# After 50 years, why are Protein X-ray Crystallographers still in business?

SCHOLARONE™
Manuscripts

Principled Models
⇩
Representational Models
⇩
Specific Hypotheses and Generalizations
⇧
Models of Experiments and Data
⇧
The World, Experiments and Data

*Figure 1*: Giere's Hierarchical Approach

**After Fifty Years, Why are Protein X-ray Crystallographers Still in Business?**

Sandra D. Mitchell

Department of History and Philosophy of Science

1017 Cathedral of Learning

University of Pittsburgh

Pittsburgh, PA 15260 USA

smitchel@pitt.edu


Angela M. Gronenborn

Department of Structural Biology

University of Pittsburgh School of Medicine

Pittsburgh, PA 15260 USA, and

Institute for Advanced Study

Berlin, Germany

amg100@pitt.edu

Acknowledgements

**Abstract**

It has long been held that the structure of a protein is determined solely by the interactions of the atoms in the sequence of amino acids of which it is composed, and thus the stable, biologically functional conformation should be predictable by *ab initio* or *de novo* methods. However, except for small proteins *ab initio* predictions have not been successful.  We explain why this is the case and argue that the relationship among the different methods, models and representations of protein structure is one of integrative pluralism. Our defense appeals to specific features of the complexity of the functional protein structure and to the partial character of representation in general. We present examples of integrative strategies in protein science.

**1.    Introduction**

The British chemist John Kendrew in his 1963 Nobel lecture stated:

'… the polypeptide chain, once synthesized, should be capable of folding itself up without being provided with additional information; this capacity has, in fact, recently been demonstrated by Anfinsen *in vitro* for one protein, namely ribonuclease. If the postulate is true it follows that one should be able to predict the three-dimensional structure of a protein from a knowledge of its amino acid sequence alone. Indeed, in the very long run, it should only be necessary to determine the amino acid sequence of a protein, and its three-dimensional structure could then be predicted; in my view this day will not come soon, but when it does come the X-ray crystallographers can go out of business, perhaps with a certain sense of relief, and it will also be possible to discuss the structures of many important proteins which cannot be crystallized and therefore lie outside the crystallographer's purview' (Kendrew [1964]).

In this quote Kendrew characterizes what was to be called the 'problem of protein folding' as 'to predict the three-dimensional structure of a protein from a knowledge of its amino acid sequence alone,' i.e., *ab initio*, using only information about the molecular constituents[1]. Fifty years after Kendrew's Nobel lecture, more than 100,000 protein structures have been deposited in the Protein Data Bank (PDB) [http://www.rcsb.org], *all* of them determined experimentally. Meanwhile, the success of *ab initio* or *de novo* model predictions has been modest at best and restricted to very small proteins. Clearly progress has been made, but X-ray crystallographers are not yet out of business, indeed, they have been joined by other experimentalists who engage an increasing diversity of *in vitro* empirical methods including Nuclear Magnetic Resonance (NMR) spectroscopy, cryo-electron microscopy, small-angle X-ray scattering and a variety of spectroscopic techniques for determining the static and dynamic features of protein structures. In addition, *in vivo* experimental investigations, have exposed difficulties for the unifying goal of predicting functional protein structure solely by models of the physical interactions of the atoms that compose the primary structure. .

The path anticipated by Kendrew, from identification of a complex problem

through analysis of simpler component problems, to arriving at a canonical, typically

mathematical representation that models the system *ab initio* has many precedents in the

history of science.[2] However the current reality of the sciences engaged in elucidating

protein structure from sequence to structure to function[3] is one of multiple methods,

models and representations, investigating different features of this phenomenon in a

variety of contexts.  In this sense, the original 'problem' of protein folding, i.e. predicting

functional structure from the physical interactions of the atoms, which constitute the

string of amino acids in a protein, has not be 'solved'. We argue that a reductive,

eliminativist characterization of a 'solution' to protein structure prediction misdescribes

the role of multiple scientific models in solving complex scientific problems. Indeed, we

present arguments why the development of protein science not only did not, but also was

very unlikely to have turned out as Kendrew predicted it would.  Moreover, we conclude

that pursing integrative, pluralist research strategies will promote rather than hinder

scientific understanding.

Contrary to a widely held belief that science aims at globally unified, true,

complete representations of nature (or some part, thereof), thus relegating anything less

than that to, at best, some immature stage of scientific development, our claim is that

model plurality and partiality are necessary for scientific understanding.  The complexity

of the phenomena investigated conspires with the inherent partiality of scientific

representation to generate pluralities of explanatory and predictive models.  This raises

the question of what relationships hold among the many models.  In what follows we

argue that the correct relationship among multiple models will be one of integration that

maintains pluralism, rather than unification that eliminates all but one fundamental,

4

complete model. We acknowledge that integration can take many forms. The case of functional protein structure we describe below illustrates two kinds of integration of multiple models. The first is in the construction of predictive hypotheses, and the second is in the refinement of empirical models (for other types of integration see Mitchell [2003], [2009]).

In section two we argue that representation must be partial to be scientifically useful. Thus the search for a unique, complete representation for explanation, prediction, and intervention is a fool's errand. Partially of representation grounds the desirability for multiple models. In addition, the complexity of phenomena further supports the variability in methods and models which, in section three, we demonstrate in the case of protein structure. Thus we argue that model pluralism, and the associated multiple methods and representations, is not just a descriptive fact of contemporary scientific practice, it is entailed by partiality of representation and complexity of phenomena.

## *2. Partiality of Representation*

Scientific models of the world involve abstractions from or idealizations of nature.[4] They do not map one-to-one onto the entirety of the undescribed world. Every representational model of a phenomenon, whether it is logical, linguistic, mathematical or graphical, leaves some features out (and correlatively, highlights what is left in). That a representation is a model *of* a given phenomenon is determined by an agent, by someone deciding to use it *as standing for* the phenomenon (Giere [2010], VanFraassen [2010]). Scientific models are judged for their ability to help us explain and predict what goes on in nature. To be successful they need to capture (by similarity, isomorphism, structural or

causal mirroring etc.) features that are relevant to the processes and events we want to

understand and on which we might be able to intervene in order to produce or prevent

effects of interest (see e.g. Bokulich [2011]). What count as relevant, for most scientific

purposes, will be features of the phenomena that are causally involved in the processes

we investigate. Much philosophical work has been directed to identifying what makes a

feature causally relevant (statistically relevant, physically conserved, mechanistically

productive, or counterfactually a difference-maker). What is true for all these accounts is

that not every describable feature of a system in every possible degree of precision is

required for identifying that which permits prediction, explanation and intervention on

that system. We do not need to have a complete representation, in that sense, for

successful science.

Indeed, if we met this strong standard for completeness, the description or model

would fail to be a representation, it would be a duplicate. For the purposes of facilitating

explanation and prediction, it would be no better than engaging directly with the very

system we are trying to understand (see also Truesdell, [1980], p. 72). To identify causes

and generate explanations and predictions, a representation cannot be complete in the

sense of including all that could be described. It must be partial. What is represented and

what is left out are usually tailored to meet some explanatory or pragmatic goal.[5]

For example, the features represented to model a system differ, if your interest is

to produce an effect or to prevent an effect. Take an organism exhibiting symptoms of

AIDS. To model how the HIV virus can generate the symptoms of AIDS requires

articulating sufficient processes and details of the viral life-cycle and the disease: such as

the virus binding and entering the host cell, reverse transcription of the viral genome,

entry into the nucleus, integration of the viral genes into the host genome, transcription and translation of viral proteins, generation and release of new viruses and their maturation, eventually leading to the breakdown of the host's immune system. For a successful model of preventing AIDS, one needs to identify one or more necessary steps in the process, on which one could intervene therapeutically.  For example, if one could block the viral entry into the host cell this would effectively prevent the disease. However, attempts to do this have been futile.  The most successful therapy at present involves blocking of enzymes in the immediate post-entry viral life cycle.

Now consider the (not uncommon) use of different, partial models of the same target phenomenon, one built from principles of fundamental physics and another derived from experiment.  Even if building a complete model is not the explicit goal, many scientists have the intuition that these types of models ultimately can be unified and eliminate the need for empirically provided information from the second model.  One of the motivations for believing that unification through reduction models based on fundamental physical interactions is always possible is the consensus assumption that the objective world, studied by science, is made up in its entirety of physical matter. Everything in nature, simple or complex, is built from this basic material.  To conclude from this core *ontological* assumption that there will be a corresponding unification of *knowledge* in terms of representations from physics is not warranted.  Why not?  The problem is that such a conclusion conflates a commitment to compositional materialism (there is one kind of 'stuff' from which all things are created) with descriptive monism (there is a privileged, complete description of the world in the language of fundamental physics) (Dupre [1993], Mitchell [2009]).  Of course, we agree that everything at its core

is physical, but *representations* and *explanations* of physical matter when engaged in

chemical interactions or biological activities are unlikely to be exhaustively described in

models of physics with no need for models of higher scale behavior.

Developing and using multiple scientific models has been characteristic of the

study of complex phenomena that are sensitive to contextual interactions. In complex

phenomena, there are always multiple features, often at multiple scales or levels of

organization, involved in causal behavior.  Minimally, to carve out a phenomenon to be

investigated we divide the world into the system under investigation and the context,

boundary or background conditions.   We then further simplify by distinguishing within

the system a subset of features to measure and represent to produce a model of the

phenomena. Nature affords many different ways to represent it, both in terms of the

content of the abstraction, i.e. the features selected and those left out and the degree of

abstraction, i.e. how coarse or fine-grained are the decompositions of the selected

features. One might maintain that what is left out is irrelevant to explaining or predicting

the target behavior.  But as the case of protein structure illustrates, this is not always, or

even usually, the case.[6]  Indeed, any system, subject to evolution by natural selection will

always involve a changing, interacting and open-ended set of internal and external

features that contribute historically and proximately to the function of the system.

Especially in multi-component, multi-level complex systems, what is labeled context, or

treated as boundary conditions in a model, often contains causal variables.  The explicit

representation of these variables in other models (with different targets or of coarser or

finer grain) can allow increased understanding and opportunities for intervention.  Which

degree of precision, and to some extent which features to target, will reflect pragmatic

concerns.  But even once those are specified, multiple, and always partial, modeling

options remain.

In this section we have presented philosophical reasons for why we should expect

multiple, compatible models to explain and predict complex phenomena. Given this

plurality, and the partiality of representation, we argue that multiple empirically adequate

models should be and, in fact, are related by integration, even though they are not

completely inter-translatable or reducible.   In what follows we will explore these features

of scientific models in the case of protein structure prediction.  First, we will establish

how complexity requires a plurality of models of functional protein structure.   We will

then turn to two examples of integrative strategies.


*3.  Functional Complexity*

The intricate details of the crowded cellular environment in which a protein folds, are

typically not included in *ab initio* algorithms for predicting structure. For example, the

myriad of small molecules such as water, ions, metabolites etc. and macromolecules such

as proteins, nucleic acids, lipids, carbohydrates, etc. are not included in these models.

There are also a variety of complexities in the functional protein itself that will escape

explanation or prediction from a model of the atoms of its defining sequence of amino

acids alone.

In the early days of molecular protein science, it was assumed that the

thermodynamically most stable conformation of an amino acid chain would be the

biologically functional conformation of a protein.  That is, the sequence (of amino acids)

was thought to determine structure, which in turn would determine function. While

functioning in a cell in an organism does require a form of stability, this is a dynamically realized and maintained state that sometimes involves molecules beyond the single protein.

There is always more involved in the causal network resulting in function than a single polypeptide chain.  Minimally, there are boundary and context conditions, like temperature and cell location, that must be satisfied.  Many proteins may require, in addition, covalent or non-covalent association with other molecules to be functional. Sometimes those other molecules are proteins, and sometimes they are other types of molecules.[7] We describe below protein structures of increasing complexity, containing additives that are required for fulfilling their functional roles,

The ability to induce chemical reactions, transmit signals and build sub-cellular and cellular structures depends critically upon proteins acquiring and retaining stable *functional* conformations in a cell. Function[8] is always the result of a protein interacting with other components of the cell and its structure-in-isolation will not always depict its functional structure-in-context.

An example of a relatively simple structure-function relationship is lysozyme, the first enzyme to have its atomic structure solved by X-ray crystallography. This is a small protein of 129 amino acids, composed of a single chain, whose function is to provide protection against bacterial infection.  It accomplishes this by breaking the carbohydrate chains in the bacterial cell wall. Lysozyme's structure allows it to interact with the bacterial substrate and achieve its function.  Under the right conditions it, and it alone, does the job.

But proteins often require assistance of other molecules to function.  Sometimes a

small molecule needs to be bound. For example, consider hemoglobin, the protein in blood responsible for supplying oxygen to tissues. Functional hemoglobin requires the embedding of the non-proteinatious heme group into the center of the structure formed by the globin chains.  Oxygen binds to the iron in the heme, and both the heme and the globin chains are necessary for hemoglobin to fulfill its biological functions. There are many examples, in addition to hemoglobin, of functionally significant small molecule parts attached to proteins.  Another group is chlorophyll that is bound non-covalently by a number of different proteins that function in photosynthesis.

Besides these non-covalent additives to proteins, covalent modifications of proteins also occur frequently.  Such post-translational modifications can occur at any step in the life cycle of a protein.  Common post-translational modifications include phosphorylation, glycosylation, ubiquitination, sumoylation methylation, acetylation, myristoylation, etc., all modulating or enabling specific protein functions.

Sometimes a protein requires the interaction with other proteins to function. For example, T7 DNA polymerase is composed of two proteins and neither component protein can perform the biological function of the complex alone. Furthermore, the two proteins in the complex are derived from two *different* organisms: the bacteriophage protein recruits the thioredoxin protein from the host E. coli and forms a 1:1 heterodimer. It is only when this complex is formed that the T7 DNA polymerase can function to synthesize double-stranded DNA from a single-stranded template.

There is additional complexity to be considered in protein function beyond a single stable protein structure or even a protein complex. Functioning may require an ensemble of different conformers. Despite proteins often being depicted as static

structures, it is well known that they experience significant fluctuations, yielding

ensembles of closely related conformational sub-states (Frauenfelder *et al.* [1988]), only

one (or few) of which are functionally competent (Baldwin and Kay [2009]).  NMR can

characterize such ensembles, both structurally and with respect to their motions at atomic

resolution (Palmer [2004]).  NMR has also revealed flexible regions in many solved

structural ensembles. The detectable motions may range over an array of spatial and

temporal scales, with some regions more constrained than others. For example, flexible

linkers between relatively rigid domains permit these domains to re-arrange in space in

different orientations to recruit diverse binding partners.  Alternatively, plastic regions

may lose flexibility and become more rigidly structured upon binding to other proteins.

Protein plasticity plays a pivotal role in a large number of cellular processes, such

as transcription, translation, signaling and so on. In addition, it explains the multiplicity

of functions that any given protein can perform.  While one protein can have many

functions that are realized by different conformations, it is also the case that many

proteins (or many protein complexes) may be necessary to achieve a single particular

function.

The variety and complexity described above exposes the range of model

components required to capture the functional structure of proteins.  Sometimes the

primary structure of amino acids is enough to confer functional capacity.  Other times,

additional molecules, small and large, are recruited and contribute to the functional

structure.  And other times there is no one fixed structure that is functional, but rather it is

the product or products of interactions that shape the functional structure of a protein.

The structure-in-isolation, modeled by the basic physics of the atoms in the amino acid

sequence does not always fully capture functional complexity.

If all functional protein structure was simply determined by the amino acids composing it, and other molecules in the cell were irrelevant, a single modeling strategy would be adequate for explaining and predicting protein structure. However, this is not the case. We have described the diversity of components required for functional protein structure. What is relegated to the undescribed context in some models is represented by other models. The complexity of the phenomena requires multiple, compatible models.

*4. Modeling Protein Structure*

We have argued in section 2 that some selection of what is represented and what is left out is required by the partial character of abstract and idealized models. We also illustrated in section 3 how the complexity of functional protein structure and diverse goals of inquiry require a plurality of models. In addition, different methods target different features of a protein in developing models of functional structure. In *ab initio* models, the structure is predicted based on the relative total potential energy of alternate conformations. In models derived by X-ray crystallography the ability of the electrons to diffract X-rays allows atomic positions to be specified. In NMR models the properties of nuclei to absorb radio frequencies permit the extraction of atomic distances based on differential magnetization transfer. Models at the cellular level, derived from the analysis of cellular location of proteins and multi-protein interactions, relevant for protein function, aim at a scale that leaves out the details of the atomic level, while focusing on inter-atomic interactions will simplify, idealize or omit certain entities and activities beyond the boundaries of the targeted amino acid sequence. Different approaches,

characterized by a constellation of methods, models and representations investigate

protein structure in different 'natural habitats'. Studying protein structure *in silico*, *in vitro* and *in vivo* engage increasingly complex contexts.

Kendrew recognized the model pluralism found in protein science in 1964, but suggested it was a temporary stage in a relatively young science and that eventually one representation at one scale, targeting inter-atomic interactions alone would be sufficient. Given enough computational power, it was assumed, everything that occurs at every level of organization would be represented by the minimum energy of the atomic configuration. This is an assumption, not a fact. There is no *a priori* reason to think that representations of protein structure from any single modeling approach, partial and tailored to particular goals of investigation as it must be, will map completely without loss onto any of the others. Pluralism is here to stay. Science lives in a world of multiple models. They cannot always, or perhaps even often, be reduced or unified into one, 'complete' model.

Currently predicting functional protein structure, engages multiple models representing the same phenomenon. The features they individually target are not independent, and the methods share a scientific goal. Our thesis is that these multiple models are integrated in the service of that shared goal. In what follows we will examine two cases of integration of multiple models that display the features discussed in sections and 3 above. We will consider first the relationship between principle models based on thermodynamics and the empirical models inferred from X-ray crystallographic and NMR experimental data. We will argue that a hierarchical, confirmation relationship between theory-driven and experiment-driven models fails to capture their integrative,

constructive relationship.  Then we turn to the relationship between models derived from

the two experimental protocols themselves.  Using multiple empirical models produces

superior results than concentrating on improving a single protocol.

### 4.1 Integrating ab initio and experimental models

The relationship between *ab initio* predictive models and the experimentally derived

models is not one-directional. Experimental models are not used solely to test the

correctness of theory; they are also used to refine, improve, and directly contribute

content to ever more advanced theoretical models.  In short, there is two-way traffic; both

experiment and theory inform each other resulting in an integrative feedback loop (see

also Peschard [2011]).  Questions about and limitations on what is relevant to model,

what is possible to measure, what is computationally accessible necessitate what we

identify as integrative practices.

Let's return to Kendrew's Nobel lecture. His prognostication that empirical

models of structure will become redundant with the development of sophisticated *ab

initio* algorithms is reasonable, only if we assume that models of protein structure derived

from first principles (e.g. Newton's and Maxwell's laws) stand primarily in a

confirmation relationship to the empirical models inferred from experiments.  Thus, once

mathematical methods for generating predictions are made reliable, there would be no

further need for the experimental determination of structure.  Indeed, models from

principles could predict structures for proteins that are not amenable to the preparations

necessary for experimental determination.

Giere has developed an account of the relationships of theory to observation mediated through a hierarchy of models (see figure 1). Giere's aim was to explicate the ways in which explicitly false, or idealized models can nevertheless be understood realistically by means of their mediated relationship to target phenomena, which we endorse (Giere [2004], [2006], [2010], see also Suppes [1969], Teller [2010]).[9] However, as we argue below, the focus on confirmation has occluded non-hierarchical constructive relationships between models that are important for understanding the case of protein structure prediction.

**Figure 1 Giere's hierarchy**

Giere proposes that there is an indirect, imperfect relationship between models and the world, but a connection nonetheless. The relationship is one that allows principles to be tested by comparison of the hypotheses that can be derived from them (via representational models) to the hypotheses developed from models of experimental data. Principles, like F=ma, are abstract, and to know 'where in the world to look to see whether or not the laws apply' (Giere [2004], p. 745) requires introducing specific interpretations, yielding a representational model. These representational models are still abstract. For example, F=-kx is a specification of Newton's 2nd law for simple harmonic oscillators, where x is displacement from equilibrium. To be tested, actual springs and masses need to be observed, and the results of those observations (a model developed from the experimental data of observed quantities) can be compared with the hypothesis arrived at from the representational model.

In short, Giere's view is that representational models come from principles via specification, and that experimental models come from observations via generalization.[10] They meet in the middle, so to speak, where each generates a hypothesis or prediction about the phenomenon of interest. Comparing the two intermediate models constitutes 'testing' the principles by observations. While this is certainly an important way in which multiple models are related, a strict hierarchy leaves out other important integrative relationships.[11] Giere's picture does not rule out these other relationships, but it does not make them explicit. Our account of integrative strategies in constructing predictive models extends the set of model-theory relationships beyond confirmation. As we show below, while experimental models of protein structure can and are used to test principle models, they also are used more directly in the construction of predictive hypotheses.

*Ab initio* all-atom models of proteins are, typically, specifications of Newtonian and Maxwellian principles,[12] which, on Giere's framework, are designated representational models. However, in developing predictive hypotheses of protein structure, these representational models standardly involve experimental models in a constructive, not confirming way. Models of protein structure inferred from X-ray crystallography or NMR are used directly in the algorithms that constitute the representational models that then lead to predictions of specific protein structures. This role of experimental models does not conform to Giere's hierarchy and partially explains why experimental models of protein structure have not become redundant.

The integration of *ab initio* and experimental models is due, in part, to computational intractability[13] . Although steric considerations limit the number of

possible conformations to some degree, i.e. no two atoms can occupy the same place, and certain bond rotations are restricted, the number of conformations for a protein is enormous, thus making finding the native structure a daunting task. *Ab initio* or *de novo* physics-based methods that explicitly model the interactions of all atoms are based on the thermodynamic hypothesis - assuming that the functional, native structure of the protein is at its free energy minimum (Anfinsen [1973], Bryngelson *et al.* [1995])[14]. Correctly calculating all interaction energies requires accurate potential functions and any uncertainties and errors of any terms included accumulate to produce large overall uncertainties in the calculated total energy. This effect of error propagation distorts the calculated potential energy surface of a protein, making it impossible to equate the lowest energy structure with the native one (Unger and Moult [1993], Freddolino *et al.* [2010]).

Another difficulty for all-atom models of large proteins is the sampling problem. Starting from an extended polypeptide chain, a major computational challenge is to generate enough conformations sufficiently close to the native free-energy minimum to reliably pick out the lowest energy conformation.  Given the hundreds of thousands of non-native conformations, not every possible conformational energy can be calculated (Alm and Baker [1999]).

Given these problems, prediction by *ab initio* methods has not been very successful except for small molecules and small proteins. Therefore, a range of integrative, semi-empirical computational approaches have been developed that have had greater predictive success.  These make direct use of experimentally derived models, deposited in the Protein Data Bank, in generating a predictive hypothesis.

Homologous structures and other experimental information are now incorporated into computational algorithms to help overcome the sampling problem, making it possible to more reliably locate the native energy minimum by adding experimentally derived constraints (Baker and Sali [2001]). Homology modeling uses sequence similarity between proteins to infer structure similarities. Specifically, if a particular sequence is shared between a protein that has an experimentally determined structure and protein with unknown structure, a predictive model can be developed that assumes that the unknown structure is the same as the experimentally known one. This reduces the problem so that an exhaustive search in conformational space is only required for those parts of the target protein's primary structure that are not similar to known proteins. In this manner, the native structure of the target protein can be predicted by integrating models from X-ray and NMR experiments with *ab initio* models. (Das and Baker [2008]).

Every two years since 1994, the Critical Assessment of Protein Structure Prediction (CASP) provides the opportunity to evaluate different predictive algorithms against new experimental structure determinations that have not yet been made public. In reflecting on trends from the first to the 2012 CASP results, Moult, *et al.* [2014] report that:

The accuracy of homology models, as monitored by CASP, has improved dramatically, through a combination of improved methods, larger databases of structure and sequence, and feedback from the CASP process. *Ab initio*[15] modeling methods have also improved substantially, from a very low base in the first CASP experiment. It is now not unusual to see topologically accurate models

for small (<100 residues),[16] regular, and single domain non-template proteins.

While the success of homology models clearly exemplifies the confirmation relationship Giere characterized in his hierarchical picture of models, more is going on. Giere's picture does not explicitly represent the constructive relationship between experimental data models (deposited structures in this case) and representational models (computational algorithms based on physical and chemical principles applied to proteins) exhibited in the homology modeling strategy. Not only does the integrative, constructive relationship account for the continued importance of experiments, contrary to a prediction of their demise in the wake of increased success of *ab initio* models, it also raises new issues for holding to a simple confirmation relation. The representational models that have been developed for protein structure prediction, are not purely 'top-down' procedures for generating hypotheses, but are at the same time 'bottom-up'. Not all experimental results and data can be used for confirmation.[17]

Not only do experimental results shape representation models, but also the use of principles shapes the models derived from experiments in other integrative practices found in protein structure prediction. NMR parameters, such as resonance frequencies can guide the selection of protein fragments to build up the low-resolution conformation in the initial stage of constructing the experimental model of a protein. Subsequent incorporation of sparse NOE-derived distance restraints[18] (Shen *et al.* [2009]) or residual dipolar coupling-derived orientational restraints permits the determination of protein structures from data in a fraction of the time and with considerably less effort than using traditional NMR structure determination approaches alone. Not surprisingly, this combined methodology is becoming increasingly popular in the NMR community (van

der Schot *et al.* [2013]).

Structure prediction is increasingly the result of integrating *ab initio* and experimental models.  The experimentalists, contra Kendrew, are still in business, but the business is now not just using experimental technologies to confirm hypotheses of particular protein structures, it is also to provide constructive components in the development of computational predictive algorithms.

*4.2 Integrating multiple experimental models*

X-ray crystallography and NMR are both applied to proteins that are removed from their native, cellular environments, which have been purified, physically and chemically to produce homogeneous materials. For X-ray crystallography, the proteins have to be in a crystalline state, i.e. an ordered three-dimensional array of molecules, whereas for solution NMR, the proteins are dissolved and move randomly in aqueous solution. Thus in the two types of experiment they are in physically different states, solid and liquid. Furthermore, X-ray crystallography and NMR target different atomic features of a protein: X-ray crystallography relies on the scattering of X-rays by the electrons while in NMR interactions of nuclear spins with a magnetic field are exploited. In addition, the nature of the data that is collected is also different: in X-ray crystallography, reflections (spots on a film or CCD camera) caused by diffracted X-rays are recorded, the positions and intensities of which contain information about atomic positions; in NMR resonance frequencies of nuclear spins and their interaction energies (cross-peaks on a piece of paper or a computer screen) are recorded, the magnitude of which contains information about interatomic distances or angles. Both, atomic positions and inter-

atomic distances and angles, are used to develop three-dimensional molecular models of the proteins.

Naturally, at each of the above stages, preparation of the material for the investigation, the application of a specific instrumental method and the use of the associated equipment, the type and method of data collection, as well as the computational algorithms applied for transforming the data into molecular models, noise, error, and incompleteness is present.[19]

Since there is no direct access to the atomic three-dimensional structure of a protein *in vivo*, there is no direct way to determine which model from *in vitro* experiments most *accurately* reflects the protein structure in the context in which it performs its biological functions. On the other hand, statistical and empirical approaches can estimate the *precision* of the respective models compared to hypothetical error-free or complete models (Brünger [1997]). In X-ray crystallography the R-factor is a measure of the agreement between the data and the model and in NMR models, violations of distance and angular constraints indicate model uncertainty.

While there are other experimental methods that can probe protein structure, with cryo-electron microscopy becoming more important, X-ray crystallography and NMR are the predominant methodologies for deriving empirical models of protein structure at the atomic level. Both methods clearly target different features of atoms in the folded protein, generate different kinds of data, and rely on different algorithms that are used to infer a model of a given protein structure. Each method is blind to some features of the protein that are relevant for its function, each method abstracts away or idealizes different potentially significant factors, and each method introduces different sources of bias in

generating data and inferring models from data. Multiple models of the same protein can be produced by multiple experimental protocols. Our question is, how are these multiple models related to each other.

*A priori*, any two partial models of a phenomenon, which encode different abstractions or idealizations, if they are comparably accurate, can contribute to scientific understanding of the phenomenon. There is reason to believe that simple addition of partial models will not yield more accurate models if the features they target or the subdomains of the phenomenon are causally independent, though they will represent more features of the phenomenon. For example conjoining the models of the different stages of HIV infection may add up to knowing more of the path of the disease, but not of knowing any one stage better. However, if the targeted features are causally relevant to the outcome of interest, and not independent of each other, using two or more models together can mutually correct inaccuracies due to systematic noise or bias in the single models. Indeed, continuing improvement in one method may yield more precise models, but not necessarily more accurate ones. Alternatively, if multiple models are used jointly, though not additively, models of the phenomenon can improve in accuracy.

There are several avenues available for integrating X-ray crystallography and NMR and their resultant models in the case of protein structure. An NMR structure can be used to aid solving a crystal structure by molecular replacement (Brünger *et. al* [1987]) and a crystal structure can be used as an input model for an NMR structure determination (Delaglio et. al [2000]). However, crystallographic and NMR models of a given protein can diverge, exhibiting differences that are beyond the experimental errors of the measurements. Where do these differences come from? They certainly can be

caused by differences in the environment (a protein in a crystal lattice, ordered in a three-dimensional array, versus a protein tumbling randomly in solution). Alternatively, though not as significant, they may be due to details in the different computational algorithms , which generate the model from the experimental data.

Relying on a single method can lead to erroneous inferences that can be exposed by comparison with another method. The protein APOBEC2 (A2) is a muscle specific family member of the APOBEC/AID (Activation Induced Deaminase) family of cytidine deaminases. A truncated version of A2 was purified, crystallized and its X-ray structure was solved, revealing an extended V-shaped homotetramer (Prochnow, et al. [2007]. This A2 structure had considerable impact in the HIV field since it was considered a good model for the structure of A3G, another member of the APOBEC/AID family and an important HIV restriction factor. The tetrameric structure and the arrangement of the monomers of A2 in the crystal were taken as representative for A3G and used to explain the latter's enzymatic and anti-HIV activity.  However, in contrast to the findings in the crystal, A2 is monomeric in solution, and the NMR solution structure of the full-length A2 revealed that the N-terminal tail component of the complete protein that was removed for crystallization, is positioned such, that it interferes with the tetramer interface (Krzysiak et al. [2012]). Therefore, any conclusions as to the functioning of A3G on the basis of the protein conformation found in the A2 crystal were erroneous. The disagreement of models of the same protein generated by X-ray and NMR protocols does not mean that one or the other method produced flawed results.  Each got the 'right' model of its prepared sample of the A2 protein.  But given the variation in the different *in vitro* contexts, there was no longer warrant for inferring that the homologous A3G protein

has features similar to the X-ray model of A2. It may instead be closer to the NMR A2 model, or unlike either.  As a consequence, given the dependence of the A2 model on the different *in vitro* experiments, the original inferences from the X-ray model had to be questioned.

Systematic noise or model biases in the two techniques may occur in the development of a model for a protein target.  Using a joint refinement approach, an overall better model of a protein structure can be derived by combining X-ray and NMR data ((Shaanan, *et al.* [1992]). This type of integration reduces the underdetermination in the models inherent to each methodology. Using data from both methods in refinement reduces the total range of possible models by mutually correcting individual model bias. No matter how many technical advances will contribute to improving the accuracy and the precision of generated data and inferential models, there will always be relevant factors that cannot be represented in either X-ray or NMR based models. While it is clear that a given protein will have, under specific conditions, a given structure, neither experimental/inferential protocol is expected to perfectly, or completely detect it and different environments can further modulate the molecular behavior that influences the targeted atomic properties. Integrating multiple models from different experimental protocols provides a means to reach more accurate results than relying on any single method.  We have shown this to be the case not only when both methods and models converge, providing support of accuracy through consilience, but also when they diverge, providing mutual correction.

*5 Conclusions*

The goal of solving the problem of predicting protein structure articulated by Kendrew has been confounded by the complexity and variability of functional proteins. The envisioned solution in terms of *ab initio* models from first principles has faced so far insoluble computational difficulties and the best models have come from semi-empirical methods incorporating both principles of basic atomic physics and experimentally determined empirical models. Acquiring knowledge necessary for dealing with real world protein function and malfunction by limiting the investigation to proteins-in-isolation will fail. We suggest that to understand protein structure in ways that can contribute to achieving our pragmatic goals requires multiple models, methods and representations. The diversity of approaches includes atomic interactions modeled *in silico*, position and distance measures experimentally modeled *in vitro*, and functional interactions modeled i*n vivo*.

We have presented both philosophical (sections 1 and 2) and scientific (sections 3 and 4) reasons why the relationships among multiple methods, models, and representations engaged in a quest to understand protein structure is one of integrative pluralism. Each method, model and representation provides a partial view of the phenomena. There is no reductive or unified complete model that can replace this pluralism. No single theoretical or empirical approach targets all the features that are relevant to the structure of functional proteins. What initially looked like a single question, how do we predict the structure required for function from sequence, proliferated into a multitude of sub-questions, the answers to which cannot simply be added together to form a unified, single solution (Dill and MacCallum [2012]). The descendent problems of protein science have formed new fields of research, populated by

a rich, variegated plurality of integrative scientific practices. Should there be concern

that the state of protein science has not reached the kind of 'final result' that Kendrew

anticipated? We think not. The philosophical framework of integrative pluralism we

present here, that acknowledges the partiality of representation and the pragmatics that

shape the selection of relevant features, explains why the current pluralism of methods,

models, and representations will and should endure.

References

Alm, E. and Baker, D. [1999]: 'QPrediction of protein-folding mechanisms from

free-energy landscapes derived from native structures',Q Proceedings of the National

Academy of Science U S A 96, pp. 11305-10.


Anfinsen C. B. [1973]: 'Principles that govern the folding of protein chains', *Science*

**181** (4096)c pp. 223–30.


Anfinsen, C. B., Haber, E., Sela, M. and White, Jr. F. H. [1961]: 'The Kinetics Of

Formation Of Native Ribonuclease During Oxidation Of The Reduced Polypeptide

Chain', Proceedings *of the National Academy of Science*, **47**(9), pp. 1309–14.

Bailer-Jones, D. [2009]: *Scientific models in philosophy of science*, Pittsburgh: University of Pittsburgh Press.

Baker, D. and Sali, A. [2001]: 'Protein structure prediction and structural genomics', *Science* **294** (5540) pp. 93–6.

Baldwin, A. J. and Kay, L. E. [2009]: 'NMR spectroscopy brings invisible protein states into focus', *Nature Chemical Biol*ogy, 5, pp. 808-14.

H.M. Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N., and Bourne, P.E.: [2000)] 'The Protein Data Bank' Q*Nucleic Acids Research,* 28, pp. 235-42.

Bogen, J. and Woodward, J. F. [1988]: 'Saving the Phenomena', *Philosophical Review*, **97** (3), pp. 303-52.

Bokulich, A. [2011] 'How Scientific Models Can Explain', *Synthese*, **180** (1), pp. 33-45.

Borges, J. L. [1998]: 'On exactitude in science', *Collected Fictions*, New York: Penguin.

Brünger, A. T. [1997]: 'X-ray Crystallography and NMR: Complementary Views of

Structure and Dynamics', *Nature Structural Biology* **4**, pp. 862-5.

Bryngelson, J. D., Onuchic, J. N., Socci, N. D. and Wolynes, P.G. [1995]: 'Funnels, Pathways, and the Energy Landscape of Protein Folding: A Synthesis', *Proteins: Structure, Function and Genetics.* **21** (3), pp. 167–95.

Cartwright, N. D. [1999]: *The Dappled World: A Study of the Boundaries of Science*, Cambridge: Cambridge University Press.

Chang, H. [2004]: *Inventing Temperature: Measurement and Scientific Progress.* Oxford Studies in the Philosophy of Science, New York: Oxford University Press.

Das, R., and Baker D. [2008]: 'Macromolecular modeling with Rosetta', *Annual review of biochemistry*, **77**, pp. 363-82.

Delaglio, F., Kontaxis, G. and Bax, A. [2000]: 'Protein structure determination using molecular fragment replacement and NMR dipolar couplings', *Journal of the American Chemical Society,* 122, pp. 2142-3.

Dill, K. A. and MacCallum, J. L. [2012]: 'The protein-folding problem, 50 years on', *Science* **338**:  pp. 1042-6.

Driscoll, P.C., Gronenborn, A.M., Beress, L. and Clore, G. M. [1989]: 'Determination

of the three-dimensional structure of the anti-hypertensive and anti-viral protein BDS-I from the sea anemone Anemonia sulcata: a study using nuclear magnetic resonance and hybrid distance geometry-dynamical simulated annealing', *Biochemistry* 28, p. 2188.

Dupre, J. [1993]: *The Disorder of Things. Metaphysical foundations of the disunity of science*, Cambridge: Harvard University Press.

Epstein, B. and Forber, P. [2013]: 'The perils of tweaking: how to use macrodata to set parameters in complex simulation models', *Synthese* **190**: pp. 203-18.

Ernst, R. R. [1992]: 'Nuclear magnetic resonance Fourier transform spectroscopy (Nobel lecture)', Angewandte *Chemie International Edition,* 3, 805-23.

Frauenfelder, H., Parak, F. and Young, R. D. [1988]: 'Conformational Substates in Proteins', *Annual Review of Biophysics and Biophysical Chemistry* **17**: pp. 451-79.

Freddolino, P. L., Harrison, C. B, Liu, Y and Schulten, K. [2010]: 'Challenges in protein-folding simulations', *Nature Physics* **6**, pp. 751−8.

Giere, R. N. [1999]: *Science without Laws*, Chicago: University of Chicago Press.

Giere, R. N. [2004]: 'How Models are Used to Represent Reality', *Philosophy of Science* **71**(5): pp. 742-52.

Giere, R. N. [2006]: *Scientific Perspectivism*, Chicago: University of Chicago Press.

Giere, R. N. [2010]: 'An Agent-based conception of models and scientific representation', *Synthese*, **172**, Issue 2, pp. 269-81.

Hon, G. [1989]: 'Towards a Typology of Experimental Error: an Epistemological View', *Studies in History and Philosophy of Science* 20, pp. 469–504.

Hüttemann, A. and Love, A. C. [2011] 'Aspects of Reductive Explanation in Biological Science: Intrinsicality, Fundamentality, and Temporality', *British Journal for the Philosophy of Science*, **62** (3):519-549.

Kendrew, J. C. [1964]: 'Nobel Lecture: myoglobin and the structure of proteins', *Nobel Lectures, Chemistry 1942-1962*, Elsevier, Amsterdam.

Krzysiak, TC, Jung, J., Thompson, J., Baker, D., Gronenborn, A.M. [2012] 'APOBEC2 is a monomer in solution: Implications for APOBEC3G models', Biochemistry 51, pp. 2008-17.

Levinthal, C. [1968]: 'Are there pathways for protein folding?', *Journal de Chimie Physique et de Physico-Chimie Biologique* **65**: pp. 44–5.

Levinthal, C. [1969]: 'How to Fold Graciously'. *Mossbauer Spectroscopy in Biological Systems: Proceedings of a meeting held at Allerton House, Monticello, Illinois*: pp. 22–4.

Levitt, M. and Warshel, A. [1975]: 'Computer Simulation of Protein Folding', *Nature* **253**, pp. 694-8.

Mitchell, S. D. [2003]: *Biological Complexity and Integrative Pluralism* Cambridge: Cambridge University Press.

Mitchell, S. D. [2009]: *Unsimple Truths: Science, Complexity and Policy*, Chicago: University of Chicago Press.

Morgan, M. S. and Morrison, M. (eds). [1999]: *Models as Mediators. Perspectives on Natural and Social Science*, Cambridge University Press.

Morrison, M. [2000]: *Unifying Scientific Theories*, Cambridge: Cambridge University Press.

Moult, J., Fidelis, K., Kryshtafovych, A., Schwede, T. and Tramontano, A. [2014]: 'Critical assessment of methods of protein structure prediction (CASP) — round X', *Proteins*, **82**: pp. 1–6.

Ong, D. S. T. and Kelly, J. W. [2011]: 'Chemical and/or Biological Therapeutic

Strategies to Ameliorate Protein Misfolding Diseases', *Current Opinion Cell Biol.* 23 (2): pp. 231-8.

Palmer, A. G. 3rd [2004]: 'NMR characterization of the dynamics of biomacromolecules', *Chemical Review*, **104**(8), pp. 3623-40.

Peschard, I. [2011]: 'Making sense of modeling: beyond representation', *European Journal for Philosophy of Science*, **1**(3), pp. 335-52.

Prochnow, C., Bransteitter, R., Klein, M. G., Goodman, M. F., and Chen, X. S. [2007]: 'The APOBEC-2 crystal structure and functional implications for the deaminase AID', *Nature* **445**, pp. 447-51.

Russel, D., Lasker, K., Webb, B., Velazquez-Muriel, J., Tjioe, E., Schneidman-Duhovny, D., Peterson, B., and Sali. A. [2012]: 'Putting the pieces together: integrative structure determination of macromolecular assemblies', *PLoS Biology* **10**, p.e1001244.

Seck, C. [2012]: 'Metaphysics within Chemical Physics: The Case of ab initio Molecular Dynamics', *Journal for General Philosophy of Science* **43**, pp. 361-75.

Shen, Y, Vernon, R., Baker, D, and Bax, A. [2009]: 'De novo protein structure generation from incomplete chemical shift assignments', *Journal of Biomolecular*

*NMR* **43**, pp. 63-78.

Shaanan, B., Gronenborn, A.M., Cohen, G. H., Gilliland, G. L., Veerapandian, B., Davies, D. R., Clore, G. M. [1992]: 'Combining experimental information from crystal and solution studies: joint X-ray and NMR refinement', *Science 257*, pp. 961-4.

Shen, Y, Lange, O. F., Delaglio, F, Rossi, P., Aramini, J.M., Liu, G., Eletsky A., Wu, Y., Singarapu, K.K., Lemak, A., Ignatchenko, A., Arrowsmith, C.H., Szyperski, T., Montelione, G.T., Baker, D. and Bax, A. [2008]: 'Consistent blind protein structure generation from NMR chemical shift data', *Proceedings of the National Academy of Sci*ence USA **105**, pp. 4685–90.

Suppes, P. [1969]: *Studies in the Methodology and Foundations of Science: Selected Papers from 1951 to 1969*. Dordrecht: Reidel.

Tal, E. [2013]: 'Old and New Problems in Philosophy of Measurement', *Philosophy Com*pass 8 (12), pp. 1159–73.

Teller, P. [2010]: '"Saving the Phenomena" Today', *Philosophy of Science*, Vol. 77, No. 5 pp. 815-26.

Unger, R. and Moult, J. [1993]: 'Finding the lowest *free energy conformation of a protein is a NP-Hard problem: proof and implications', Bulletin of Mathematical*

*Biology,* 55, pp. 1183–98.

Truesdell, C. [1980]: 'Statistical Mechanics and Continuum Mechanics', in *An Idiot's Fugitive Essays on Science.* New York: Springer-Verlag, pp. 72–9.

van der Schot, G., Zhang Z., Vernon R., Shen Y., Vranken W. F., Baker D., et al. [2013]: 'Improving 3D structure prediction from chemical shift data', *Journal of Biomolecular NMR,* **57**(1), pp. 27-35.

Van Fraassen, B. C. [2010]: *Scientific Representation: Paradoxes of Perspective.* Oxford: Oxford University Press.

Weisberg, M. [2012]: *Simulation and Similarity: Using Models to Understand the World*, Oxford: Oxford University Press.

Winsberg, E. [2010]: *Science in the Age of Computer Simulation*, Chicago: Chicago University Press.

Wüthrich, K. [2003]: 'NMR Studies of Structure and Function of Biological Macromolecules (Nobel Lecture)', *Angewandte Chemie International Edition*: 42, pp. 3340–63.

**Figure 1.** Giere's Hierarchical Approach

[1] Proteins are macromolecules involved in all functions necessary for life. The sequence of a protein is encoded in DNA, which is transcribed into messenger RNA, which is translated by the ribosome into the chain of amino acids.  A specific number and sequence of amino acids constitutes the primary structure a protein. The three-dimensional structure of the chain of amino acids in space in called the tertiary structure of a protein. This is what is called "protein structure" here and throughout this article. The terms "ab initio" and "de novo" are used throughout this article as equivalent, following current usage in protein science.

[2] Indeed, there was evidence from the experiments performed by Anfinsen (Anfinsen et al. [1961]), as Kendrew references, that the denatured protein ribonuclease spontaneously refolds, suggesting that no other molecules besides the string of amino acids is required for acquiring the native conformation.

[3] While Kendrew cast the problem in terms of inference from sequence to structure, even then it was recognized the ultimate goal is understanding how the structure, into which a proteins folds, is key to understanding how it functions *in vivo*.  See also Hüttemann and Love [2011].

[4] Much has been written on this, see, for example, *Fictions in Science: Essays on Idealization and Modeling*, ed. Mauricio Suárez, London: Routledge, 2009.

[5] This type of selection encodes what is sometimes referred to as a "perspective", see Giere 2006 and Van Fraassen 2010.

[6] One can always shrink the explanatory goal to match the scope of the model to avoid this criticism. However, looking only within the abstract, ideal world defined by a model will not address the real world problems science aims to solve.

[7] The first three dimensional atomic structure of a protein that was solved by X-ray crystallography by Kendrew – myoglobin – is one that embeds a heme or iron-containing non-polypeptide component in its structure.

[8] Examples of protein function include controlling sugar levels (by hormones), combatting infection (by antibodies), digestion (by enzymes), etc.

[9] There has been substantial philosophical research on scientific models – from the semantic view of theories being best understood as models, to their autonomy from theories, to the many functions they perform in science (Van Fraassen [1980], Morgan and Morrison [1999]). We are considering Giere's approach, in part, because we share with him commitments to the pluralism, partiality of models and pragmatism, but find the hierarchy picture inadequate. Giere does discuss a constructive relationship between data modes and representational models, but only in cases where there are no principles available that could generate a representational model. The problem we consider is one in which principles play a role and it was hoped that with minimal interpretation (atoms in a polymer, force equations) would be sufficient to predict protein structure. Indeed, we interpret Kendrew's comments that with those minimal specifications we would not need to derive models from experiments at all. We are arguing, contra Giere that even when there are well-established principles, there are still constructive relationships that require the use of data models to derive a hypothesis contrary to a strict hierarchy.

[10] Of course, there is much more to generating either of these models. Other assumptions constrain the interpretations of a formal principle in application to a particular context. And, in producing a model from data, decisions are made about outliers, test conditions, sources and degrees of error, theory of instruments, and other things, see Bogen and Woodward [1988].

[11] See Epstein and Forber [2013] for consideration of the pros and cons of tweaking micro simulations by using macrodata.

[12] Due to computational limitations, quantum mechanical principles have been used only to characterize small molecules, although there have been recent developments to combine quantum mechanical and molecular dynamics approaches for studies of protein structure.

[13] Levinthal calculated that for a fairly small protein, composed of 100 amino acids, with each amino acid capable of adopting three possible conformations (an underestimate), $3^{100}$ or $5 \times 10^{47}$ different three-dimensional protein structures can potentially be realized (Levinthal [1968], [1969]).

[14] In its simplest form, the total potential energy is expressed as $E = E_{covalent} + E_{noncovalent}$, which can be further decomposed into $E = E_{bond} + E_{angle} + E_{dihedral} + E_{electrostatic} + E_{van\ der\ Waals}$ (Levitt and Warshal [1975]).

[15] Here *ab initio* is used in the sense of '*de novo*'; i.e. there is no closely related homologous structure available in the protein data base

[16] A residue is an amino acid.

---

[17] Other experimentally derived constraints can also be used to guide the search towards the native folded structure, such as provided by chemical crosslinking, small angle X-ray scattering, or fluorescence energy transfer experiments (Russel *et al.* [2012]).

[18]The Nuclear Overhauser Effect (NOE) is the transfer of magnetization from one nuclear spin to another via cross-relaxation. Its size is related to the interatomic distance r with a $1/r^6$ dependence.

[19] Hon [1989] provides a more nuanced elaboration of sources of error in experiment, rejecting the claim that systematic vs. random error is sufficient. However, he does not address the issues arising from the partiality of experimental models. Recent work on experimental measurement, for example, Chang [2004] and Tal [2013], discusses how two different experiments that target different features are taken to be measurements of the same phenomenon, appealing to robustness of outcomes, and the role of models in securing reliability. While different experiments generating the same result has been taken as grounds for reliability, we explore here the ways in which different experiments getting different results might contribute to knowledge of the target phenomenon.