less robust than an MM potential, because there is only a minor attraction between the ion and its ligand at long distances. Therefore, COMQUM–N cannot be used to construct the metal site during early phases of the refinement. Instead, a more robust method has to be used to construct the starting structure for the final local refinement with COMQUM–N. On the other hand, COMQUM–N allows the site to disrupt if it is not supported by the NMR restraints.

Another advantage with the MM method is of course the speed. With only an MM potential, the NMR refinement is as fast as standard refinement, meaning that an ensemble of 200 structures can be constructed within a few hours on a standard PC. However, the QM optimisation of the metal site and the frequency calculation takes appreciably longer time, typically several days. Sometimes, the QM system may be so large (over  $\sim 50$  atoms) that it becomes hard to perform the frequency calculation. The COMQUM–N calculations, on the other hand, typically take one or two days each, meaning that a full ensemble of 200 structures would probably take a prohibitively long time. Therefore, the COMQUM–N calculations have to be restricted to the  $\sim 10$  best structures obtained by other methods. However, COMQUM–N involves only energy and force calculations and therefore avoid the frequency calculations, which have a much larger demand of memory and disk space. Therefore, COMQUM–N can be run on appreciably larger systems than the frequency calculation (up to  $\sim 200$  atoms).

Finally, it can be noted that additional experimental information can easily be included in the Hess2FF MM potential. For example, data from crystal structures (either proteins or small molecules) can be included. On the other hand, QM calculations using standard density functional methods typically give geometries of metal sites of an accuracy that is better than in a single protein crystal structure (Ryde and Nilsson, 2003a), so this is normally not advantageous except when accurate small-molecule crystallographic data are present for exactly the metal site of interest (the type of metal ligands strongly affects the bond lengths also of the other ligands to the metal) (Nilsson et al., 2003).

In conclusion, we suggest the following approach for the treatment of metal sites in NMR structure refinement: If the main interest is the general fold of the protein, a simple MM potential with only metal–ligand bonds is probably the best way to model the metal site. In this case, QM calculations are not necessary; instead the ideal bond length can be estimated from crystal structures of similar sites or from chemical intuition, and dummy force constants ( $\sim 100$  kJ/mole/Å<sup>2</sup>) can be employed. However, if a more detailed picture of the metal site is intended, a better method is needed. Hess2FF is recommended when a large number of different systems is to be tested, whereas COMQUM–N is the best method when accurate results is needed, e.g. at the end of an investigation with Hess2FF. The COMQUM–N calculations can be started from a refinement with a simple bonded potential.

The present QM calculations have been performed with density functional theory and medium-sized (6-31G\*) basis sets. We think this is an appropriate level of theory for the general use of our methods, even if it is quite costly (a few days of CPU time). However, for simpler systems (e.g. normal organic molecules), a lower level of theory could be used, e.g. the semiempirical PM3 method (Stewart, 1989), or even accurate MM methods, such as MMFF (Halgren, 1996), could be used (Nilsson et al., 2003). Such calculations can be performed within an hour for most systems of interest.

It is important to note that the presented methods, Hess2FF and COMQUM–N, are not restricted to metal sites. On the contrary, they are fully general and can be used to any site of interest. However, for the normal amino acids and nucleic acids, quite accurate target values for bond lengths and angles exist (Engl and Huber, 1991), reducing the need of more accurate methods. Yet, for unusual molecules (hetero-compounds), such as substrates and inhibitors, no MM potentials exist and for such sites, the present methods could be used. In particular, they could be useful for high-throughput NMR structure determination, where automatic methods are needed for hetero-compounds.

An important use of the present methods is to test conflicting structural hypotheses, such as whether a ligand binds in a mono- or bidentate mode in the present investigation. This is done by refining both candidates and

comparing their energies and structures. By similar methods, it has been possible to decide the protonation state of metal-bound solvent molecules by COMQUM-X (Ryde and Nilsson, 2003b; Nilsson and Ryde, 2004).

Another possible application of COMQUM-N is for structures of proteins that contain paramagnetic metal ions. Such metals lead to a significant broadening of the NMR signals around the metal site so that the local structure is hard to determine (Banci et al., 2002; Arnesano et al., 2003). By the use of COMQUM-N, accurate information about the missing local structure around the metal ion could be obtained. Finally, COMQUM-N can provide ideal starting structures for QM investigations of the structure, function, and reaction mechanism of proteins for which only NMR structures are available, providing an optimum compromise between experiments and quantum chemistry.

## 5. Acknowledgements

This investigation has been supported by grants from the Swedish research council (VR) and by computer resources of Lunarc at Lund University.

## References

- Ahlrichs, R., M. Bär, H.-P. Baron, R. Bauernschmitt, S. Böcker, M. Ehrig, K. Eichkorn, S. Elliott, F. Haase, M. Häser, H. Horn, C. Huber, C. Kölmel, M. Kollwitz, C. Ochsenfeld, H. Öhm, A. Schäfer, U. Schneider, O. Treutler, M. von Arnim, F. Weigend, P. Weis, and H. Weiss: 2000, 'TURBOMOLE Version 5.6'. Universität Karlsruhe, Germany.
- Appella, E., I. Weber, and F. Blasi: 1988, 'Structure and function of epidermal growth factor-like regions in proteins'. *FEBS Lett.* **231**, 1–4.
- Arnesano, F., L. Banci, I. Bertini, I. C. Felli, C. Luchinat, and A. R. Thompson: 2003, 'A Strategy for the NMR Characterization of Type II Copper(II) Proteins: the Case of the Copper Trafficking Protein CopC from *Pseudomonas Syringae*'. *J. Am. Chem. Soc.* **125**, 7200–7208.
- Banci, L., R. Pierattelli, and A. J. Villa: 2002, 'Nuclear magnetic resonance spectroscopy studies on copper proteins'. *Adv. Protein Chem.* **60**, 397–449.
- Becke, A.: 1988, 'Density-functional exchange-energy approximation with correct asymptotic behavior'. *Phys. Rev. A* **38**, 3098–3100.
- Brunger, A. T., P. D. Adams, G. M. Clore, W. L. Delano, P. Gros, R. W. Grosse-Kunstleve, J.-S. Jiang, J. Kuszewski, M. Niges, N. S. Pannu, R. J. Read, L. M. Rice, T. Simonson, and G. L. Warren: 2000, 'Crystallography & NMR System CNS, Version 1.1'. Yale University.
- Campbell, I. D. and P. Bork: 1993, 'Epidermal growth factor-like modules'. *Curr. Opin. Struct. Biol.* **3**, 385–392.
- Case, D. A., D. A. Pearlman, J. W. Caldwell, T. E. C. III, J. Wang, W. S. Ross, C. L. Simmerling, T. A. Darden, K. M. Merz, R. V. Stanton, A. L. Cheng, J. J. Vincent, M. Crowley, V. Tsui, H. Gohlke, R. J. Radmer, Y. Duan, J. Pitera, I. Massova, G. L. Seibel, U. C. Singh, P. K. Weiner, and P. A. Kollman: 2002, 'Amber Version 7'. University of California, San Francisco, USA.

- Cavanagh, J., W. J. Fairbrother, A. G. Palmer, and N. J. Skelton: 1996, *Protein NMR spectroscopy. Principles and practice*. Academic Press, London.
- Cruickshank, D. W. J.: 1999, 'Remarks about protein structure precision'. *Acta Crystallogr. D* **55**, 583–601.
- da Silva, J. J. R. F. and R. J. P. Williams: 1991, *In the biological chemistry of the elements*. Oxford: Clarendon.
- Doreleijers, J. F., M. L. Raves, T. Rullmana, and R. Kaptein: 1999, 'Completeness of NOEs in protein structures: A statistical analysis of NMR data'. *J. Biomol. NMR* **14**, 123–132.
- Downing, A. K., V. Knott, J. M. Werner, C. M. Cardy, I. D. Campbell, and P. A. Handford: 1996, 'Solution Structure of a Pair of Calcium-Binding Epidermal Growth Factor-like Domains: Implications for the Marfan Syndrome and Other Genetic Disorders'. *Cell* **85**, 597–605.
- Drakenberg, T., H. Ghasriani, A. Muranyi, A.-M. Thmlitz, and J. Stenfb: 2004. *Manuscript in preparation*.
- Eichkorn, K., O. Treutler, H. Öhm, M. Häser, and R. Ahlrichs: 1995, 'Auxiliary basis sets to approximate Coulomb potentials'. *Chem. Phys. Lett.* **240**, 283–290.
- Eichkorn, K., F. Weigend, O. Treutler, and R. Ahlrichs: 1997, 'Auxiliary basis sets for main row and transition metals and their use to approximate Coulomb potentials'. *Theor. Chim. Acta* **97**, 119–124.
- Engh, R. A. and R. Huber: 1991, 'Accurate bond and angle parameters for X-ray protein structure refinement'. *Acta Crystallogr.* **A47**, 392–400.
- Fields, B. A., H. H. Bartsch, H. D. Bartunik, F. Cordes, J. M. Guss, and H. C. Freeman: 1994, 'Accuracy and precision in protein crystal structure analysis: two independent refinements of the structure of poplar plastocyanin at 173 K'. *Acta Crystallogr.* **D50**, 709–730.
- Halgren, T. A.: 1996, 'Merck molecular force field. 1. Basis, form, scope, parameterization, and performance of MMFF94'. *J. Comput. Chem.* **17**, 490–641.
- Hehre, W. J., L. Radom, P. v. R. Schleyer, and J. A. Pople: 1986, *Ab initio molecular orbital theory*. Wiley-Interscience, New York.
- Henikoff, S., E. A. Greene, S. Pietrokovski, P. Bork, T. K. Attwood, and L. Hood: 1997, 'Gene familis: The taxonomy of protein paralogs and chimeras'. *Science* **278**, 609–614.
- Hertwig, R. H. and W. Koch: 1997, 'On the parameterization of the local correlation functional. What is Becke-3-LYP?'. *Chem. Phys. Lett.* **268**, 345–351.
- Jalilvand, F., D. Spånberg, P. Lindqvist-Reis, K. Hermansson, I. Persson, and M. Sandström: 2001, 'Hydration of the calcium ion. An EXAFS, large-angle X-ray scattering, and molecular dynamics simulation study'. *J. Am. Chem. Soc.* **123**, 431–441.
- Jensen, F.: 1999, *Introduction to computational chemistry*. John Wiley & Sons, Chichester.
- Klamt, A., V. Jonas, T. Bürger, and J. C. W. Lohrenz: 1998, 'Refinement and parametrization of COSMO-RS'. *J. Phys. Chem. A* **102**, 5074–5085.
- Klamt, A. and J. Schürmann: 1993, 'COSMO'. *J. Chem. Soc., Perkin Transact. II* **5**, 799–805.
- Kleywegt, G. J. and T. A. Jones: 1995, 'Where freedom is given, liberties are taken'. *Structure* **3**, 535–540.
- Kleywegt, G. J. and T. A. Jones: 1998, 'Databases in protein crystallography'. *Acta Crystallogr.* **D54**, 1119–1131.
- Maseras, F. and K. Morokuma: 1995, 'IMOMM – a new integrated ab-initio plus molecular mechanics geometry optimization scheme of equilibrium structures and transition-states'. *J. Comput. Chem.* **16**, 1170–1179.
- Monard, G. and K. M. Merz: 1999, 'Combined quantum mechanical/molecular mechanical methodologies applied to biomolecular systems'. *Acc. Chem. Res.* **32**, 904–911.
- Mulholland, A. J.: 2001, 'The QM/MM approach to enzymatic reactions'. In: L. A. Eriksson (ed.): *Theoretical biochemistry – Processes and properties of biological systems*

- (*Theoretical and computational chemistry*, Vol. 9). Amsterdam: Elsevier Science, pp. 597–505.
- Nicoll, R. M., S. A. Hindle, G. MacKenzie, I. H. Hillier, and N. A. Burton: 2001, 'Quantum mechanical/molecular mechanical methods and the study of kinetic isotope effects: modelling the covalent junction region and application to the enzyme xylose isomerase'. *Theor. Chim. Acta* **106**, 105–112.
- Nilsson, K., D. Lecerof, E. Sigfridsson, and U. Ryde: 2003, 'An automatic method to generate force-field parameters for hetero-compounds'. *Acta Crystallogr.* **D59**, 274–289.
- Nilsson, K. and U. Ryde: 2004, 'Protonation status of metal-bound ligands can be determined by quantum refinement'. *J. Inorg. Biochem.* **98**, 1539–1546.
- Norrby, P.-O. and T. Liljefors: 1998, 'Automated molecular mechanics parameterization with simultaneous utilization of experimental and quantum mechanical data'. *J. Comput. Chem.* **19**, 1146–1166.
- Olsson, M. H. M. and U. Ryde: 2001, 'Geometry, reduction potential, and reorganisation energy of the binuclear Cu<sub>A</sub> site, studied by theoretical methods'. *J. Am. Chem. Soc.* **123**, 7866–7876.
- Perdew, J. P.: 1986, 'Density-functional approximation for the correlation energy of the inhomogeneous electron gas'. *Phys. Rev. B* **33**, 8822–8824.
- Rao, Z., P. Handford, M. Mayhew, V. Knott, G. G. Brownlee, and D. Stuart: 1995, 'The structure of a Ca<sup>2+</sup>-binding epidermal growthfactor-like domain: its role in protein–protein interactions'. *Cell* **82**, 131–141.
- Ryde, U.: 1996, 'The coordination of the catalytic zinc ion in alcohol dehydrogenase studied by combined quantum chemical and molecular mechanical calculations'. *J. Comput.-Aided Mol. Design* **10**, 153–164.
- Ryde, U.: 1999, 'Carboxylate binding modes in zinc proteins, a theoretical study'. *Biophys. J.* **77**, 2777–2787.
- Ryde, U.: 2003, 'Combined quantum and molecular mechanics calculations on metalloproteins'. *Curr. Opin. Chem. Biol.* **7**, 136–142.
- Ryde, U. and K. Nilsson: 2003a, 'Quantum chemistry can improve protein crystal structures locally'. *J. Am. Chem. Soc.* **125**, 14232–14233.
- Ryde, U. and K. Nilsson: 2003b, 'Quantum refinement – a combination of quantum chemistry and protein crystallography'. *J. Mol. Struct. (Theochem)* **632**, 259–275.
- Ryde, U., L. Olsen, and K. Nilsson: 2002, 'Quantum chemical geometry optimisations in proteins using crystallographic raw data'. *J. Comput. Chem.* **23**, 1058–1070.
- Ryde, U. and M. H. M. Olsson: 2001, 'Structure, strain, and reorganisation energy of blue-copper models in the protein'. *Int. J. Quantum Chem.* **81**, 335–347.
- Ryde, U., M. H. M. Olsson, B. O. Roos, and A. B. Carlos: 2001, 'The dependence of the geometry of blue copper protein models on the method, basis sets, and model size'. *Theor. Chim. Acta* **105**, 452–462.
- Saha, S., J. Boyd, J. M. Werner, V. Knott, P. A. Handford, I. D. Campbell, and A. K. Downing: 2001, 'Solution Structure of the LDL Receptor EGF-AB Pair: A Paradigm for the Assembly of Tandem Calcium Binding EGF Domains'. *Structure* **9**, 451–456.
- Schäfer, A., A. Klant, D. Sattel, J. C. W. Lohrenz, and F. Eckert: 2000, 'COSMO Implementation in TURBOMOLE: Extension of an efficient quantum chemical code towards liquid systems'. *PhysChemChemPhys* **2**, 2187–2193.
- Siegbahn, P. E. M. and M. R. A. Blomberg: 2000, 'Transition-Metal Systems in Biochemistry Studied by High-Accuracy Quantum Chemical Methods'. *Chem. Rev.* **100**, 421–437.
- Sigfridsson, E., M. H. M. Olsson, and U. Ryde: 2001, 'A comparison of the inner-sphere reorganisation energies of cytochromes, iron–sulphur clusters, and blue copper proteins'. *J. Phys. Chem. B* **105**, 5546–5552.

- Stenberg, Y., S. Linse, T. Drakenberg, and J. Stenfb: 1997a, 'The high affinity calcium-binding sites in the epidermal growth factor module region of vitamin K-dependent protein S'. *J. Biol. Chem.* **272**, 23255–23260.
- Stenberg, Y., A. Muranyi, C. Steen, E. Thulin, T. Drakenberg, and J. Stenfb: 1997b, 'EGF-like module pair 3–4 in vitamin K-dependent protein S. Modulation of calcium affinity of module 4 by module 3 and interaction with factor X'. *J. Mol. Biol.* **293**, 653–665.
- Stenfb, J., Y. Stenberg, and A. Muranyi: 2000, 'Calcium-binding EGF-like modules in coagulation proteinases: function of the calcium ion in module interactions'. *Biochim. Biophys. Acta* **1477**, 51–63.
- Stewart, J. J. P.: 1989, 'Optimization of parameters for semiempirical methods. I. Method'. *jcc* **10**, 209–220.
- Svensson, M., S. Humbel, R. D. J. Froese, T. Matsubara, S. Sieber, and K. Morokuma: 1996, 'ONIOM: A Multilayered Integrated MO + MM Method for Geometry Optimizations and Single Point Energy Predictions'. *J. Phys. Chem.* **100**, 19357.
- Torrent, M., D. G. Musaev, and K. Morokuma: 2001, 'The flexibility of carboxylate ligands in methane monooxygenase and ribonucleotide reductase: A density functional study'. *J. Phys. Chem. B* **105**, 322–327.
- Tossavainen, H., P. Permi, A. Annala, I. Kilpeläinen, and T. Drakenberg: 2003, 'NMR solution structure of calerythrin, and EF-hand calcium-binding protein from *Saccaropolyspora erythraea*'. *Eur. J. Biochem.* **270**, 2505–2512.
- Wang, X., M. X. Li, L. Spyropoulos, N. Beier, M. Chandra, R. J. Solaro, and B. D. Sykes: 2001, 'Structure of the C-domain of Human Cardiac Troponin C in Complex with the Ca<sup>2+</sup> Sensitizing Drug EMD 57033'. *J. Biol. Chem.* **276**, 25456–25466.

**Scheme 1.** The program flow in the COMQUM–N program. Tasks performed by the QM program are shown in **bold face**, those performed by the CNS software are shown in *italics*, and those performed by the COMQUM–N interface routines are underlined. S1 and S2 denotes systems 1 and 2.

**Evaluate QM wavefunction of S1 including electrostatics of S2**

Repeat

**Evaluate QM forces from S1 + electrostatics of S2 onto S1**

*Evaluate CNS forces (from S1 and S2 onto S1), no electrostatics*

Add the QM and CNS forces

**Relax the geometry of S1 using these forces<sup>a</sup>**

Change coordinates of S1 in CNS representation

**Calculate charges of S1**

Insert these charges into CNS representation

*Relax S2 by an NMR-restrained minimisation with S1 fixed*

Change the coordinates of S2 in the QM representation (point charges)

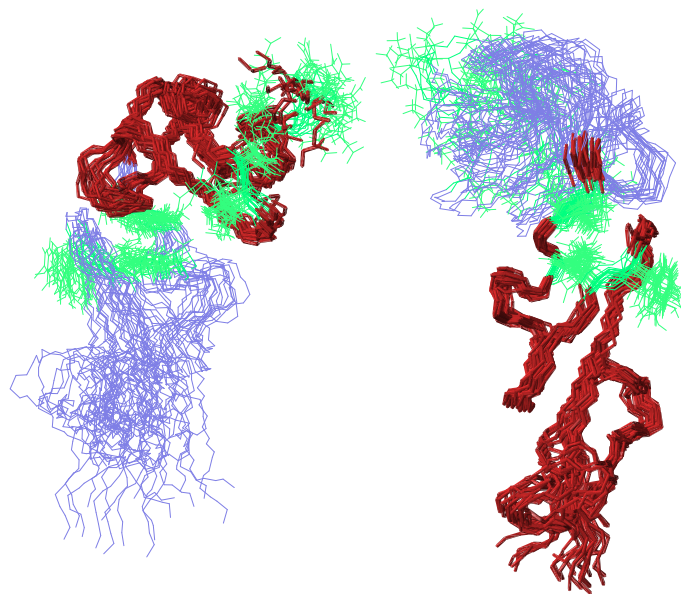
**Evaluate QM wavefunction and energy of S1 including S2 electrostatics**

*Evaluate CNS energy function*

Add energies

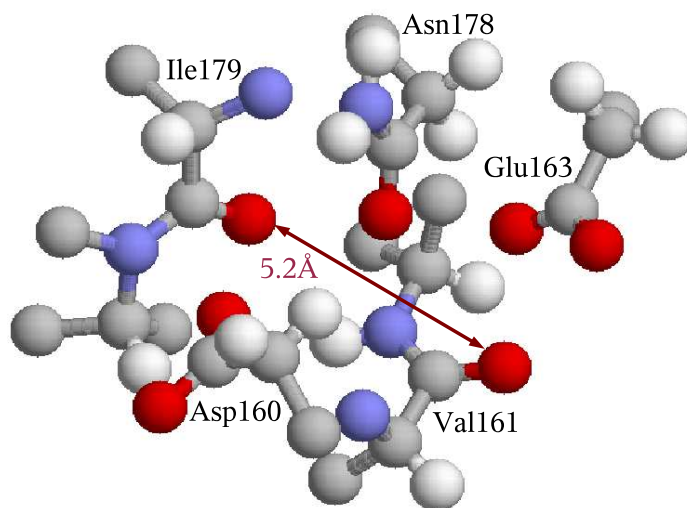
Until convergence

<sup>a</sup> The geometry optimisation of S1 can be performed by any program, but for convenience, we have used the QM program.

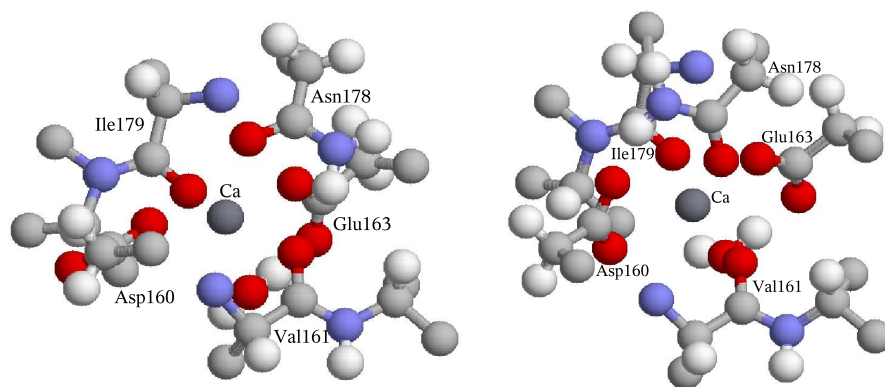


*Figure 1.* General structure of the EGF34 fragment with the ligands of the two  $\text{Ca}^{2+}$  sites emphasized in green. The best 20 structures are used, employing the data in Table S2 (no  $\text{Ca}^{2+}$  restraints). The left-hand side image was obtained by superimposing domain 3 (with Ca site 1), whereas the right-hand side image was obtained by superimposing domain 4 (with Ca site 2). The two domains are connected by a flexible hinge.

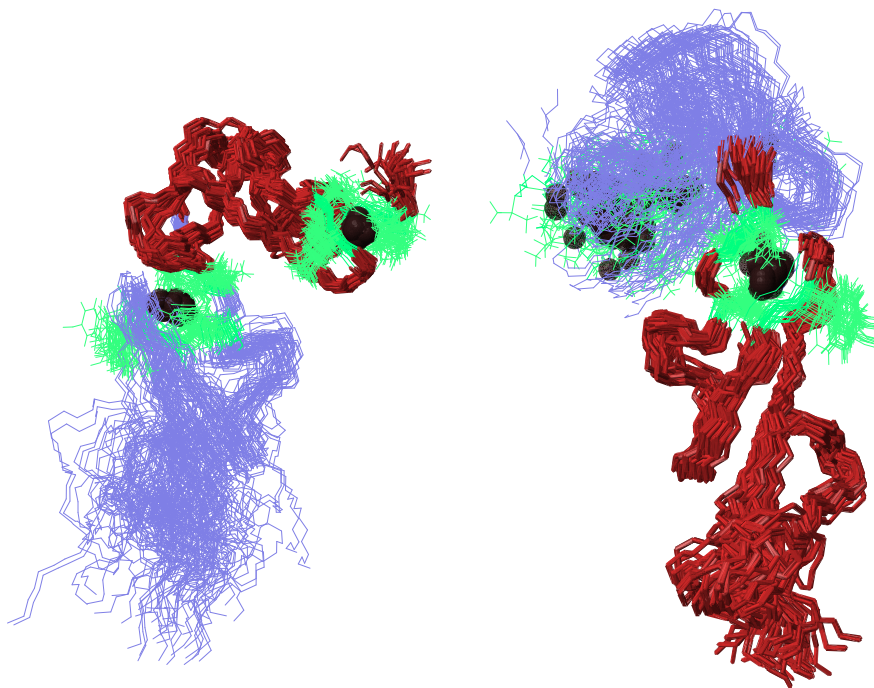




*Figure 2.* The structure of  $\text{Ca}^{2+}$  site 1, obtained with O–O restraints, showing that the orientation of some of the ligands are not proper for Ca binding.



*Figure 3.* The final COMQUM–N structures of  $\text{Ca}^{2+}$  site 1 with a monodentate (left) or bidentate with Asp (right) binding. The best structures (in terms of total energy) from Tables S12 and S13 were used.



*Figure 4.* General structure of the EGF34 fragment with the ligands of the two Ca sites emphasized in green and with the  $\text{Ca}^{2+}$  ions in black. All 60 COMQUM-N structures are superimposed, using the data in Tables S12 and S13. Note that this is three times as many structures as in Figure 1 (and Figures S1 and S2). The left- and right-hand side images were obtained by superimposing domain 3 and 4, respectively.

## **6. Supplementary material**

Table I. The result of quantum chemical geometry optimisations of the  $\text{Ca}(\text{CH}_3\text{COO})_2(\text{CH}_3\text{CONHCH}_3)_2(\text{CH}_3\text{CONH}_2)(\text{H}_2\text{O})$  model, using various methods and starting either from the crystal structure of the EGF-like domain in human clotting factor IX (Rao et al., 1995) or from preliminary NMR structures of EGF34 (site 1 or 2). The various Ca–O distances are listed as well as the relative energy of models obtained with the different starting structures ( $\Delta E$ ). The geometry of the two  $\text{Ca}^{2+}$  sites in the crystal structure are also given for comparison.

Method	Start	$\Delta E$	Ca–O distance (Å)							
			Asp	CO1	Glu1	Glu2	Asn	CO2	Wat	Av.
BP86/6-31G*	Crystal	0.0	2.39	2.40	2.35	3.88	2.39	2.39	2.44	2.40
	NMR1	-13.4	2.38	2.43	2.41	3.53	2.40	2.43	2.42	2.41
	NMR2	44.8	2.40	2.38	2.34	3.90	2.42	2.46	2.43	2.41
COSMO	Crystal		2.41	2.37	2.35	3.87	2.39	2.38	2.45	2.39
Bidentate	NMR1	18.6	2.37	2.49	2.51	2.55	2.58	2.48	2.46	2.49
Extra water	NMR2		2.36	2.38	2.54	2.54	2.55	2.52	2.49	2.48
BP86/6-311+G(2d,2p)	Crystal	0.0	2.35	2.38	2.32	3.84	2.36	2.37	2.42	2.37
	NMR1	-11.5	2.33	2.42	2.36	3.56	2.39	2.41	2.38	2.38
	NMR2	38.6	2.35	2.35	2.29	3.85	2.41	2.44	2.42	2.38
B3LYP/6-31G*	Crystal		2.38	2.40	2.34	3.88	2.39	2.39	2.45	2.39
Crystal structure, site 1			2.26	2.37	2.57	2.60	2.36	2.21		2.40
Crystal structure, site 2			2.30	2.27	2.40	2.60	2.41	2.29	2.36	2.38

Table II. The result of a refinement of EGF34 using only NMR restraints and no information at all about the  $\text{Ca}^{2+}$  sites. The table shows the total (Tot) and NMR energy, as well as the individual NOE, dihedral (Dih), and SANI energy terms. In addition, the two  $\text{Ca}^{2+}$  sites are described by giving the distance between the two carbonyl oxygen atoms involved in the site (CO), the maximum O–O distance of the ligand atoms, and the maximum Ca–O distance, assuming that the  $\text{Ca}^{2+}$  ion resides at the midpoint between the two carbonyl atoms. Data is given for the best five structures in terms of total energy, the five best structures for each  $\text{Ca}^{2+}$  site in terms of the maximum O–O distance, as well as the average for all 200 and the 20 best structures (in terms of total energy).

Struct	Energy terms kJ/mole					Site 1 dist. (Å)			Site 2 dist (Å)		
	Tot	NMR	NOE	Dih	ANI	CO	O–O	Ca–O	CO	O–O	Ca–O
1	529	75	68	3	4	12.17	15.33	10.87	6.76	10.40	6.79
2	576	80	74	2	3	10.90	14.36	10.41	6.85	10.20	7.18
3	619	88	80	1	7	10.34	15.32	10.75	6.70	9.00	6.36
4	619	82	73	5	4	11.82	13.19	9.51	6.32	9.73	7.06
5	617	87	75	2	10	12.64	13.79	9.15	6.81	8.56	6.33
av 20	654	92	80	6	6	11.88	14.76	9.96	6.56	9.34	6.32
av 200	2281	582	479	78	26	11.47	15.55	10.28	6.20	9.39	6.49
136	2755	1079	1064	8	7	7.90	8.78	6.41			
179	4597	1210	1015	138	57	10.21	11.08	7.28			
124	2255	642	525	90	27	11.41	11.41	7.03			
67	1078	172	103	50	19	11.58	11.58	6.86			
85	1257	241	207	13	21	11.82	11.82	7.22			
170	4031	1136	970	121	44				5.56	5.56	3.83
179	4597	1210	1015	138	57				6.07	6.07	3.14
189	5702	871	738	96	37				5.34	6.35	4.25
89	1303	312	193	98	22				6.43	6.43	4.44
185	5141	734	654	63	18				6.54	6.54	4.80

Table III. The result of a refinement of EGF34 using O–O restraints for the  $\text{Ca}^{2+}$  site. The table shows the same entries as Table S2, but the best sites are sorted after the maximum Ca–O distance.

Struct	Energy terms kJ/mole					Site 1 dist. (Å)			Site 2 dist (Å)		
	Tot	NMR	NOE	Dih	ANI	CO	O–O	Ca–O	CO	O–O	Ca–O
1	762	111	99	8	5	5.17	5.49	4.30	5.17	5.38	4.42
2	796	122	106	10	6	5.20	5.44	4.15	5.22	5.35	4.17
3	806	112	90	14	9	5.12	5.41	3.78	5.20	5.23	3.98
4	807	122	111	6	6	5.14	5.60	4.29	5.23	5.43	3.97
5	825	121	100	13	8	5.13	5.48	3.25	5.18	5.42	3.94
av 20	859	134	116	11	8	5.17	5.50	3.89	5.21	5.39	3.80
av 200	2895	744	632	84	28	5.09	5.60	4.10	5.19	5.46	3.82
11	884	143	124	10	9	5.16	5.50	3.04			
180	6056	1976	1658	182	36	5.27	5.33	3.14			
43	1018	186	129	43	15	5.16	5.38	3.19			
18	907	152	128	13	10	5.13	5.49	3.25			
6	825	121	100	13	8	5.13	5.48	3.25			
14	890	154	126	17	12				5.19	5.19	2.79
41	994	160	128	18	15				5.24	5.44	2.84
89	1576	336	247	50	39				5.18	5.38	2.84
194	8271	2791	2723	46	22				5.22	5.42	2.86
18	907	152	128	13	10				5.23	5.23	2.88

Table IV. The result of a refinement of EGF34 using flat-bottomed Ca–O NMR restraints. The table shows the same entries as Table S2, but the best sites are sorted after the maximum Ca–O distance.

Struct	Energy terms kJ/mole					Site 1 dist. (Å)			Site 2 dist (Å)		
	Tot	NMR	NOE	Dih	ANI	CO	O–O	Ca–O	CO	O–O	Ca–O
1	650	93	84	4	5	6.00	6.00	3.07	6.10	6.10	3.07
2	656	98	90	4	4	6.05	6.05	3.04	6.07	6.07	3.05
3	735	108	94	9	6	5.88	5.88	3.05	5.38	5.68	3.05
4	741	121	107	8	7	6.05	6.05	3.21	6.09	6.09	3.05
5	749	93	80	7	6	5.98	5.98	3.04	6.09	6.09	3.06
av 20	827	131	113	10	8	5.91	5.95	3.09	5.97	5.97	3.06
av 200	3504	1008	833	125	50	5.88	6.03	3.12	5.74	5.92	3.08
182	7608	1915	1670	195	49	5.53	5.53	2.98			
198	11132	3428	2758	501	170	5.67	5.67	3.00			
99	2531	687	473	157	57	5.73	5.73	3.01			
114	3092	766	618	117	31	5.44	5.44	3.02			
194	9345	3288	2806	212	271	5.91	5.91	3.02			
145	4760	977	8006	93	84				5.04	5.05	2.93
73	1624	328	988	70	60				5.53	5.53	2.99
186	8212	3604	3318	225	60				5.98	5.98	2.99
87	2033	397	288	79	31				5.38	5.38	3.00
69	1585	330	217	60	54				5.04	5.34	3.00



Table V. Equilibrium bond lengths ( $r_0$  in Å) and force constants ( $k$  in kJ/mole/Å<sup>2</sup>) used in the refinements with an MM potential for Ca<sup>2+</sup>. They were obtained from the first and sixth structures in Table S1. Values of similar interactions (i.e. water, carboxylate, and carbonyl groups) were averaged. For the bidentate site, different parameters were used for the mono- and bidentate carboxylate groups.

Interaction	Monodentate		Bidentate	
	$r_0$	$k$	$r_0$	$k$
Carboxylate (monodentate)	2.371	146.10	2.364	144.38
Carboxylate (bidentate)			2.540	85.86
Carbonyl	2.394	115.60	2.483	72.71
Water	2.441	129.28	2.490	98.29

Table VI. The result of a refinement of EGF34 using Ca–O MM bonds for a monodentate site. The table shows the same entries as Table S2, but the best sites are sorted after the maximum Ca–O distance.

Struct	Energy terms kJ/mole						Site 1 dist. (Å)			Site 2 dist (Å)		
	Tot	NMR	NOE	Dih	ANI		CO	O–O	Ca–O	CO	O–O	Ca–O
1	761	121	103	12	7	4.93	4.93	2.55	4.91	4.91	2.57	
2	770	110	100	5	5	5.03	5.03	2.56	5.45	5.45	2.75	
3	791	93	83	6	5	4.75	4.75	2.52	5.44	5.44	2.79	
4	801	112	94	10	8	5.09	5.09	2.57	5.40	5.40	2.71	
5	850	139	115	16	8	4.98	4.98	2.52	4.94	4.94	2.53	
av 20	946	144	117	17	10	5.05	5.07	2.60	5.25	5.25	2.68	
av 200	2511	633	504	90	39	5.16	5.21	2.68	5.12	5.14	2.65	
48	1272	249	120	60	69	4.80	4.80	2.45				
111	1701	446	318	81	47	4.89	4.89	2.48				
76	1440	348	188	121	39	4.90	4.90	2.49				
41	1226	224	180	19	25	4.94	4.94	2.49				
37	1189	256	185	36	34	4.81	4.81	2.50				
157	3332	671	513	137	31				4.83	4.83	2.43	
62	1393	343	281	16	46				4.74	4.77	2.48	
99	1605	322	170	112	40				4.84	4.84	2.48	
97	1579	377	254	91	31				4.54	4.77	2.48	
108	1658	501	311	136	54				4.94	4.94	2.48	

Table VII. The result of a refinement of EGF34 using Ca–O MM bonds for sites with the Glu residue bidentate. The table shows the same entries as Table S2, but the best sites are sorted after the maximum Ca–O distance.

Struct	Energy terms kJ/mole					Site 1 dist. (Å)			Site 2 dist (Å)		
	Tot	NMR	NOE	Dih	ANI	CO	O–O	Ca–O	CO	O–O	Ca–O
1	866	169	146	17	6	5.42	5.42	2.81	5.74	5.42	2.81
2	877	131	106	7	17	5.59	5.59	2.84	5.92	5.59	2.84
3	876	145	133	7	4	5.59	5.59	2.82	5.75	5.59	2.82
4	910	131	118	7	6	5.42	5.42	2.74	5.72	5.42	2.74
5	940	159	136	14	9	5.63	5.63	2.85	5.95	5.63	2.85
av 20	986	155	129	15	11	5.44	5.44	2.81	5.74	5.74	2.91
av 200	3394	882	735	108	40	5.53	5.69	3.05	5.44	5.47	2.82
177	6543	1806	1650	119	37	4.53	4.77	2.57			
24	1159	252	170	54	28	5.09	5.09	2.62			
67	1647	416	189	176	50	5.19	5.19	2.62			
50	1450	324	236	76	12	5.15	5.15	2.63			
171	5869	1294	1070	159	66	5.25	5.25	2.63			
145	4778	1390	1044	245	101				4.13	4.78	2.49
189	8192	2242	1904	158	180				4.85	4.85	2.51
109	2823	724	576	82	65				5.02	5.02	2.57
40	1354	297	192	68	37				4.99	4.99	2.61
90	1969	392	284	61	47				5.07	5.09	2.61

Table VIII. The result of a refinement of EGF34 using Ca–O MM bonds for sites with the Asp residue bidentate. The table shows the same entries as Table S2, but the best sites are sorted after the maximum Ca–O distance.

Struct	Energy terms kJ/mole					Site 1 dist. (Å)			Site 2 dist (Å)		
	Tot	NMR	NOE	Dih	ANI	CO	O–O	Ca–O	CO	O–O	Ca–O
1	727	94	83	6	6	5.24	5.24	2.69	5.76	5.76	2.92
2	747	87	76	6	5	5.52	5.52	2.79	5.88	5.88	3.00
3	794	116	100	8	8	5.40	5.40	2.74	5.80	5.80	2.94
4	793	112	97	9	5	5.69	5.69	2.88	5.65	5.65	2.85
5	845	117	94	14	8	5.40	5.40	2.83	5.78	5.78	2.96
av 20	880	121	102	10	8	5.36	5.36	2.78	5.64	5.64	2.90
av 200	3002	772	635	106	31	5.61	5.68	3.01	5.43	5.45	2.84
121	2888	1876	944	917	16	5.53	5.03	2.56			
185	6732	4088	2071	1790	227	4.09	4.89	2.61			
26	1036	347	180	149	19	5.19	5.11	2.65			
70	1397	555	292	186	77	5.15	5.05	2.65			
37	1148	431	224	118	89	5.25	5.20	2.65			
147	4261	1387	708	562	117				4.88	4.93	2.54
63	1352	522	282	158	83				4.95	4.95	2.56
129	3261	1632	872	468	292				5.02	5.07	2.57
142	3855	1677	869	700	109				5.12	5.12	2.61
101	2018	730	382	336	12				5.04	5.04	2.61

Table IX. The number MM minimisation steps in the COMQUM–N calculations. We list the number of geometry optimisation iterations (It.), the Ca–O distances in site 1, as well as the total and NMR energy and the individual NMR energy terms). The total energy is relative to the one obtained with 100 000 steps (–1917.7102 H).

Step	It.	Ca–O distance (Å)							Energy(kJ/mole)				
		Asp	Val	Glu	Asn	Ile	Wat	Av.	Tot	NMR	NOE	Dih	ANI
1	212	2.23	2.40	2.35	2.48	2.39	2.68	2.42	116	300	242	28	31
50	61	2.35	2.44	2.35	2.46	2.45	2.53	2.43	27	300	241	28	31
200	69	2.38	2.39	2.36	2.47	2.41	2.50	2.42	8	299	241	27	31
1 000	65	2.39	2.37	2.37	2.48	2.40	2.51	2.42	4	302	243	28	31
2 000	61	2.39	2.37	2.37	2.48	2.40	2.52	2.42	4	302	243	27	31
5 000	64	2.39	2.36	2.37	2.48	2.40	2.51	2.42	3	301	242	27	31
10 000	54	2.38	2.37	2.38	2.46	2.41	2.50	2.42	7	300	242	27	31
20 000	64	2.39	2.37	2.37	2.48	2.39	2.50	2.42	1	301	242	28	31
30 000	64	2.39	2.37	2.37	2.48	2.40	2.51	2.42	2	301	242	28	31
40 000	67	2.39	2.37	2.36	2.48	2.39	2.51	2.42	0	301	242	28	31
100 000	67	2.39	2.37	2.36	2.48	2.39	2.51	2.42	0	301	242	28	31

Table X. The influence of the NOE weight on the COMQUM-N results. Ca-O distances in site 1, as well as the total and NMR energy and the individual NMR energy terms) are listed. The total energy is relative to that obtained with the NOE weight 75 (-1917.59 H).

NOE	Ca-O distance (Å)							Energy(kJ/mole)				
	Asp	Val	Glu	Asn	Ile	Wat	Av.	Tot	NMR	NOE	Dih	ANI
75	2.39	2.37	2.38	2.47	2.41	2.49	2.42	0	301	241	27	33
125	2.39	2.38	2.37	2.48	2.42	2.50	2.42	131	277	211	30	37
300	2.39	2.39	2.37	2.46	2.43	2.49	2.42	420	214	138	32	44
750	2.35	2.42	2.34	2.45	2.45	2.47	2.42	446	155	74	33	48

Table XI. Test of the treatment of electrostatics in the COMQUM–N method. Three options were tested: electrostatics included in the QM system (a point-charge model of the protein; QM), electrostatics included also in the NMR-minimisation of the protein (Prot), and a COSMO continuum solvation model in the QM calculations. Ca–O distances in site 1 as well as the NMR energy terms (total as well as the individual terms) are listed, based on two different starting structures.

Elstat. in			Ca–O distance (Å)							Energy(kJ/mole)			
QM	Prot	COSMO	Asp	Val	Glu	Asn	Ile	Wat	Av.	NMR	NOE	Dih	ANI
Starting structure 1													
no	no	yes	2.35	2.40	2.60	2.47	2.52	2.32	2.44	324	261	29	35
no	no	no	2.39	2.40	2.34	2.52	2.46	2.47	2.43	325	261	29	35
yes	no	no	2.35	2.58	2.27	2.78	2.39	2.60	2.50	332	268	29	35
yes	yes	no	2.33	2.60	2.26	2.83	2.41	2.63	2.51	488	382	48	58
Starting structure 2													
no	no	yes	2.35	2.42	2.38	2.55	2.48	2.56	2.46	312	249	30	33
no	no	no	2.36	2.49	2.32	2.58	2.42	2.52	2.45	317	253	31	34
yes	no	no	2.36	2.68	2.40,2.57	4.77	2.33	2.36	2.45	321	257	31	33
yes	yes	no	5.02	4.36	2.47,2.43	2.37	2.32	2.48	2.41	501	382	46	72

Table XII. Result of the COMQUM–N calculations, starting from a monodentate site. In addition to the energy terms described in the other tables, we here also report the QM energy of the quantum system (relative to the structure with the lowest energy, –1917.59285 H and the total QM/MM energy (again relative to the structure with the lowest energy, –1917.96853 H).

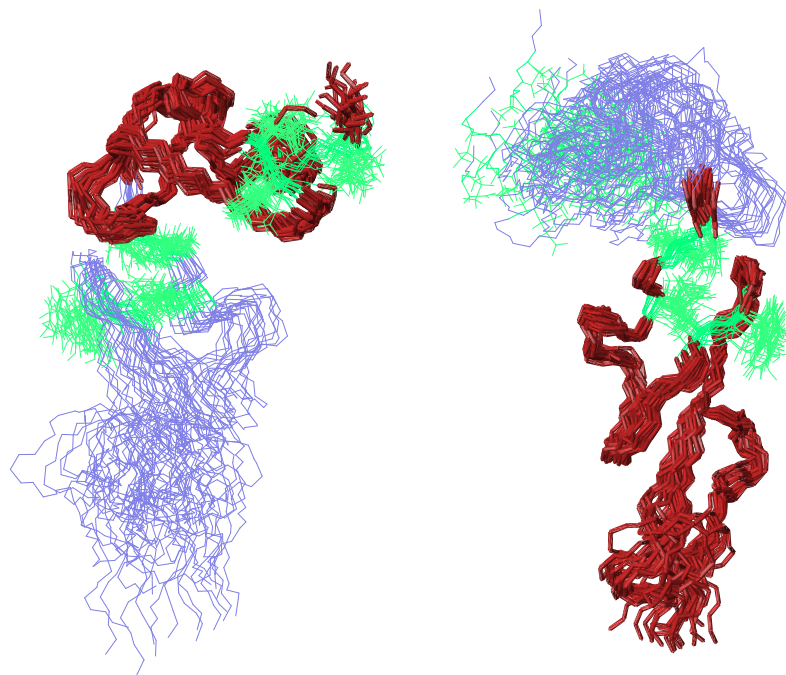
Asp	Ca–O distance (Å)						Energy(kJ/mole)					
	CO(1) Glu		Asn	CO(2) Wat		Av.	NMR	NOE	Dih	ANI	QM	Tot
Monodentate, Site 1												
2.31	2.56	2.31	2.42	2.49	2.42	2.42	127	100	10	16	31	159
2.45	2.45	2.31	2.31	2.58	2.48	2.43	128	108	5	14	55	225
2.33	4.46	2.23	2.36	2.36	2.50	2.71	108	83	8	17	194	300
2.27	2.64	2.37	2.39	2.52	2.53	2.45	141	111	9	21	29	24
2.52,2.61	2.62	2.36	2.39	2.25	2.53	2.51	177	144	16	16	71	332
2.40	2.53	2.40	2.37	2.49	2.52	2.46	135	112	7	15	31	215
2.29	2.60	2.38	2.34	2.48	2.57	2.44	175	117	31	27	56	342
2.27	2.53	2.29	2.42	2.60	2.45	2.43	140	114	7	19	25	413
2.32	2.34	2.26	2.36	3.67	2.45	2.57	170	122	18	30	76	387
2.38	2.60	2.56,2.46	2.36	2.69	2.69	2.53	150	104	27	19	44	305
Monodentate, Site 2												
2.64,2.44	2.68	2.30	2.43	2.43	2.51	2.49	139	119	7	13	46	145
2.29	2.59	2.26	2.28	3.04	2.41	2.48	117	95	7	14	92	221
2.34	2.55	2.32	2.38	2.51	2.44	2.42	103	83	10	13	97	84
2.43,2.41	2.64	2.42	2.42	2.48	2.44	2.46	127	97	11	20	45	0
2.34	2.41	2.53	2.31	2.45	2.45	2.41	139	127	23	19	70	358
2.46,2.44	3.50	2.29	2.33	2.51	2.52	3.58	117	98	6	12	91	175
2.40,2.50	2.44	2.39	2.41	3.54	2.44	2.59	172	117	34	22	79	335
2.29	3.52	2.35	2.28	2.43	2.54	2.57	120	103	4	13	111	432
2.32	3.38	2.33	2.35	2.46	2.42	2.55	143	102	21	20	97	295
2.37,2.52	2.52	2.50	2.32	2.50	2.47	2.46	158	103	34	21	30	317



Table XIII. Result of the COMQUM–N calculations, starting from a bidentate site. The items are the same as in Table S12.

Asp	Ca–O distance (Å)					Energy(kJ/mole)					
	CO(1)	Glu	Asn	CO(2)	Wat Av.	NMR	NOE	Dih	ANI	QM	Tot
Glu bidentate, Site 1											
2.27	2.41	2.44,2.36	4.51	2.45	2.45	2.70	188	156	18	14	106 359
2.27	2.43	2.35,2.55	2.49	4.86	2.52	2.50	127	107	8	12	117 304
2.34	4.86	2.30,4.50	2.31	2.62	2.43	2.77	157	140	6	12	50 349
2.27	2.47	4.61,2.38	2.35	2.44	2.53	2.72	159	137	10	12	50 371
2.30	2.55	2.38,2.48	2.45	2.49	2.58	2.46	165	135	15	15	70 499
2.26	2.37	2.34,2.54	2.45	3.79	2.44	2.60	146	129	4	12	91 355
2.32	2.33	2.35,2.50	2.45	3.46	2.41	2.55	176	136	20	19	41 447
2.27	2.45	2.53,2.52	2.40	3.10	2.42	2.53	146	93	16	37	66 335
2.27	2.35	4.42,2.32	2.38	3.19	2.44	2.77	174	137	10	27	94 410
2.27	2.55	3.64,2.32	2.40	2.51	2.43	2.59	163	122	18	23	35 251
Glu bidentate, Site 2											
2.37,2.49	2.72	2.46,2.53	2.45	2.50	4.48	2.75	185	153	21	11	52 286
2.28	2.53	2.38,3.81	2.37	2.58	2.44	2.63	117	98	7	11	44 136
2.28	3.42	2.44,2.69	2.36	2.39	2.48	2.58	141	122	8	11	119 363
2.30	2.47	2.80,2.39	2.43	2.69	2.67	2.53	146	123	12	11	106 359
2.35,2.39	3.43	2.35,3.85	2.43	2.35	2.46	2.70	156	128	15	14	84 443
2.27	3.29	2.57,2.63	2.45	2.38	2.47	2.58	128	112	5	11	104 283
2.27	2.62	2.42,2.46	2.35	2.44	4.58	2.73	159	119	22	18	95 440
2.33	2.58	2.39,3.04	2.33	3.03	2.73	2.63	114	94	7	13	144 301
2.26	2.69	2.56,2.62	2.43	2.61	2.51	2.53	168	131	10	28	112 396
2.37,2.44	2.69	2.51,2.45	2.47	2.49	4.34	2.72	157	115	20	21	192 344
Asp bidentate, Site 1											
2.41,3.72	2.65	2.27	2.35	2.52	2.40	2.62	126	106	5	16	0 41
2.40,2.50	4.45	2.24	3.10	2.37	2.38	2.78	108	85	8	15	135 272
2.49,2.53	2.52	2.28	4.23	2.47	2.45	2.71	135	111	6	19	51 247
2.61,2.52	2.63	2.27	4.35	2.51	2.48	2.77	135	112	9	14	91 270
2.49,2.59	2.68	2.34	2.35	2.51	2.59	2.51	130	102	11	18	39 319
4.50,2.27	2.89	2.26	2.36	2.53	2.56	2.77	121	104	4	13	87 281
2.48,2.62	2.53	2.27	2.48	2.48	2.50	2.48	143	117	11	18	89 377
2.47,2.58	2.67	2.31	2.41	2.55	2.54	2.50	130	110	6	14	17 280
2.52,2.53	2.51	2.27	2.44	2.66	2.46	2.48	117	103	4	9	37 278
2.48,2.55	4.61	2.26	2.33	2.35	2.52	2.73	137	113	10	14	58 223
Asp bidentate, Site 2											
2.55,2.41	2.84	2.26	2.43	2.49	2.67	2.52	116	99	4	13	79 77
2.42,2.43	2.75	2.32	2.39	3.98	2.41	2.67	98	77	9	13	49 119
2.42,2.43	2.41	2.24	2.32	4.80	4.66	3.04	123	101	6	15	82 174
2.43,2.46	2.58	2.26	2.33	2.52	4.24	2.69	121	100	9	12	63 205
2.42,2.42	2.43	2.25	2.36	5.00	2.41	2.73	129	101	11	17	93 320
2.45,2.46	2.51	2.29	2.43	2.59	3.05	2.54	116	97	7	12	120 292
2.42,2.46	2.76	2.32	2.45	2.53	2.59	2.50	137	103	16	18	61 295
2.46,2.46	2.48	2.28	2.45	3.38	2.42	2.56	121	105	6	10	64 290
2.35,2.43	2.60	2.27	2.37	3.96	2.39	2.63	105	91	6	8	80 228
2.51,2.45	2.58	2.26	2.43	2.73	2.71	2.52	130	107	12	12	78 247

**Figure S1.** General structure of the EGF34 fragment with the ligands of the two  $\text{Ca}^{2+}$  sites emphasized in green. The best 20 structures are superimposed, using the data in Table S3 (O–O restraints). The left- and right-hand side images were obtained by superimposing domain 3 and 4, respectively.



**Figure S2.** General structure of the EGF34 fragment with the ligands of the two Ca sites emphasized in green and the Ca<sup>2+</sup> ions in black. The best 20 structures are superimposed, using the data in Table S6 (monodentate Ca MM restraints). The left- and right-hand side images were obtained by superimposing domain 3 and 4, respectively.

