

[IBM Developer](#)[Topics](#)[Community](#)[More open source at IBM](#)

We've done a little recoding. Introducing IBM Developer.
The next level of coding, content, and community.

[Learn
more](#)
[Learn](#) > [Linux](#)

Contents

Protein modeling

Introduction

Real-world scientific advances through
Proteomics: The protein economy

Commercial and academic motivations
T.J. Christopher Ward and Ruhong Zhou

Published on June 09, 2009
So what are we modeling?

Equipping the laboratory

In 2001, IBM's research scientists started
Running the model
Gene®. These servers have been available
What does all this give us?

Predicting the future
The Blue Gene family of supercomputers

Environment, they're also designed to provide
universities, government, and commercial

Related topics
protein folding, climate research, cosmology, and drug development. The system is making
magnitude, in the way science can be done, because it offers a more cost-effective tool
alternative versions of complex models.

In this article, we present some of the progress that has been made by one of the projects
modeling. Figure 1 shows the scale of work we can do now, thanks to the power of Blue
starts from the lysozyme crystal structure (see [Related topics](#) for source).

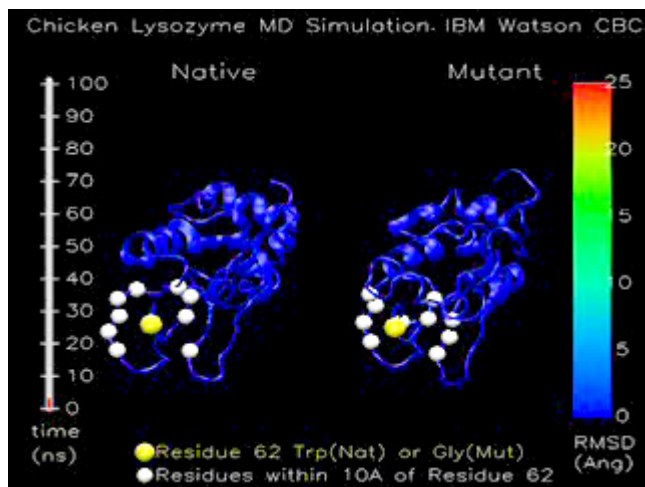
Figure 1. Part of the total ten microseconds of life inside a living cell; [watch the video](#)

with Blue Gene/

deling and data visualization on a

sign of a new family of servers, now made
since 2004—first the Blue Gene/L (wh

ned to deliver ultra-scale performance
iciencies in power, cooling, and floor-s
ch labs use Blue Gene for computation



Contents Proteomics: The pr

Introduction

Proteins are biological macromolecules within cells. Many proteins are enzymes that catalyze biochemical reactions; some are involved in immune responses; many others have structural and mechanical functions for muscles and bones. Commercial and academic motivations for studying proteins illustrate how pervasive and important they are.

So what are we modeling?

One protein is responsible for the "clotting" of blood; it carries oxygen from the lungs to the rest of the body.
 Equipping the laboratory
 Another protein is responsible for the "digestion" of food.
 Running a normally harmful.

What does all this give us?

There are hundreds of thousands of proteins in a cell, and how their diversity affects the function of the cell.
 Predicting the future
 How proteins are made, how they are modified, and how they interact with each other.
 Downloadable resources

DNA is the information storage component of every cell in every plant and animal. It stores the genetic information that is passed from parent to offspring.

Related topics

Chemical building blocks (nucleotides) we call **A**, **C**, **T**, and **G** (for adenine, cytosine, thymine, and guanine). From a distance, these building blocks look very similar, so they form the same overall shape—the famous Watson-Crick Double Helix.

To read out the information in the DNA, the DNA untwists and another molecule called RNA is synthesized. The RNA then serves as a blueprint to the ribosome, a protein that behaves like an all-purpose factory. The ribosome reads the A/C/T/G code in groups of three, allowing us to derive a 64-letter "alphabet."

Twenty of these "letters" correspond to amino acids, the building blocks for proteins. The ribosome uses the code to assemble the protein. The food we eat (humans cannot synthesize all the amino acids we need and therefore need "essential" amino acids, from food). Each amino acid has a "head" and a "tail." The ribosome

n economy

Proteins are an essential component of organisms and they catalyze biochemical reactions; some are involved in structural and mechanical functions for muscles and bones. Proteins are:

"clotting" of blood; it carries oxygen from the lungs to the rest of the body.
 Equipping the laboratory
 Another protein is responsible for the "digestion" of food.
 Running a normally harmful.

Proteomics is the study of the entire set of proteins involved in life on Earth. Proteomics is the study of the proteins that are involved in the evolution of the living organisms and how they affect the systems they inhabit.

Proteomics is the study of the proteins that are involved in life on Earth. Proteomics is the study of the proteins that are involved in the evolution of the living organisms and how they affect the systems they inhabit.

acid for each "letter" and assembles them head-to-tail in sequence; other "letters" indicate the chemical properties of the amino acids. The resulting linear sequence of amino acids is a newly minted protein molecule, formerly a blueprint, now a reality. The protein is then folded into the shape it takes in the cell, determined by the sequence of amino acids in the section of DNA that was used.

Stresses and strains between the atoms in the protein molecule, interactions with the surrounding environment, and random vibrations that you would call *heat* then cause the protein molecule to "fold" into its final shape.

Protein molecules are quite stable; some of them can exist unchanged for hundreds of years. They are stable at a wide range of temperatures, which would kill the organism that made them. They stay roughly the same shape even when denatured by strong chemicals, high pressure, heat or cold, or by becoming food for another organism.

Contents

The shape and the way it varies with time
will do—whether it will transport oxygen
Introduction
tiny scale.

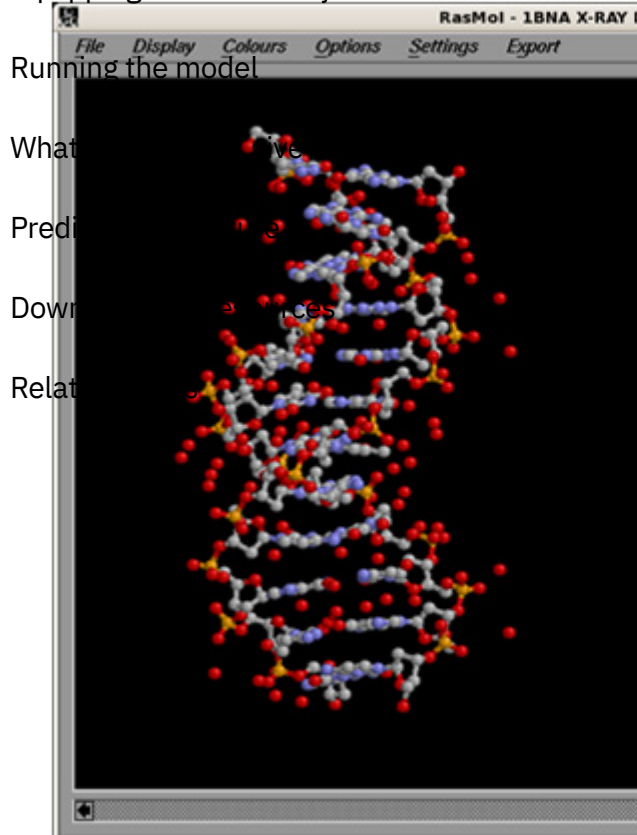
Proteomics: The protein economy

Figure 2 demonstrates the familiar ball-and-stick model of DNA (image is a stereo pair; source):

So what are we modeling?

Figure 2. The ball-and-stick model of DNA

Equipping the laboratory



What is the model?

Predicting the model

Download the model

Related to the model

temperature, and surrounding molecules determine the shape of a protein. You can see a poison-ivy allergy, or do any of the other things that a protein can do.

Figure 2 demonstrates the familiar ball-and-stick model of DNA (image is a stereo pair; source):

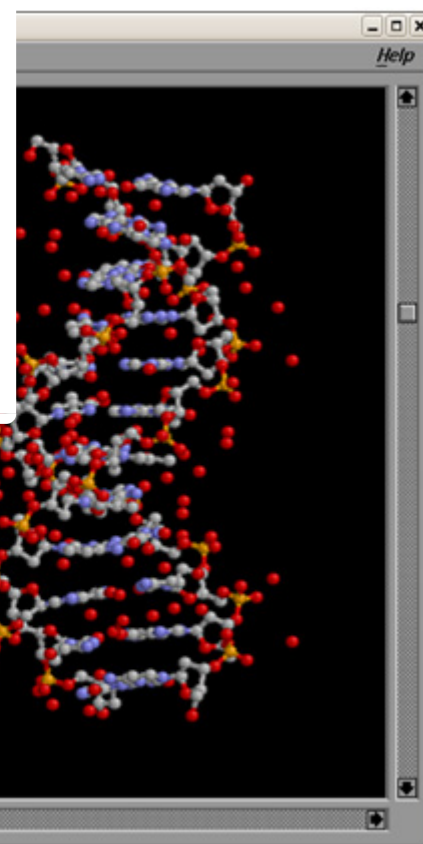
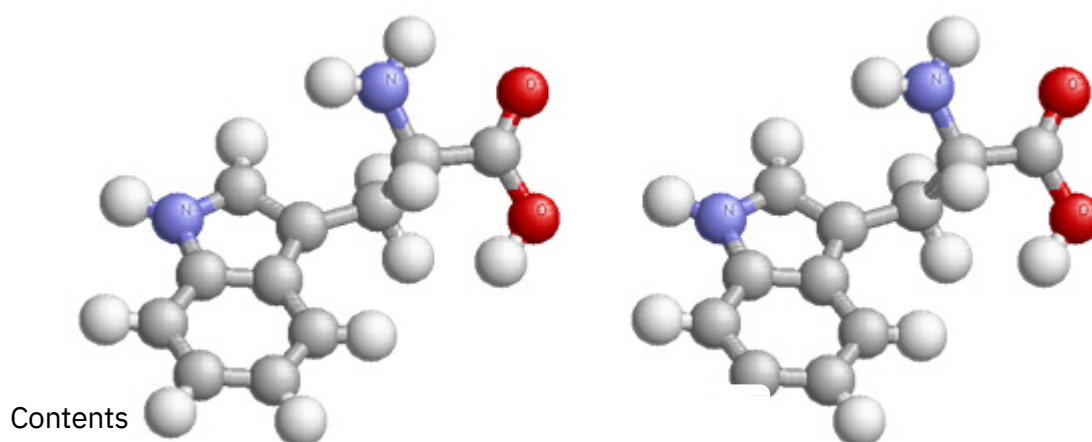


Figure 3 shows tryptophan, one of the 20 standard amino acids (image is a stereo pair; source).

Figure 3. Tryptophan, one of the 20 standard amino acids



Contents

Introduction

Amino acids are assembled into proteins. The protein (top of Figure 3) of an amino acid is the part of the molecule that is away from the N (top of Figure 3) of an amino acid. All amino acids have a common structure. Commercial and academic motivations

Figure 4 offers a visual look at the protein structure of hemoglobin.

Equipping the laboratory
Figure 4. The protein hemoglobin

Running the model

What

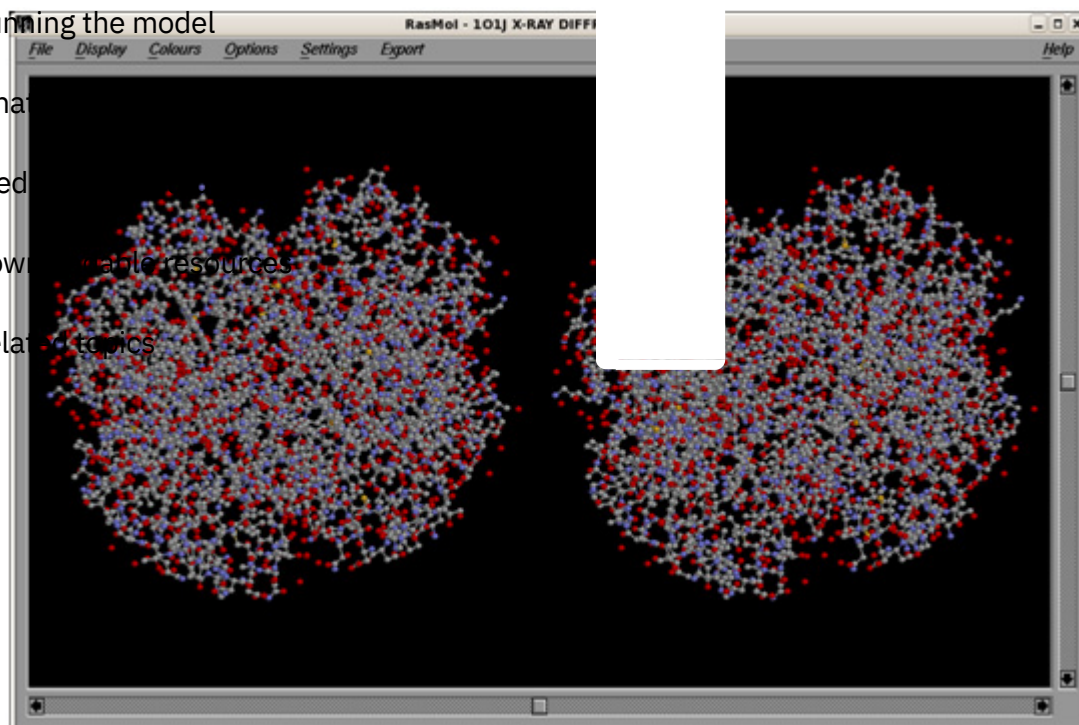
Pred

Down [related resources](#)

Related topics

...tting off the O-H group (right side of Figure 3) of the molecule, and splicing the molecules together to form a characteristic atomic grouping (top right of Figure 3).

...oglobin (image is a stereo pair; see [Related topics](#))



Hemoglobin is a total of 574 amino acid molecules in 4 subunits. Hemoglobin, with its assembly into the protein is beyond the scope of this article), transports oxygen around the body. The transport system is possible with just the iron atoms, but it is very much more effective with the help of hemoglobin provides. If you put this image into a stereo viewer, you can pick out the iron atoms.

more complex than this, we need a different way to visualize what is going on.

Commercial and academic motivations

Increasingly, advances in designing pharmaceuticals and protecting public health are coming from a better understanding of the basic building blocks of life such as proteins. The current topic is *protein misfolding and aggregation*—if a protein folds into a shape other than the intended one, the result often produces inactive proteins with different properties, which can lead to such neurodegenerative diseases as Alzheimer's, Creutzfeldt-Jakob, bovine spongiform encephalopathy (Mad Cow), Huntington's and Parkinson's, cystic fibrosis and other amyloidosis.

Contents

Introduction

Understanding what can cause proteins to take a different folded form is an active research area with significant diseases. Recent experiments at Cambridge University (see [Related topics](#)) show that amyloids can be formed not only from the traditional proteins (such as lysozyme) given the right conditions (W62A) on lysozyme protein can cause amyloids compared to the wildtype (see sidebar) amyloids in urea solution due to the loss of hydrophobic interactions. What does all this give us?

Scientists do not yet know how this single residue can play a key role in the folding process and then shift to the surface for better understanding of the single residue effects, as well as the mechanism behind the aforementioned diseases related to protein misfolding and aggregation.

Proteomics: The protein economy
Commercial and academic motivations
So what are we modeling?
Equipping the laboratory
Running the model
What does all this give us?

! residue can play a key role in the folding process and then shift to the surface for better understanding of the single residue effects, as well as the mechanism behind the aforementioned diseases related to protein misfolding and aggregation.

The Blue Gene/L technology offers a powerful way to study these types of diseases, becoming an effective (and faster) way to model the effects of protein folding and misfolding.

So what are we modeling?

The [video](#) from which [Figure 1](#) was captured is a visualization of part of a sequence of a single mutation. Lysozyme is a protein that is part of the human immune system; when it encounters the cell wall of an invading bacterium and destroys it.

A single mutation, a different sequence in the DNA, causes the ribosome to use a different lysozyme molecule. The theory is that this different amino acid affects the shape that the differently shaped lysozyme molecule is differently effective in puncturing bacterial cell change, we may be able to design pharmaceuticals or other forms of therapy that will assist in recovering from bacterial disease.

As part of the work, we store the positions and velocities of every atom in one molecule approximately 10,000 water and urea molecules (this simulation is done in an 8 molar concentration experiments), in the computer's memory. There are many ways to model the forces between atoms: a *ball and spring* model for bonded forces with an inverse-square-law model for electrostatic forces and an attract/repel model for atoms that are near each other but not covalently bonded.

Contents

At each time step, we calculate the forces on each atom, then we update the velocities according to Newton's second law.

At each time step, we calculate the forces on each atom, then we update the velocities according to Newton's second law.

Introduction

Proteomics: The protein economy
At each time step (very small, on the order of femtoseconds), there are in principle hundreds of calculations. This large number of calculations (microseconds) to model motions of individual atoms means this approach has only recently become feasible. So what are we modeling?
"Destruction of long-range interaction: Equipping the laboratory"

Running the model

Equipping the laboratory
What does all this give us?

At IBM Watson Research Lab in Yorktown Heights, New York, we have 20 racks of BlueGene/L processors. Each microprocessor is attached to 512MB of memory. A single microprocessor connected to a 1Gbps Ethernet link. To standard IBM Power Systems machines. Related topics: compilers, and job-control software.

At each time step (very small, on the order of femtoseconds), there are in principle hundreds of calculations. This large number of calculations (microseconds) to model motions of individual atoms means this approach has only recently become feasible. So what are we modeling?
"Destruction of long-range interaction: Equipping the laboratory"

Equipping the laboratory

At IBM Watson Research Lab in Yorktown Heights, New York, we have 20 racks of BlueGene/L processors. Each microprocessor is attached to 512MB of memory. A single microprocessor connected to a 1Gbps Ethernet link. To standard IBM Power Systems machines. Related topics: compilers, and job-control software.

This lysozyme modeling work has used an average of four racks of BlueGene/L processors to produce an aggregate of more than 10 microseconds of molecular dynamics data. Periodically, the positions and velocities of all the atoms under simulation (part of this stream of information is shown in the [synthetic video](#) mentioned above). Whenever it is necessary to restart the simulation run, the positions and velocities can be reloaded. Restarting may be needed after a planned shutdown, after a hardware failure, or in order to replay a model event of scientific interest with a different time step granularity.

Running the model

The application is booted onto the Blue Gene/L nodes by a mechanism similar to MPICH—available, portable implementation of MPI, the message-passing interface; see [Related](#) the cluster provides a POSIX-file-system environment to the application. Data can be seen in the System (GPFS) file system for the application to read; when the application writes results for external use.

For time-series modeling applications such as this, it is normal to read the initial conditions and write periodic "snapshots" of the model state to the file system.

What does all this give us?

Contents

Introduction

The video is a glimpse into a world that scientists always need to compare what Proteomics: The protein economy misfolds in reality is still a dream; even them under an electron microscope or Commercial and academic motivations X-ray diffraction spectroscopy. However, So what are we modeling? might move.

Equipping the laboratory

Therefore, the current large-scale simulation critical changes involved in disease-related happen will push the envelope and address What does all this give us? generation of scientists to solve these

Research

Predicting the future

Downloadable resources

Predicting the future

Related topics

never before been visible. Of course, we don't know what the model shows with what they can see in the real world; "part of the "fixed" conformations mean that they are not even causing large numbers of lysozyme denaturation. Experimental techniques typically do not

offer a unique window to look into the details of protein misfoldings. Hopefully, the availability of the new state of the art in amyloidosis studies. I think this will be a problem in this new way as their primary

Actually, we would not be so bold as to attempt to divine tomorrow, but we would venture that the future of computing will continue to follow a development path (we use version L; the available Blue Gene/L processors per chip, 10Gbps Ethernet, and a host of other improvements). The cost of doing arithmetic and the cost of more and faster storage (both heavily associated with the data in this article) will most likely continue to fall—as they must, because there is several world-class scientists need to be doing, both for public research and for businesses to bring products to market.

The lysozyme model we describe only scratches a molecule off the surface of the new frontier. There are more than 50,000 proteins whose structures are catalogued in the public Protein Data Bank (link); there are millions of potential pharmaceutically useful compounds to be analyzed;

diseases known to be related to proteins and their defects. And we're not even consider that can benefit from modeling on this scale. Blue Gene's work has just begun.

Downloadable resources

 [PDF of this content](#)

Related topics

Contents

Introduction A roundup of Blue Gene/L-related components and resources include

Proteomics: The protein economy
The IBM [Blue Gene/P solution](#)

Commercial and academic motivations
The IBM [General Parallel File S](#)

So what are [Everything?](#) You could want to know

Equipping the laboratory
IBM [Redbooks on Blue Gene te](#)

Running the model
And a picture of the [Blue Gene](#)

What does all this give us?
The [RCSB Protein Data Bank](#) (PDB)

Predicting the future
Include such cool things as the [Mol](#)

Download [Source code for Figure 1](#) is from the

Related topics
Source data for [Figure 2](#) is from the

[Figure 3](#) is courtesy of the [MathMol library](#) hosted at New York University.

Source data for [Figure 4](#) is from the PDB, [Deoxy hemoglobin \(A-GLY-C:V1M,L29F,H5\)](#)

"[Destruction of long-range interactions by a single mutation in lysozyme](#)" (R. Zhou, N. Berne; Proc. Natl. Acad. Sci., 2007) gives more information about the modeling approach

"[Parallel implementation of the replica exchange molecular dynamics algorithm on E](#)
Rayshubski, J. W. Pitera, B. G. Fitch, R. Zhou, R. S. Germain; IEEE, 2006) explains so
used for the simulation.

[MPICH2](#) is the next stage of MPICH, the high-performance, widely portable (and free
Passing Interface (MPI) standard.

The [Argonne Leadership Computing Facility](#) has a collaborative program that provide

can be found at the [IBM Blue Gene proje](#)

ut the IBM [XL C/C++ Advanced Edition fo](#)
[ies](#)

one of the authors

chive for the study of biological macromolecules, nucleic acids, and complex assem
[the Month.](#)

[re 1.33 A structure of tetragonal hen egg](#)
[tructure of a B-DNA dodecamer: conform.](#)

computational science community.

"High-performance Linux clustering" is a two-part series providing background on hi Linux. [Part 1](#) (developerWorks, September 2005) covers HPC fundamentals, types o choosing a cluster configuration, and the role of Linux in HPC. [Part 2](#) (developerWork programming using MPI, covers cluster management and benchmarking, and shows open source software.

"[Port Fortran applications](#)" (developerWorks, April 2009) helps you overcome comm applications among various high performance computing systems.

Some of the tools integrated into the applications described in this article include th [for Blue Gene/L](#) and author Chris W [Custom Math Functions for High Performan](#)

Contents

— In the [developerWorks Linux zone](#),
Introduction tutorials.

Protein modeling: the protein economy
See [all Linux tips and Linux tutorial](#)

Commercial and academic motivations

So what are we modeling?

Equipping the laboratory

Running the model

☐ Subscribe me to comment notificatio
What does all this give us?

Predicting the future

[Custom Math Functions for High Performan](#)
re resources for Linux developers, and sc.

eloperWorks.

ments.

IBM Developer

About

Site Feedback & FAQ

Submit content

Report abuse

Third-party notice

Follow us





Select a language

English

中文

日本語

Русский

Português (Brasil)

Español

한국어

Code Patterns

Articles

Tutorials

Recipes

Open Source Projects

Videos

Newsletters

Events

Cities

Developer Answers

Contact

Privacy

Terms of use

Accessibility

Feedback

Cookie Preferences