



Please download and read the instructions before proceeding to the peer review

Advances in distributed computing with modern drug discovery

Journal:	<i>Expert Opinion On Drug Discovery</i>
Manuscript ID	EODC-2018-0116.R1
Manuscript Type:	Review (Invited)
Keywords:	Distributed computing, Drug discovery, High performance computing, Computational chemistry

SCHOLARONE™
Manuscripts

Advances in distributed computing with modern drug discovery

Antonio Jesús Banegas-Luna¹, Baldomero Imbernón¹, Antonio Llanes Castro¹, Alfonso Pérez-Garrido¹, José Pedro Cerón-Carrasco¹, Sandra Gesing², Ivan Merelli³, Daniele D'Agostino⁴, Horacio Pérez-Sánchez^{1,*}

¹*Bioinformatics and High Performance Computing Research Group (BIO-HPC), Universidad Católica de Murcia (UCAM), Spain.*

²*Center for Research Computing, University of Notre Dame, Notre Dame, IN, USA.*

³*Institute for Biomedical Technologies, National Research Council of Italy, Segrate (Milan), Italy.*

⁴*Institute for Applied Mathematics and Information Technologies "E. Magenes", National Research Council of Italy, Genoa, Italy.*

*Corresponding author: Horacio Pérez-Sánchez.: Tel: (0034) 968278819. E-mail: hperez@ucam.edu

ABSTRACT

Introduction: Computational chemistry dramatically accelerates the drug discovery process and high-performance computing (HPC) can be used to speed up the most expensive calculations. Supporting a local HPC infrastructure is both costly and time consuming and therefore many research groups are moving from in-house solutions to remote distributed computing platforms.

Areas covered: The authors focus on the use of distributed technologies, solutions and infrastructures to gain access to HPC capabilities, software tools and datasets to run the complex simulations required in computational drug discovery.

Expert opinion: The use of computational tools can decrease the time to market of new drugs. HPC has a crucial role in handling the complex algorithms and large volumes of data required to achieve specificity and avoid undesirable side-effects. Distributed computing environments have clear advantages over in-house solutions in terms of cost and sustainability. The use of infrastructures relying on virtualization reduces set-up costs. Distributed computing resources can be difficult to access, although web-based solutions are becoming increasingly available. There is a trade-off between cost effectiveness and accessibility in using on-demand computing resources rather than free/academic resources. Graphics processing unit computing, with its outstanding parallel computing power, is becoming increasingly important.

HIGHLIGHTS

- Virtual screening is a computational chemistry approach that is currently used to reduce the number of wet laboratory experiments required in drug discovery campaigns. However, it may require a considerable amount of computing power, the availability of which is an issue for many research groups and small- or medium-sized companies.
- Distributed computing is a cost-effective solution in drug discovery research and can be implemented in a number of ways. Grid computing is the most efficient way of

1
2
3 distributing the demand involved in calculations across a network of available computing
4 units. Cloud computing is deployed through virtual hardware resources and is a more flexible
5 approach than grid computing because it can be configured and scaled depending on the
6 complexity of the task.
7
8
9

10
11
12 ● Within distributed computing platforms, graphics processing unit computing
13 increasingly has a key role as a result of the parallel nature of the hardware, which increases
14 the throughput in most scientific computing tasks. This can be useful for the computationally
15 demanding tasks involved in drug discovery. However, not all computational chemistry codes
16 are amenable to efficient parallelization.
17
18
19

20
21
22 ● Cloud applications are now available to support the deployment of simulations to
23 non-expert users at an early stage. Web servers are available for virtual screening
24 calculations, which implement many different computational drug discovery techniques,
25 although most of these are for ligand-based methods because these are computationally
26 cheaper than structure-based methods.
27
28
29

30
31
32 ● Even when using distributed computing approaches, the size of the chemical
33 space is still too large (millions of compounds) to be covered by one single virtual screening
34 technique. Therefore a hierarchical virtual screening approach is often adopted: inexpensive
35 methods (e.g., similarity searching or pharmacophore modeling) are applied first to create
36 small, focused libraries, followed by more computationally expensive methods (e.g.,
37 molecular docking and molecular dynamics) to achieve the best results.
38
39
40

41
42
43 ● Structure-based virtual screening approaches, such as molecular docking or
44 molecular dynamics, generally require more computational resources than ligand-based
45 methods. However, ligand-based techniques are less accurate than structure-based methods
46 and therefore a trade-off is required when performing computations in distributed computing
47 platforms to obtain the best results according to the resources available.
48
49
50
51
52
53
54
55
56
57
58
59
60

KEYWORDS

Cloud computing; computational chemistry; distributed computing; drug discovery; grid computing; high performance computing; virtual screening.

1. Introduction

It is now widely accepted that the discovery of new drugs can be aided by the use of computational drug discovery (CDD) techniques and many approved drugs have reached and passed clinical trials with their help. New terms have been added to the vocabulary of researchers, including computer-aided drug design (CADD) [1], computer-aided molecular design (CAMD) [2] and computer-aided molecular modeling (CAMM).

Most drug discovery studies focus on enhancing specific parts of the processes involved in the development of new drugs, which can be divided into three phases: (1) a discovery phase, in which millions of candidate compounds are screened; (2) a selection phase, in which the candidate drugs undergo preclinical research; and (3) an assessment phase, in which the drug is developed and extensive clinical trials conducted. *In silico* solutions are usually carried out in the discovery phase, whereas screening was previously carried out in a laboratory over several years with significant economic costs.

Many computational techniques are now available to study the molecular interactions relevant to drug discovery, such as virtual screening (VS) [3], which is used to simulate a large number of interactions between proteins (also known as receptors and/or enzymes) and small molecule drug candidates (ligands). Docking software is usually tested on protein families where there are the most crystal structures and therefore it is common practice to test many proteins in parallel [4–6]. Side-effects caused by off-target bindings should be avoided and therefore the most promising compounds are usually tested against many other proteins. The docking

1
2
3 conformations that describe the interactions between each compound and the corresponding
4
5 target are optimized through molecular dynamics (MD) simulations to relax the system and
6
7 improve the accuracy with which the binding energy is calculated. Molecular dynamics is a
8
9 physics-based simulation method in which Newton's equations of motion are solved for each
10
11 atom of the system considering all the forces involved in their interactions [7]. Depending on the
12
13 number of atoms involved, it can be computationally demanding. Other cheaper techniques, such
14
15 as chemical similarity and calculating the proximity matrix, are also used in the drug discovery
16
17 process [8–10].
18
19
20

21
22 CDD methods can be used in either a predictive or a prospective way. In predictive CDD,
23
24 calculations are carried out to process a database of compounds and to anticipate which
25
26 compounds are most likely to be of interest before their characterization in the laboratory. In
27
28 prospective CDD, the experimental results obtained after screening compounds in the laboratory
29
30 are analyzed by CDD methods to try to understand why the compounds were selected. In both
31
32 these approaches, the complexity of the computations depends on the size of the database and/or
33
34 the accuracy of the methods used. The use of high performance computing (HPC) techniques is
35
36 now mandatory in the development of efficient and scalable tools for CDD. Many different HPC
37
38 approaches can be used, including graphics processing units (GPUs), in-house clusters of
39
40 computers, remote supercomputers and distributed computing infrastructures.
41
42
43

44
45 The creation and maintenance of a local computing facility capable of managing the huge
46
47 amount of data obtained from state-of-the-art acquisition instruments and simulation tools in
48
49 drug discovery projects is too expensive for many small- or medium-sized biotechnology
50
51 laboratories [11] and the use of remote computational services is a cost-effective solution.
52
53
54 Remote computing infrastructures provide researchers with the ability to adjust the
55
56
57
58
59
60

1
2
3 computational resources according to their actual requirements, whereas doubling the size of a
4 traditional cluster is expensive and can result in many idle resources during off-peak periods. In a
5 distributed environment, requests to double the amount of access to resources simply involve
6 paying for twice the capacity. Usability is the most crucial aspect in this scenario because end-
7 users require solutions that simplify their activities rather than making them more complex.
8
9

10 We review here the current trends and advances in the application of distributed
11 computing infrastructures to drug discovery. Section 2 reviews the most commonly used tools in
12 ligand- or structure-based virtual screening. The tools discussed here provide a diversity of
13 screening services, such as quantitative structure–activity relationship (QSAR) modeling,
14 docking and molecular docking, and MD. Section 3 introduces the currently available distributed
15 computing environments, including advances in the use of GPU-based HPC platforms. Section 4
16 gives several examples and case studies showing how large-scale distributed platforms based on
17 the grid and cloud computing paradigms have been used successfully in CDD projects. Section 5
18 presents our conclusions and discusses the future perspectives for drug discovery in combination
19 with advanced computational techniques.
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38

39 **2. Virtual screening methods**

40 This section reviews some of the methods—such as QSAR, docking and molecular
41 docking, and MD—that are run over distributed computing infrastructures and used routinely in
42 VS calculations.
43
44
45
46
47
48

49 **2.1. QSAR**

50 QSAR models can be defined as regression or classification models that relate several
51 variables, called descriptors, with bioactivity values to predict the activity of new compounds.
52
53
54
55
56
57
58
59
60

1
2
3 These descriptors codify several chemical features of compounds, including their
4
5 physicochemical properties and experimental measurements. QSAR models are developed using
6
7 different computational strategies [12]—such as statistical methods, which include multivariate
8
9 linear regression analysis (MLR) [13], principal components analysis (PCA) [14], partial least-
10
11 squares (PLS) analysis [15] and linear discriminant analysis (LDA) [16]—or artificial
12
13 intelligence approaches, such as extreme learning machines (ELMs) [17], neural networks (NN)
14
15 [18] and support vector machines (SVM) [19].
16
17
18

19 The applicability of QSAR models in drug discovery requires that the QSAR model has
20
21 good predictability and provides a physicochemical interpretation of the possible mechanism of
22
23 action. The development of QSAR models, their application to large datasets, the calculation of
24
25 quantum chemical descriptors and the use of advanced statistical methods or artificial
26
27 intelligence approaches demands large amounts of computing resources [20] and mandatory
28
29 access to HPC resources.
30
31
32

33 Tetko et al. [21] reviewed various public modeling environments for the development of
34
35 QSAR models, such as OpenMolGRID [22,23] and its successor Chemomentum [24] for grid
36
37 computing, and their application to aquatic toxicity [25], acute toxicity [26] and the discovery of
38
39 HIV-1 protease inhibitors [27]. Other QSAR tools running on distributed computing
40
41 environments include: Simplex Representation of Molecular Structure–Structural and
42
43 PhysiCochemical Interpretation [28], which can exploit ensemble predictions, classification
44
45 techniques and non-linear methods, although the model-building parameters are kept fixed;
46
47 ISIDA/QSPR [29,30], which also exploits ensemble predictions; and DTC Lab. Software Tools
48
49 [31], which includes a validation method for the selection of models, is also valid for non-linear
50
51 techniques.
52
53
54
55
56
57
58
59
60

2.2. *Molecular docking*

Ligand-based methods (e.g. QSAR, similarity searching, pharmacophore modeling and docking) represent worthwhile solutions in drug discovery. However, QSAR and similarity searching do not take into account knowledge about the binding site within the protein target and this can reduce the accuracy of the calculations. To overcome this issue, structure-based methods are the preferred choice when the 3D structure of the target is known, although they are usually computationally more expensive than ligand-based approaches. In such cases, it is studied how the activity of proteins may be altered when small ligands dock into the well-defined cavities of protein receptors. These ligands can act as molecular switches and control the activity of the protein. For proteins involved in a metabolic pathway related to a disease, artificial ligands can act as drugs [32]. As more metabolic pathways and their associated key proteins are identified, the search for artificial ligands has intensified as a method of improving the treatment of various diseases. The number of known protein structures continues to grow exponentially, a trend increasingly complemented by initiatives in structural genomics [33]. Molecular docking identifies the lead compounds that can bind to a target protein with high affinity [34]. This is achieved by calculating the optimum binding position for each molecule in a large database of potential targets using heuristics and then ranking the database with a scoring function according to the estimated affinity [35].

Docking methods have been investigated for many years and several compounds have been identified and developed as drugs [36]. Several docking methods are currently available—including AutoDock4 [37], AutoDock Vina [38], Glide [39] and Lead Finder [40]—each of which has different scoring functions and optimization methods. All of these methods use an atomistic representation of the protein and the ligand and allow the exploration of thousands of

1
2
3 possible binding positions and ligand conformations in the coupling process [41,42]. As a result,
4
5 the binding modes for many complexes are reliably predicted. However, unbiased comparative
6
7 evaluations of the estimations of affinity by molecular docking show little correlation between
8
9 the measured and predicted affinities over a wide range of receptor–ligand complexes [43],
10
11 suggesting that more advanced approaches are required to increase the accuracy of the total
12
13 binding energies. The use of explicit solvent molecules and the addition of dynamic effects to the
14
15 system may partially circumvent the limitations associated with classical docking simulations.
16
17

18 19 **2.3. Molecular dynamics**

20
21 The dynamic nature of real biological environments needs to be considered if meaningful
22
23 predictions are to be made. As it is not possible to apply the most demanding simulations directly
24
25 to large libraries of compounds, an efficient workflow should start by first using inexpensive
26
27 techniques, such as QSAR, and then proceed with more expensive molecular docking
28
29 techniques. The resulting best-ranked drug candidates can then be selected and implemented in
30
31 more advanced simulations, such as MD. This approach can include the combined effects of the
32
33 solvent and temperature in the evolution of the system over relatively long time trajectories—
34
35 that is, on the scale of nanoseconds to microseconds. MD has been used to rank a series of
36
37 biologically active ligands docked into the herceptin antibody, an efficient biological molecule
38
39 able to localize malignant cells in patients with breast cancer [44]. The MD simulations allowed
40
41 ligands that produced an early release from the binding site (during first 100 ns) to be discarded
42
43 because they were incompatible with a stable interaction.
44
45
46
47
48

49 MD has both advantages and disadvantages in drug discovery research. It has the
50
51 advantage that the stability of binding sites can be validated before cell-based assays are carried
52
53 out, but the production of MD trajectories requires large amounts of computational resources,
54
55
56
57
58
59

usually defined as the number of nanoseconds computed per day. Huge efforts have been devoted to speed up simulations and most of the currently available molecular dynamics codes include GPU-accelerated versions. Unfortunately, such codes require previous expertise from the user because molecules with chemical features that are not implemented in the force field parameters need to be optimized/implemented by the user. As an alternative, some web servers running under distributed computing infrastructures are available that may help non-expert users to perform short molecular dynamics simulations in a friendly framework, which might be used as a first proof of concept. A series of representative examples of online solutions are MDWeb [45], the Gromacs server as implemented in Haddock [46], Vienna-PTM [47], CABS-flex [48], Protein structure REFinement via MD [49], MoMA-LigPath [50] and MoSGrid [51]. Most of these solutions offer predefined protocols to guide the preparation of the structure (i.e. the experimental PDB file), which should be “cured” and then transformed to a specific format to run standard molecular dynamics simulations. These listed tools also allow the trajectories produced to be analyzed by monitoring the stability of the whole system using the root-mean-square deviation (RMSD) as well as the geometrical parameters of the binding site (the non-covalent interactions between the target protein and ligand).

One of the major drawbacks of classical MD calculations is the assessment of entropic terms, thermodynamic parameters that are required to determine reliable absolute free energies. There is no universal solution for the refinement of the molecular dynamics protocol and several methods have been proposed to better capture the ligand–protein problem in drug discovery, including free energy perturbations, umbrella sampling, the potential of the mean force and metadynamics. All these additional descriptors may be used to produce a more accurate prediction. However, their computational cost restricts their application to small and rigid

1
2
3 systems, which, in turn, prevents the implementation of molecular dynamics techniques in the
4
5 servers currently available. Further development is needed before molecular dynamics can be
6
7 systematically applied to complex systems because it requires sampling at large dimensions to
8
9 produce accurate values for the entropic and solvation contributions to the free energy [52].
10
11
12

13 14 **3. Distributed computing environments**

15
16 Biomedicine was one of the first areas of research [53] to move from the use of in-house
17
18 computing facilities or single supercomputing centers to distributed infrastructures, in particular
19
20 grid and cloud platforms. These infrastructures are based on the vision of providing services to
21
22 users through the sharing of capabilities and resources. The core idea is that any simulation can
23
24 be achieved using the concept of service: a workstation or a supercomputer represents a
25
26 computing service, but a database, a domain-specific application or an authentication mechanism
27
28 are also services. It is possible to find dozens of published definitions of these platforms, which
29
30 have been modified during their development. Foster [54]—one of the first developers of grid
31
32 computing technologies—defined a computational grid as an infrastructure of both hardware and
33
34 software that provides dependable, consistent, pervasive and inexpensive access to high-end
35
36 computational capabilities. Foster [54] suggested that the essence of grid computing can be
37
38 described by three concepts: (1) the coordination of resources that are not subject to centralized
39
40 control (see Figure 1); (2) access to, and the management of, resources using standard, open,
41
42 general purpose protocols and interfaces; and (3) the delivery of non-trivial qualities of service.
43
44 Many years after the spread of this technology began, we can now describe grid computing as a
45
46 wrapper with which to freely access remote multi-institutional resources, with either dedicated or
47
48 shared computational time. Grid computing paved the way for cloud computing.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Cloud computing has been defined by the US National Institute of Standards and
4
5 Technology as a model for enabling convenient, on-demand network access to a shared pool of
6
7 configurable computing resources (e.g., networks, servers, storage, applications and services)
8
9 that can be rapidly provisioned and released with minimal management effort or interaction with
10
11 the service provider. It is based on three service models (see Figure 2): (1) the software as a
12
13 service (e.g., Dropbox); (2) the platform as a service (e.g., the Google App Engine); and (3) the
14
15 infrastructure as a service (i.e., the possibility of managing virtual machines). In contrast with the
16
17 definition of grid computing, cloud computing can be defined as a way to use the internet to
18
19 deliver on-demand, scalable, pay-per-use services that rely on computational facilities hosted by
20
21 a single institution.
22
23
24
25

26 Cloud infrastructures can be used for business, commercial or research purposes, but the
27
28 concept is always based on access to on-demand resources on a pay-per-use basis. By contrast,
29
30 grid platforms are used for research by virtual organizations [55], which are dynamic sets of
31
32 individuals and/or institutions sharing a goal to be pursued using the grid resources without any
33
34 additional charge. The users of a cloud infrastructure are customers, but the users of a grid
35
36 platform are members of one or more virtual organizations.
37
38
39

40 The following sections analyze the relevant tools and infrastructures based on the use of
41
42 non-local computational resources, with a particular focus on architectures that have been used
43
44 for CDD projects.
45
46
47

48 **3.1. Grid computing**

49

50 Although initially exploited in the field of high-energy physics, the biomedical
51
52 community showed an immediate interest in grid computing—for drug discovery applications in
53
54 particular [56]—and many computational challenges have been met in this context. In addition to
55
56
57
58
59

1
2
3 large campaigns for the analysis of biosequences [57], grid computing has been widely exploited
4
5 in structural biology for large-scale virtual screenings of neglected diseases [58,59] and to
6
7 determine the molecular dynamics of huge biomolecular systems [60].
8
9

10 11 **3.1.1. Service-based grid computing**

12
13 Most of the grid infrastructures rely on middleware toolkits—that is, a software layer that
14
15 lies between the operating system and the applications. The first example of middleware was the
16
17 Globus Toolkit [61], followed by others such as the Advanced Resource Connector and the
18
19 Uniform Interface to Computing Resources. Many different platforms have been proposed on top
20
21 of these middleware components for drug discovery projects and these usually follow one of two
22
23 approaches.
24
25

26
27 The first approach relies on general purpose grid service environments coordinated by
28
29 virtual organizations. This paradigm can be exploited through both service and software
30
31 frameworks, such as the distributed infrastructure with remote agent control (DIRAC) [62], the
32
33 distributed analysis environment (DIANE) [63] and the gLite workload management system
34
35 (WMS) [64], which are able to provide high-level job submission services at the front end and
36
37 the possibility of interacting with heterogeneous systems at the back end (see Figure 1). The
38
39 second approach aims to standardize the creation, representation and sharing of the so-called
40
41 computational workflows that tie a number of software tools together into a single analysis. This
42
43 is built on application-oriented grid service infrastructures where the focus is on the software
44
45 executed over the clusters provided within a grid (see Figure 1). This is the case for myGrid [65],
46
47 a widely used infrastructure for composing pipelines in the field of bioinformatics and structural
48
49 biology using the Taverna Workbench or the Virtual Laboratory for e-Science (VL-e) [66],
50
51
52
53
54
55
56
57
58
59
60

1
2
3 which rely on a grid service infrastructure and use Nimrod-G semantics to compose complex
4
5 workflows.
6

7
8 Service-based grid computing is usually free, but the user needs to belong to an academic
9
10 institution. Association with a virtual organization is required for grids relying on such
11
12 organizations.
13

14 **3.1.2. Desktop-based grid computing**

15
16
17 The other grid approach that gained progressive success in the scientific community is
18
19 known as desktop grid or volunteer computing because it often relies on the general public to
20
21 donate resources. Desktop computing and, more generally, volunteer computing, is usually free
22
23 and only registration is required (see Figure 1).
24

25
26 This paradigm relies on middleware oriented to peer-to-peer architectures. This means
27
28 that there is no clear distinction between the grid entities owned by users (e.g., a workstation) to
29
30 access the infrastructure (e.g., a cluster shared via grid middleware) and every resource can act as
31
32 both a client and a server at any one time. Any user can bring resources into the grid because the
33
34 installation and maintenance of the software is intuitive, requiring no special expertise, which
35
36 enables a large number of donors to contribute to the pool of shared resources.
37
38

39
40 The most popular middleware for desktop grid applications is the Berkeley Open
41
42 Infrastructure for Network Computing (BOINC) [58], which was originally developed in the
43
44 context of the SETI@home project [67] to search for extraterrestrial intelligence. The main
45
46 feature of BOINC is represented by its simple API, which allows easy interaction with other
47
48 environments, and its great community support. For example, the EDGeS project [68], which
49
50 aims to create an infrastructure combining the advantages of service (the EGEE infrastructure)
51
52 and desktop grids, is using BOINC.
53
54
55
56
57
58
59
60

3.2. *Cloud computing*

Grid computing infrastructures are not completely satisfactory when running complex applications, in particular for industrial companies [69], because they do not provide a flexible and reliable environment. Most of the users submit batch jobs to remote clusters with little or no interactivity or the possibility of customizing the environment. The management of the storage of distributed data is complex, especially the administration of geographically dispersed databases. Misconfigured nodes on the grid are common, which results in jobs that fail continuously, emptying the grid queue, an effect known as shrink hole.

By contrast, cloud computing aims to provide the elements of a classical computational infrastructure (from a single workstation running the company website to a fully operational HPC cluster with high-speed network connections) as an on-demand service using virtual machines and virtual clusters, with development and execution frameworks and applications as services (see Figure 2). This paradigm addresses the key demands of creating an easily accessible and flexible environment (see Figure 2), able to support data processing and, in general, to provide everything required to perform complex analyses [70]. Cloud computing vendors such as Amazon [71] and Google [72] provide specialized environments to ease biomolecular data processing in the cloud.

Cloud computing allows the adoption of novel programming models—such as MapReduce and Spark, both of which rely on the Hadoop file system, an open source framework that enables the distributed processing of large datasets—which are exploited for sequence alignment [73–75], drug discovery [76–78] and many other applications [79]. In MapReduce, the input of a computation is split into independent chunks, which are then processed by the map tasks in a parallel manner. The results are then sorted and processed by the reduce tasks to provide the final output. By contrast, Spark works as a toolset for distributed programs and offers

1
2
3 a (deliberately) restricted form of distributed shared memory, facilitating the implementation of
4
5 iterative algorithms that visit their dataset multiple times.
6

7
8 Some initiatives support the execution of virtual machines on top of grid resources, such
9
10 as the Worker Nodes on Demand Services (WNoDes) framework [80]. This approach has been
11
12 exploited for macromolecular characterization and the estimation of free energy in protein–
13
14 protein docking [81].
15
16

17 18 **3.3. GPUs and distributed computing** 19

20 GPUs can be used to reduce the execution time of scientific applications [82,83] within
21
22 distributed computing scenarios. In particular, GPUs have been ranked as one of the platforms
23
24 with the highest projection for implementing algorithms that simulate complex scientific
25
26 problems. Since their first appearance, the development of GPUs has been marked by the world
27
28 of video games, which have reached high levels of popularity as more realism is achieved. In
29
30 2006, NVIDIA, the leader in the manufacturer of GPUs, made a breakthrough in the world of
31
32 HPC when they released the Compute Unified Device Architect (CUDA) development kit. This
33
34 architecture makes it possible to use GPUs for the development of scientific applications. Six of
35
36 the ten most powerful supercomputers in the world [84] currently have coprocessor (GPU or
37
38 similar vector-based devices) accelerators.
39
40
41
42

43 In the limited subset of problems for which enough effort can be invested to ensure that
44
45 specific drug discovery software supports GPU computing, these devices can greatly accelerate
46
47 calculations. Not all algorithms are good candidates for acceleration [85,86], but scientists are
48
49 now aware of the benefits that HPC can bring to their research, either by including these
50
51 hardware platforms in their studies or by directly designing and implementing specific software
52
53 that can be used on these platforms.
54
55
56
57
58
59

4. Examples of applications

This section discusses examples of drug discovery using different distributed computing architectures and ligand-based virtual screening (LBVS) servers.

4.1 Grid-based examples

One of the most interesting grid projects is Wide in Silico Docking on Malaria (WISDOM) [87], which has the goal of using the computing resources offered by the Enabling Grids for E-Science in Europe (EGEE) project (the largest grid project in Europe developed for the Large Hadron Collider at CERN in Geneva) to select drug-like molecules active on a biological target to fight malaria [88]. Malaria was chosen because tropical diseases often suffer from a lack of research due to the cost of bringing new drugs to market. The project developed *in silico* approaches to VS, mainly using resources from the EGEE project, but also from EUChinaGrid and TWGrid in Asia, EUMEDGRID in Africa, and OSG and Digital Ribbon in the USA [59]. Several of the drug-like molecules selected *in silico* have been confirmed by *in vitro* tests to be active inhibitors and the most promising molecules have been patented [89].

An infrastructure similar to that used in WISDOM was also used to perform a large VS project (see Figure 3) aimed at identifying new drugs for potential variants of the influenza A virus [58]. About 300,000 compounds selected from the ZINC database and an *ad hoc* chemical combinatorial library were screened against eight variants of neuraminidases predicted by homology modeling. The distributed analysis environment was used as the job dispatcher and a portal was developed to visualize the achieved outcome to make the docking conformations and scoring results available to biologists. Using thousands of CPU cores on a grid, a six-week long, high-throughput screening activity was accomplished, performing >100 CPU core years of calculations, producing around 600 gigabytes of results, and identifying about 100 compounds for further biological analysis and testing (see Figure 3).

1
2
3 Among the grid initiatives oriented toward structural biology analysis and relevant to
4 drug discovery, the European Model for Bioinformatics Research and Community Education
5 (Embrace) network of excellence [90] has been particularly important in the promotion,
6
7 development and implementation of standards for the interoperability of resources in the life
8 sciences community. Another example is the BIOINFOGRID project [91], which aims to exploit
9
10 the existing EGEE infrastructure for large-scale modeling and simulations of biological problems
11 and in which docking and molecular dynamics had a primary role. Many attempts have been
12 made to perform molecular dynamics on grid [92,93] and cloud [94,95] platforms, although the
13 only approach that gained some popularity was the Hamiltonian Replica Exchange method [96].
14
15
16
17
18
19
20
21
22
23

24 In addition to dedicated virtual organizations, several *ad hoc* grid systems have been
25 proposed for the execution of bioinformatics and drug discovery workflows [97]. One of the
26 most effective service-based environments for running workflows is Taverna [98]. Many projects
27 oriented toward the annotation, modeling and simulations of biological macromolecules, such as
28 myGrid [65] and the VL-e [66], rely on this infrastructure. Taverna was designed to combine
29 distributed grid services and/or local tools into complex analysis pipelines to tackle a wide range
30 of scientific research [99]. Once constructed, the workflows become reusable resources. Many of
31 these workflows are oriented toward structural biology and drug discovery, in particular the VS
32 and identification of protein binding sites, by leveraging Taverna's technology for data storage,
33 workflow enactment, change event notification, resource discovery and provenance management
34
35
36
37
38
39
40
41
42
43
44
45
46
47 [100].
48

49 The VL-e environment, by leveraging existing grid technologies, enables molecular
50 modeling for drug design using geographically distributed resources. The concept involves
51 screening millions of compounds in a chemical database against a protein target to identify those
52
53
54
55
56
57
58
59
60

1
2
3 molecules with potential use in drug design. The VL-e uses the Nimrod-G parameter
4
5 specification language to transform the existing molecular docking application into a parameter
6
7 sweep application suitable for execution on distributed systems. New tools have been developed
8
9 to enable access to ligand records/molecules in chemical databases from remote resources. The
10
11 Nimrod-G resource broker, along with the chemical database data broker, is used for scheduling
12
13 and the on-demand processing of docking jobs on the grid resources. The results show the ease
14
15 of use and power of the Nimrod-G language and VL-e tools for drug discovery on grid
16
17 computing platforms [101].
18
19

20
21 Moving toward the desktop grid paradigm, one important example is represented by the
22
23 Drug Discovery Grid (DDGrid) [102] project. This focuses on providing drug screening services,
24
25 such as building docking processes for the virtual high-throughput screening of the avian
26
27 influenza virus [103]. The project relies on a grid environment and the grid computation
28
29 resources of the Grid@Asia project. DDGrid leverages BOINC with grid technologies to harness
30
31 the power of clusters and supercomputing systems owned by different organizations in China
32
33 and South Korea.
34
35

36
37 Rosetta@home, a distributed computing project for the prediction of protein structures,
38
39 also relies on the BOINC platform and aims to model protein-protein docking and protein
40
41 structures. Rosetta@home has been used to validate proteins that neutralize influenza [104,105]
42
43 and enzymes that break down gluten for the treatment of celiac disease, all of which are moving
44
45 through animal trials and into clinical trials [106,107]. Rosetta@home has also been applied to
46
47 research on malaria, Alzheimer's disease and other pathologies [108].
48
49

50
51 In addition to disease-related research, the Rosetta@home network serves as a testing
52
53 framework for new methods in structural biology. Such methods are then used in other Rosetta-
54
55

1
2
3 based applications, such as RosettaDock [109], RosettaDesign and Robetta [110]. Rosetta@home
4 consistently ranks as one of the foremost docking predictors and is one of the best predictors of
5 tertiary structure currently available [111]. Using hundreds of thousands of cellular telephones and
6 other mobile devices, Rosetta@home was used to validate peptide macrocycles designed with rigid
7 structures, which could be useful as peptide therapeutics [112,113].
8
9

10
11 The conformational states from Rosetta's software can be used to initialize a Markov state
12 model as starting points for simulations using Folding@home [114], which is a distributed
13 computing platform for studying protein folding and other types of problems that can be solved
14 with molecular dynamics. As Rosetta only tries to predict the final folded state, and not how
15 folding proceeds, Rosetta@home and Folding@home are complementary and address very
16 different molecular questions. Folding@home has been recognized as the most powerful
17 distributed computing network in the world [115] and, in contrast to other similar projects, it does
18 not rely on BOINC, but on a specific networking library called Cosm [116].
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 Large-scale executions of drug discovery applications on volunteer grids are no longer
34 simply research projects, but an effective way to run *in silico* simulations. The main issue is the
35 ratio of volunteer resources needed to achieve the computing power of a grid or cloud node,
36 which is, on average, greater than three [117,118].
37
38
39
40
41

42 ***4.2 Cloud-based examples***

43
44 Many examples of cloud computing rely on the platform as a service (PaaS) model,
45 where applications are built using higher level platforms/frameworks. Many MapReduce-
46 inspired frameworks developed by the Apache Software Foundation are currently available to
47 manage and process large amounts of high-throughput omic data. For example, the Cancer
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Genome Atlas project made use of Hadoop to split genome data into chunks distributed over the
4 cluster for parallel processing [119,120].
5
6

7
8 The Collaborative Genomic Data Model (CGDM) [121] uses Hadoop to boost the
9 querying speed for the main classes of query on genomic databases. MetaSparks [121] uses
10 Spark to recruit large-scale metagenomics reads to reference genomes, achieving better
11 scalability and sensitivity than programs based on a single machine [122], a principle that has
12 also been applied to drug discovery [123].
13
14
15
16
17
18

19 **4.3 GPU-based examples**

20
21
22 GPUs have recently started to be used outside their traditional role as a graphical
23 component of computers. This new application is known as general purpose computing on
24 graphics processing units (GPGPU).
25
26
27

28
29 An example of GPU sharing within a distributed infrastructure is GPUGRID.net [124], a
30 distributed computing infrastructure devoted to biomedical research. It consists of many GPUs
31 joined together to deliver high performance all-atom biomolecular simulations. The simulations
32 aim to develop better drugs by determining the mechanisms of drug resistance in cancers,
33 modeling HIV maturation and investigating the features of neurologically important proteins.
34
35
36
37
38
39

40
41 The possibility of accessing GPUs on distributed facilities is particularly appealing
42 because the most popular manufacturer, NVIDIA, markets itself as the best way to accelerate this
43 kind of application. They have reported that molecular dynamics applications, such as ACEMD
44 [125] and GROMACS [126], traditionally CPU-based, achieve increases in speed of up to eight
45 times with the incorporation of their massive parallel platforms.
46
47
48
49
50

51
52 Researchers have carried out comparisons showing the reduced times achieved by
53 performing tasks in these new platforms relative to traditional computation. Marin et al. [127]
54
55
56
57
58
59

1
2
3 compared the performance of a parallelization on multicore systems against a GPU and showed
4 that the GPU implementation was up to 33 times faster. Analyses by Chiappori et al. [94]
5
6 showed that these simulation techniques are time consuming and therefore parallel computing
7
8 and GPU computations are required to reduce computation times.
9
10

11
12 Interesting papers have been published in this area of research—for example, Hung et al.
13 [128] reviewed the two main CADD-based approaches (structure- and ligand-based drug design)
14 and concluded that both multiple computers and GPGPU approaches can significantly improve
15 the CADD performance. Vogt et al. [129] reviewed the current trends in method development,
16 including the implementation of GPUs.
17
18
19
20
21
22

23
24 Ma et al. [130] used GPUs to accelerate the chemical similarity calculation that plays a
25 major part in VS. Malhat et al. [131] saved between 76 and 99% of their computation time by
26 using GPUs in their implementation of the Ward algorithm to group similar chemical
27 compounds. Lo et al. [132] used CUDA to accelerate the prediction of protein–ligand binding
28 regions using geometrical features.
29
30
31
32
33
34

35
36 Docking applications based on the scoring function, such as MetaDock [133], use a
37 metaheuristic scheme to generate a large number of heuristic strategies for VS. MetaDock is
38 designed to take full advantage of parallel and heterogeneous architectures, including
39 multiprocessor chipsets and NVIDIA GPUs. Figure 4 shows how the MetaDock data distribution
40 model works in GPUs. In this schema, the receptor molecule is stored in the shared memory of
41 the multiprocessors to make better use of the memory accesses. The drug candidates are grouped
42 in blocks of threads and are able to share the data from the receptor between the threads of the
43 same multiprocessor. The application uses a computational molecular coupling methodology that
44 seeks to predict the non-covalent binding of molecules or, more often, a macromolecule
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 (receptor) and a small molecule (ligand). The aim is to predict the bound conformations and
4
5 affinity of the union—that is, the strength of association—which is usually measured by a
6
7 scoring function [134,135].
8
9

10 Another approach that takes advantage of parallel programming is the migration of tools
11
12 into this new parallel paradigm. McIntosh-Smith et al. [136] used this new GPU programming
13
14 paradigm to increase the overall performance of the drug screening process, porting their tools to
15
16 OpenCL. They describe the BUDE (Bristol University Docking Engine), which is a structure-
17
18 based virtual screening (SBVS) engine, and present it as one of the first applications migrated to
19
20 modern hardware [136]. Fang et al. [137] describe GeauxDock, a new molecular docking
21
22 program migrated to run over parallel platforms. Some of these research groups are also involved
23
24 in the review of new programs developed to take advantage of these parallel platforms [138].
25
26 Krige et al. [139] ported a commercial application (blazeV10) of VS to these platforms and
27
28 compared an OpenCL version on a broad range of devices. Harvey et al. [140] implemented
29
30 ACEMD, a production-class biomolecular dynamics simulation program specifically designed
31
32 for GPUs. Sukhwani et al. [141] accelerated a production mapping software program using
33
34 NVIDIA GPUs.
35
36
37
38
39

40 In addition to migrating tools to a distributed infrastructure to take advantage of the
41
42 massive parallel platforms, another approach is to implement new programs to perform some
43
44 tasks. For example, McArt et al. [142] implemented a new program to perform connectivity
45
46 mapping, which is a computational technique dedicated to drug repurposing around differential
47
48 gene expression analysis. They drastically reduced the computational times; previous
49
50 implementations took up to seven days, whereas using a GPU reduced the time required to just
51
52 ten minutes.
53
54
55
56
57
58
59

1
2
3 Studies in parallel algorithms that have led to new advances in this area include the work
4
5 of Hu et al. [143], who developed a new GPU-based analytical method focusing on risk epistasis
6
7 in a genome-wide association study of complex traits.
8
9

10 11 **4.4. Applications of virtual screening web servers**

12
13 HPC has been extensively applied in virtual screening. When full information about the
14
15 target protein is available, SBVS methods are preferred and several servers are available to
16
17 perform docking calculations. DOCK Blaster [144], Haddock [46], iScreen [145], SwissDock
18
19 [146] and VSDocker [147] are examples of this type of server. Achilles [148] is a web tool
20
21 implementing the blind docking approach. A distinguishing feature of Achilles is that it provides
22
23 detailed reports, including the most likely ligand structure clusters, after processing the whole
24
25 surface of the chosen protein.
26
27

28
29 SBVS calculations are usually very accurate, but the required information from the target
30
31 protein is not always available. In such cases, LBVS is used as an alternative approach. A
32
33 number of different servers implementing different approaches to LBVS are available on the
34
35 web. As an example, SwissSimilarity [149] performs both 2D and 3D similarity searching,
36
37 whereas HybridSim-VS [150] and USR-VS [151] focus on 3D similarity, including geometrical
38
39 information about molecules. Other servers, such as iDrug [152], ChemMapper [153] and
40
41 ZINCPharmer [154], use pharmacophore modeling to assess similarity. LBVS servers can
42
43 efficiently scan millions or billions of compounds in a short time and are often a cost-effective
44
45 solution to screening.
46
47
48

49
50 These examples have focused on the core of VS calculations. However, HPC can also
51
52 help in the early stages of the screening process. For example, the prediction of molecular
53
54 descriptors or fingerprints (e.g. molecular weight, number of rotatable bonds, number of
55
56
57
58
59
60

1
2
3 acceptor/donor hydrogens) is a sensitive task that benefits from high computing power.

4
5 BioTriangle [155] is a complete toolkit for the characterization of complex biological molecules
6
7 and their interactions. This tool is not only useful in the field of cheminformatics, but also in
8
9 bioinformatics. ChemDes is another web tool for fingerprinting and molecular descriptor
10
11 calculations. It hides the complexity of many other open source packages from the users. An
12
13 extensive compilation of tools and databases for drug discovery is maintained by the Swiss
14
15 Institute of Bioinformatics [156].
16
17
18
19
20

21 **5. Conclusions**

22
23 Distributed computing infrastructures, in particular those relying on GPUs, are very
24
25 important in drug discovery programs. Although Big Pharma can rely on proprietary
26
27 infrastructures to perform their analyses in-house, small- or medium-sized biotechnology
28
29 laboratories need to exploit distributed infrastructures to achieve a flexible and cost-effective
30
31 platform to perform their simulations. With distributed computing platforms, all players can
32
33 combine the use of third-party software made available through a web interface with their own
34
35 tools to exploit the full potential of this approach. There is a wide range of online VS servers
36
37 providing services for many of the steps in the drug discovery process, including screening and
38
39 fingerprint calculations.
40
41
42

43
44 The use of GPUs as computing platforms for drug discovery processes will further
45
46 accelerate computations, as shown by the fact that commercial tools now include GPUs in their
47
48 implementation to take advantage of their enhanced performance. Working with heterogeneous
49
50 computation resources can help to improve the performance of applications. The main issues in
51
52 these technological approaches are related to the cost-effective exploitation of the available
53
54 computing capabilities because each GPU model has different features depending on its family
55
56
57
58
59

1
2
3 and generation and variations in the number of cores and performance. There is a need to
4
5 develop applications with a load-balancing technique to correctly account for the characteristics
6
7 of each device when distributing the workload through a distributed computing environment.
8
9
10 Grid and cloud science gateways are an effective solution to providing user-friendly access to
11
12 computational power and tools, while solving security issues and the problem of moving data
13
14 between centers. These infrastructures are typically free for scientific purposes.
15
16
17

18 **6. Expert opinion**

19
20 Rapid developments in computational capabilities, combined with the explosion of
21
22 research into personalized medicine, is currently leading the third wave of drug discovery. This
23
24 new approach of personalized drug discovery should greatly improve the therapeutic effects of
25
26 drugs while reducing their side-effects. The screening and validation of functional gene and
27
28 protein targets in the early stages of drug discovery, and the development and optimization of
29
30 selected molecules, requires huge amounts of computational power, which is difficult to buy and
31
32 operate for small- and medium-sized biotechnology laboratories and academic research groups.
33
34 In addition, perhaps surprisingly, large pharmaceutical companies are increasingly outsourcing
35
36 research and computing activities to cut costs and to access state-of-the-art knowledge and
37
38 technologies. We therefore see leading scientific institutions and computational chemistry
39
40 companies shifting their business model from releasing commercial packages to providing
41
42 distributed/cloud support for drug discovery research, while pharmaceutical companies exploit
43
44 cloud services for drug development. As an example, in 2014, the Novartis Institutes for
45
46 Biomedical Research used Amazon infrastructures to build a platform for the virtual screening of
47
48 compounds against a common cancer target [157], leveraging the power and agility of cloud
49
50 computing to conduct 39 years of computing in just 11 hours, creating an on-demand computing
51
52
53
54
55
56
57
58
59
60

1
2
3 system that would cost an estimated US\$44 million to build for a cost of only US\$5000. It is
4
5 now almost impossible for researchers to enter the world of drug discovery without using these
6
7 platforms.
8
9

10 These developments lead to the key message of this paper: recent advances in distributed
11
12 computing technologies and infrastructures mean that they are now essential in the field of drug
13
14 discovery. This was not true before the emergence of the grid computing paradigm about 30
15
16 years ago and, later, the cloud computing paradigm, but is currently very clear. The possibility of
17
18 exploiting these paradigms presents great opportunities, as demonstrated by the large simulations
19
20 conducted on grids and clouds by both pharmaceutical companies and public institutions, but
21
22 also presents some weaknesses in terms of the effectiveness, accuracy and usability of software
23
24 and hardware in distributed environments.
25
26
27

28 Effectiveness and accuracy mean the use of the correct methods in QSAR, VS, lead
29
30 optimization and molecular dynamics. Therefore, after the development of new algorithms,
31
32 usually by academic researchers, it is necessary to support the software through its full
33
34 development for commercial use. The most common business model is to create spin-off
35
36 companies to consolidate the development of these approaches, closing the gap between
37
38 academia and the market, with sponsorship from pharmaceutical companies. These companies
39
40 should also support users in applying these new approaches in the correct way, preventing
41
42 failures in handling each specific tool. Some successful examples are available in the field of
43
44 drug discovery, such as AMBER, one of the most popular software packages for molecular
45
46 dynamics. Usability refers to the development of environments that lower the barriers to
47
48 applying distributed and parallel computing infrastructures for drug discovery and offer support
49
50 in handling failures.
51
52
53
54
55
56
57
58
59
60

1
2
3 Despite some open issues, drug discovery is moving towards distributed infrastructures
4 because they provide high performance and cost-effective access to software tools, data and
5 infrastructures. From the computational point of view, *in silico* analyses that were, until recently,
6 impossible, are becoming increasingly feasible by applying the most suitable distributed
7 computing infrastructure. This represents an incredible potential for this field of research, but
8 more effort is required to develop and automate the functionalities that are crucial in enabling
9 agile and flexible predictive modeling and simulation protocols. Workflow management systems
10 can aid in many of these challenges, but the currently available systems are not suitable for users
11 unfamiliar with ICT systems.
12
13
14
15
16
17
18
19
20
21
22
23

24 By considering the cost effectiveness of the different distributed infrastructures available
25 to researchers for drug discovery, a recent review of computational models in computational
26 biology evaluated the performance of a single application exploiting grid and cloud computing
27 [70]. This study reported comparable results in terms of execution time, but also demonstrated
28 that grid platforms have overheads due to failures caused by server misconfigurations and
29 waiting times in cluster queues. Cloud infrastructures should therefore be preferred over grid
30 approaches if the budget is sufficient, whereas if cost is a major issue, such as for neglected
31 diseases, grid platforms may still be a valid solution.
32
33
34
35
36
37
38
39
40
41

42 Cloud platforms have cost advantages compared with buying local facilities because
43 dynamic access to resources is less expensive than having a private infrastructure. In particular,
44 if spot instances—cheap resources that can be turned off by the vendor when there is a high on-
45 demand computational load—can be used, which may occur in virtual screening but not in
46 molecular dynamics simulations, then cloud computing could be cheaper than buying computing
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 power from a high performance computational facility. This is particularly true if the cloud
4 servers are equipped with high-end GPU devices and the applied cost model is appropriate.
5
6

7
8 Grid infrastructures and supercomputing centers require the pre-installation and tuning of
9 software, whereas cloud solutions can exploit the pre-configured virtual machines released by
10 vendors (e.g., Amazon and Google) or specialized companies (e.g., Eagle Genomics or Cloud
11 Pharmaceuticals). However, it is important to highlight the fact that containers specialized for
12 drug discovery analysis—for example, those relying on Docker [158]—can be used to rapidly
13 create custom environments on all these platforms, reducing the time needed to set up and
14 configure the software [159].
15
16
17
18
19
20
21
22
23

24 Future work in this area should be directed toward simplifying analyses and improving
25 cooperation among researchers by providing input data, sharing the parameterization settings and
26 releasing non-sensitive results to the scientific community. Open access to data and the
27 reproducibility of methods have gained much attention in recent years and this is particularly true
28 for approaches to developing new drugs. Tools are being developed to support researchers in this
29 regard, but they require improvements and the trust of the scientific community. Workflow
30 management systems and virtualization platforms are currently available to help improve the
31 reproducibility and sharing of results and methods, although this is seldom considered before a
32 publication or a patent is accepted as a result of the costly and long-running analyses required to
33 achieve significant results and due to privacy/security issues. We postulate that open sharing is
34 beneficial to the research community and research is more efficient when sharing is possible
35 while protecting sensitive data.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 A particular area of interest that, in our opinion, represents the next revolution in drug
52 discovery is related to deep learning approaches [160], both for selecting compounds and for
53
54
55
56
57
58
59
60

1
2
3 their optimization. With the support of NVIDIA, some tools have already been released to
4
5 improve the quantum mechanics energy function and to produce computationally fast and
6
7 accurate molecular energy surfaces, geometries and forces. As a result of the complexity of these
8
9 approaches, the use of HPC in a distributed infrastructure will be unavoidable.
10
11
12
13

14 **Funding:**

15
16 This work was funded by a grant from the Spanish Ministry of Economy and Competitiveness
17 (CTQ2017-87974-R) and by the Spanish MEC and European Commission FEDER (TIN2016-
18 78799-P, AEI/FEDER, UE). This research was partially supported by the supercomputing
19
20 infrastructure of Poznan Supercomputing Center, the e-infrastructure program of the Research
21
22 Council of Norway, the supercomputer center of UiT – the Arctic University of Norway and by
23
24 the computing facilities of Extremadura Research Centre for Advanced Technologies
25 (CETA–CIEMAT), funded by the European Regional Development Fund (ERDF).
26
27 CETA–CIEMAT belongs to CIEMAT and the Government of Spain. The authors also
28
29 acknowledge the computing resources and technical support provided by the Plataforma
30
31 Andaluza de Bioinformática of the University of Málaga. Powered@NLHPC research was
32
33 partially supported by the supercomputing infrastructure of the NLHPC (ECM-02).
34
35
36
37
38
39

40 **Declaration of Interest:**

41 The authors have no other relevant affiliations or financial involvement with any organization or
42
43 entity with a financial interest in or financial conflict with the subject matter or materials
44
45 discussed in the manuscript apart from those disclosed.
46
47

48 **Reviewer Disclosures:**

49 One referee is an employee of Merck & Co.
50
51
52
53

54 **Abbreviations**

55
56 CADD = Computer-Aided Drug Design
57
58
59
60

1
2
3 CAMD = Computer-Aided Molecular Design

4
5
6 CAMM = Computer-Aided Molecular Modeling

7
8 CDD = Computational Drug Discovery

9
10
11 ELM = Extreme Learning Machines

12
13 GPGPU = General Purpose Computing on Graphics Processing Units

14
15
16 GPU = Graphics Processing Unit

17
18 HPC = High Performance Computing

19
20
21 LBVS = Ligand-Based Virtual Screening

22
23
24 LDA = Linear Discriminant Analysis

25
26 MD = Molecular Dynamics

27
28
29 MLR = Multi-Linear Regression

30
31 NN = Neural Networks

32
33
34 PCA = Principal Components Analysis

35
36 PLS = Partial Least-Squares

37
38
39 QSAR = Quantitative Structure–Activity Relationships

40
41
42 SBVS = Structure-Based Virtual Screening

43
44 SVM = Support Vector Machines

45
46
47 VS = Virtual Screening

48 49 50 **REFERENCES**

51
52 1. Naylor CB, Richards WG. Computer-aided Drug Design [thesis]. Oxford (UK): University of
53 Oxford; 1984.
54
55
56
57
58
59
60

- 1
2
3 2. McCammon JA. Computer-aided molecular design. *Science*. 1987; 238(4826):486-91.
4
5
- 6
7 3. Shoichet BK. Virtual screening of chemical libraries. *Nature*. 2004; 432(7019):862-865.
8
9
- 10
11 4. Morris GM, Huey R, Lindstrom W, et al. AutoDock4 and AutoDockTools4: Automated
12 docking with selective receptor flexibility. *J Comput Chem*. 2009; 30(16):2785-91.
13
14
- 15
16 5. Friesner RA, Murphy RB, Repasky MP, et al. Extra precision glide: docking and scoring
17 incorporating a model of hydrophobic enclosure for protein-ligand complexes. *J Med Chem*.
18 2006; 49(21):6177-96.
19
20
21
- 22
23 6. Merelli I, Cozzi P, D'Agostino D, et al. Image-based surface matching algorithm oriented to
24 structural biology. *IEEE/ACM Trans Comput Biol Bioinform*. 2011 Jul-Aug;8(4):1004-16.
25
26
- 27
28 7. Merelli I, Morra G, Milanese L. Evaluation of a grid based molecular dynamics approach for
29 polypeptide simulations. *IEEE Trans Nanobioscience*. 2007; 6(3):229-34.
30
31
32
- 33
34 8. Sheridan RP, Kearsley SK. Why do we need so many chemical similarity search methods?
35 *Drug Discov Today*. 2002; 7(17):903-911.
36
37
- 38
39 9. Rudin M, Weissleder R. Molecular imaging in drug discovery and development. *Nat Rev*
40 *Drug Discov*. 2003; 2(2):123-131.
41
42
43
- 44
45 10. Wang L, Wu YJ, Deng YQ, et al. Accurate and Reliable Prediction of Relative Ligand
46 Binding Potency in Prospective Drug Discovery by Way of a Modern Free-Energy Calculation
47 Protocol and Force Field. *J Am Chem Soc*. 2015; 137(7):2695-2703.
48
49
50
- 51
52 11. D'Agostino D, Clematis A, Quarati A, et al. Cloud infrastructures for in silico drug
53 discovery: economic and practical aspects. *Biomed Res Int*. 2013; 2013:138012.
54
55
56
57
58
59

- 1
2
3 12. Roy K, Kar S, Das RN. A Primer on QSAR/QSPR Modelling. 1st ed.: Springer, Cham; 2015.
4 Chapter 2, Statistical Methods in QSAR/QSPR; p. 37.
5
6
7
8 13. Pérez-Garrido A, Girón-Rodríguez F, Helguera AM, et al. Topological structural alerts
9 modulations of mammalian cell mutagenicity for halogenated derivatives. SAR and QSAR in
10 Environmental Research. 2014; 25(1):17-33.
11
12
13
14
15 14. Yoo C, Shahlaei M. The applications of PCA in QSAR studies: A case study on CCR5
16 antagonists. Chem Biol Drug Des. 2018; 91(1):137-152.
17
18
19
20 15. Luco JM, Ferretti FH. QSAR Based on Multiple Linear Regression and PLS Methods for the
21 Anti-HIV Activity of a Large Group of HEPT Derivatives. J Chem Inf Comput Sci. 1997;
22 37(2):392-401.
23
24
25
26
27 16. Pérez-Garrido A, Helguera AM, Rodríguez FG, et al. QSAR models to predict mutagenicity
28 of acrylates, methacrylates and alpha,beta-unsaturated carbonyl compounds. Dent Mater. 2010;
29 26(5):397-415.
30
31
32
33
34 17. Pérez-Garrido A, Girón-Rodríguez F, Bueno-Crespo A, et al. Fuzzy clustering as rational
35 partition method for QSAR. Chemometr Intell Lab Syst. 2017; 166:1-6.
36
37
38
39 18. Ghasemi F, Mehridehnavi A, Pérez-Garrido A, et al. Neural network and deep-learning
40 algorithms used in QSAR studies: merits and drawbacks. Drug Discov Today. 2018; S1459-
41 6446(17):30476-2.
42
43
44
45
46 19. Darnag R, Mazouz ELM, Schmitzer A, Villemin D, Jarid A, Cherqaoui D. Support vector
47 machines: Development of QSAR models for predicting anti-HIV-1 activity of TIBO
48 derivatives. Eur J Med Chem. 2010; 45(4):1590-1597.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 20. Darvas F, Papp Á, Bágyi I, et al. OpenMolGRID, a GRID Based System for Solving Large-
4 Scale Drug Design Problems. In: Dikaiakos M.D, editor. Grid Computing. Berlin, Heidelberg:
5 Springer. 2004; 69-76.
6
7

8
9
10 21. Tetko IV, Maran U, Tropsha A. Public (Q)SAR services, integrated modeling environments,
11 and repositories on the web: state of the art and perspectives for future development. Mol Inf.
12 2017; 35:1600082.
13

14
15 *** A review of public modeling environments for QSAR development. The authors explain**
16 **the transition from the old-fashioned models to the new online approach.**
17
18

19
20 22. Sild S, Maran U, Lomaka A, et al. Open computing grid for molecular science and
21 engineering. J Chem Inf Model. 2006; 46:953-959.
22
23

24
25 23. Sild S, Maran U, Romberg M, et al. OpenMolGRID: Using Automated Workflows in GRID
26 Computing Environment. In: Sloat PMA, Hoekstra AG, Priol T, Reinefeld A, Bubak M, editors.
27 Advances in Grid Computing - EGC 2005. Berlin, Heidelberg: Springer. 2005; 464-473.
28
29

30
31 24. Schuller B, Demuth B, Mix H, et al. Chemomentum - UNICORE 6 Based Infrastructure for
32 Complex Applications in Science and Technology. In: Bougé L, Forsell M, Träff JL, Streit A,
33 Ziegler W, Alexander M, Childs S, editors. Proceedings of the Euro-Par 2007 Workshops
34 Parallel Processing; 2007; Rennes, France: Springer-Verlag. 2008;82-93.
35
36
37
38

39
40 25. Mazzatorta P, Smiesko M, Lo Piparo E, et al. QSAR Model for Predicting Pesticide Aquatic
41 Toxicity. J Chem Inf Model. 2005; 45:1767-1774.
42
43
44

45
46 26. Maran U, Sild S, Mazzatorta P, et al. Grid Computing for the Estimation of Toxicity: Acute
47 Toxicity on Fathead Minnow (*Pimephales promelas*). In: Dubitzky W, Schuster A, Sloat PMA,
48 Schroeder M, Romberg M, editors. Distributed, High-Performance and Grid Computing in
49 Computational Biology. Berlin, Heidelberg: Springer; 2007. p. 60-74.
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 27. Maran U, Sild S, Kahn I, et al. Mining of the chemical information in GRID environment.
4 Future Gen Comput Syst. 2007; 23:76–83.
5
6
7
8 28. Polishchuk P, Tinkov O, Khristova T, et al. Structural and Physico-Chemical Interpretation
9 (SPCI) of QSAR Models and Its Comparison with Matched Molecular Pair Analysis. J Chem Inf
10 Model. 2016; 56(8):1455-1469.
11
12
13
14
15 29. Varnek A, Fourches D, Horvath D, et al. ISIDA - platform for virtual screening based on
16 fragment and pharmacophoric descriptors. Curr Comput Aided Drug Des. 2008; 4:191-198.
17
18
19
20 30. Ruggiu F, Marcou G, Varnek A, et al. ISIDA Property-Labelled Fragment Descriptors. Mol
21 Inform. 2010; 29(12):855-868.
22
23
24
25 31. Ambure P, Aher RB, Gajewicz A, et al. NanoBRIDGES software: Open access tools to
26 perform QSAR and nano-QSAR modelling. Chemom Intell Lab Syst. 2015; 147:1-13.
27
28
29
30 32. Ekins S, Mestres J, Testa B. In silico pharmacology for drug discovery: Methods for virtual
31 ligand screening and profiling. Br J Pharmacol. 2007; 152(1):9–20.
32
33
34
35 33. Kanehisa M, Furumichi M, Tanabe M, et al. KEGG: New perspectives on genomes,
36 pathways, diseases and drugs. Nucleic Acids Res. 2017; 45(D1):D353–361.
37
38
39
40 34. Ferreira L, dos Santos RN, Oliva G, et al. Molecular docking and structure-based drug design
41 strategies. Molecules. 2015; 20:13384-13421.
42
43
44
45 35. Kitchen DB, Decornez H, Furr JR, et al. Docking and scoring in virtual screening for drug
46 discovery: methods and applications. Nat Rev Drug Discov. 2004; 3(11): 935-949.
47
48
49
50 36. Talele TT, Khedkar SA, Rigby AC. Successful applications of computer aided drug
51 discovery: moving drugs from concept to the clinic. Curr Top Med Chem. 2010; 10(1): 127-141.
52
53
54
55
56
57
58
59
60

- 1
2
3 37. Morris GM, Huey R, Lindstrom W, et al. AutoDock4 and AutoDockTools4: Automated
4 docking with selective receptor flexibility. *J Comput Chem.* 2009; 30(16):2785–2791.
5
6
7
8 38. Trott O, Olson AJ. AutoDock Vina: improving the speed and accuracy of docking with a new
9 scoring function, efficient optimization and multithreading. *J Comput Chem.* 2010; 31(2):455–
10 461.
11
12
13
14 39. Friesner RA, Murphy RB, Repasky MP, et al. Extra precision glide: Docking and scoring
15 incorporating a model of hydrophobic enclosure for protein–ligand complexes. *J Med Chem.*
16 2006; 49(21): 6177-6196.
17
18
19
20 40. Stroganov OV, Novikov FN, Stroylov VS, et al. Lead finder: an approach to improve
21 accuracy of protein–ligand docking, binding energy estimation, and virtual screening. *J Chem*
22 *Inf Model.* 2008; 48(12): 2371-2385.
23
24
25
26 41. Chiappori F, D'Ursi P, Merelli I, et al. In silico saturation mutagenesis and docking screening
27 for the analysis of protein-ligand interaction: the Endothelial Protein C Receptor case study.
28 *BMC bioinformatics.* 2009; 10(12):S3.
29
30
31
32 42. D'Ursi P, Chiappori F, Merelli I, et al. Virtual screening pipeline and ligand modelling for
33 H5N1 neuraminidase. *Biochem Bioph Res Co.* 2009; 383(4):445-449.
34
35
36
37 43. Bock J, Gough D. A New Method to Estimate Ligand-Receptor Energetics. *Mol Cell*
38 *Proteomics.* 2002; 11:904.
39
40
41
42 44. Cerón-Carrasco JP, Pérez-Sánchez H, Zúñiga J, et al. Antibodies as Carrier Molecules:
43 Encapsulating Anti-Inflammatory Drugs inside Herceptine. *J Phys Chem B.* 2018; 122(7): 2064-
44 2072.
45
46
47
48 45. Molecular Dynamics on Web. Available from: <http://mmb.irbbarcelona.org>
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 46. van Zundert GCP, Rodrigues JPGLM, Trellet M, et al. The HADDOCK2.2 webserver: User-
4 friendly integrative modeling of biomolecular complexes. *J Mol Biol.* 2015; 428:720-725.
5
6
7
8 47. Margreitter C, Petrov D, Zagrovic B. Vienna-PTM webserver: a toolkit for MD simulations
9 of protein post-translational modifications. *Nucleic Acids Res.* 2013; 41:W422-426.
10
11
12
13 48. CAB-S flex server. Laboratory of Theory of Biopolimers, Faculty of Chemistry. University
14 of Warsaw. Available from: <http://biocomp.chem.uw.edu.pl/CABSflex>
15
16
17
18 49. Heo L, Feig M. PREFMD: a web server for protein structure refinement via molecular
19 dynamics simulations. *Bioinformatics.* 2018; 34(6):1063-1065.
20
21
22
23 50. Devaurs D, Bouard L, Vaisset, M, et al. MoMA-LigPath: a web server to simulate protein-
24 ligand unbinding. *Nucleic Acids Res.* 2013; 41(W1):W297-302.
25
26
27
28 51. Krüger J, Grunzke S, Gesing S, et al. The MoSGrid Science Gateway - A Complete Solution
29 for Molecular Simulations. *J Chem Theory Comput.* 2014; 10(6):2232-2245.
30
31
32
33 52. Chen PC, Kuyucak S. Accurate Determination of the Binding Free Energy for KcsA-
34 Charybdotoxin Complex from the Potential of Mean Force Calculations with Restraints. *Biophys*
35 *J.* 2011; 100(10): 2466-2474.
36
37
38
39 53. Chiappori F, Merelli I, Milanesi L, et al. Static and dynamic interactions between GALK
40 enzyme and known inhibitors: Guidelines to design new drugs for galactosemic patients. *Eur J of*
41 *Med Chem.* 2013; 63: 423-34.
42
43
44
45 54. Foster I, Kesselman C. The history of the grid. *Advances in Parallel Computing.* 2011; 20:3-
46 30.
47
48
49
50
51 **** This article describes the current status and future perspective of grid computing.**
52
53
54
55
56
57
58
59
60

- 1
2
3 55. Foster I, Kesselman C, Nick JM, et al. Grid computing: Making the global infrastructure a
4 reality: Wiley Series; 2003. Chapter 8, The physiology of the grid; p. 217-249.
5
6
7
8 56. Chien A, Foster I, Goddette D. Grid technologies empowering drug discovery. Drug Discov
9 Today. 2002; 7(20 Suppl):S176-80.
10
11
12
13 57. Trombetti GA, Merelli I, Orro A, et al. BGBlast: a BLAST grid implementation with
14 database self-updating and adaptive replication. Stud Health Technol Inform. 2007; 126:23-30.
15
16
17
18 58. Anderson DP. Boinc: A system for public-resource computing and storage. In Proceedings of
19 the Fifth IEEE/ACM International Workshop on Grid Computing 2004. 2004: 4-10. IEEE.
20
21
22 **** This article describes the BOINC distributed infrastructure.**
23
24
25 59. Jacq N, Breton V, Chen HY, et al. Virtual screening on large scale grids. Parallel Comput.
26 2007; 33(4): 289-301.
27
28
29
30 60. Chiappori F, Pucciarelli S, Merelli I, et al. Structural thermal adaptation of β -tubulins from
31 the Antarctic psychrophilic protozoan Euplotes focardii. Proteins. 2011; 80(4):1154-1166
32
33
34
35 61. Foster I, Kesselman C. The globus toolkit. The grid: blueprint for a new computing
36 infrastructure. 1999; 259-278.
37
38
39
40 62. Fifield T, Carmona A, Casajús A, et al. Integration of cloud, grid and local cluster resources
41 with DIRAC. In Journal of Physics: Conference Series, IOP Publishing. 2011; 331(6): 062009.
42
43
44
45 63. Germain-Renaud C, Loomis C, Moscicki J, et al. Scheduling for Responsive Grids. J Grid
46 Comput. 2008; 6(1):15-27.
47
48
49
50 64. Cecchi M, Capannini F, Dorigo A, et al. The glite workload management system. In
51 International Conference on Grid and Pervasive Computing. Springer Berlin Heidelberg. 2009:
52 256-268.
53
54
55
56
57
58
59
60

- 1
2
3
4
5 65. Stevens RD, Robinson AJ, Goble CA. myGrid: personalised bioinformatics on the
6 information grid. *Bioinformatics*. 2003; 19 Suppl 1:i302-4.
7
8
9
10 66. Rauwerda H, Roos M, Hertzberger BO, et al. The promise of a virtual lab in drug discovery.
11 *Drug Discov Today*. 2006; 11(5-6):228-36.
12
13
14
15 67. Anderson DP, Cobb J, Korpela E, et al. SETI@home: an experiment in public-resource
16 computing. *Comm ACM*. 2002; 45(11): 56-61.
17
18
19
20 68. Fedak G, He H, Lodygensky O, et al. EDGeS: a bridge between desktop grids and service
21 grids. In *ChinaGrid Annual Conference, 2008. ChinaGrid'08. The Third*. 2008, Aug: 3-9. IEEE.
22
23
24
25 69. Merelli I, Cozzi P, Ronchieri E, et al. Porting bioinformatics applications from grid to cloud:
26 a macromolecular surface analysis application case study. *Int J High Perform Comput Appl*.
27 2017; 31(3): 182-95.
28
29
30
31
32 70. Kasson PM. *Computational Biology in the Cloud: Methods and New Insights from*
33 *Computing at scale*. *Biocomputing 2013: World Scientific*; 2012. p. 451-453.
34
35
36
37
38 71. Murty J. *Programming amazon web services: S3, EC2, SQS, FPS, and SimpleDB*.
39 *Sebastopol(CA): O'Reilly Media, Inc*; 2008.
40
41
42
43 72. Gruber K. *Google for genomes*. *Nature Research*. 2014.
44
45
46
47 73. Matsunaga A, Tsugawa M, Fortes J. *Cloudblast: Combining mapreduce and virtualization on*
48 *distributed resources for bioinformatics applications*. *eScience. IEEE Fourth International*
49 *Conference; 2008 Dec; 2008*. p. 222-29.
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 74. Forer L, Lipic T, Schonherr S, et al. Delivering bioinformatics mapreduce applications in the
4 cloud. Information and Communication Technology, Electronics and Microelectronics (MIPRO),
5 37th International Convention; 2014 May; IEEE; 2014. p. 373-77.
6
7
8
9
10 75. Schatz MC. CloudBurst: highly sensitive read mapping with MapReduce. *Bioinformatics*.
11 2009; 25(11):1363-9.
12
13
14
15 76. Ahmed L, Edlund A, Laure E, et al. Using iterative MapReduce for parallel virtual screening.
16 IEEE 5th International Conference on Cloud Computing Technology and Science 2; 2013. p. 27-
17 32.
18
19
20
21
22 77. Zhao J, Zhang R, Zhao Z, et al. Hadoop MapReduce framework to implement molecular
23 docking of large-scale virtual screening. Services Computing Conference (APSCC); 2012 Dec;
24 IEEE Asia-Pacific; 2012. p. 350-53.
25
26
27
28
29 78. Constantine RM, Batouche M. Drug discovery for breast cancer based on big data analytics
30 techniques. Information & Communication Technology and Accessibility (ICTA), 5th
31 International Conference; 2015 Dec; 2015. p. 1-6.
32
33
34
35
36 79. Zou Q, Li XB, Jiang WR, et al. Survey of MapReduce frame operation in bioinformatics.
37 *Brief Bioinform*. 2014; 15(4):637-47.
38
39
40
41 80. Salomoni D, Italiano A, Ronchieri E. WNoDeS, a tool for integrated Grid and Cloud access
42 and computing farm virtualization. *Journal of Physics: Conference Series*, IOP Publishing; 2011;
43 331(5):052017.
44
45
46
47
48 81. Ronchieri E, Cesini D, D'Agostino D, et al. The WNoDeS cloud virtualization framework: A
49 macromolecular surface analysis application case study. *Parallel, Distributed and Network-Based*
50 *Processing (PDP)*, 22nd Euromicro International Conference; 2014 Feb; IEEE; 2014. p. 218-22.
51
52
53
54
55
56
57
58
59
60

1
2
3 82. Weber R, Gothandaraman A, Hinde RJ, et al. Comparing hardware accelerators in scientific
4 applications: A case study. IEEE Trans Parallel Distr Syst. 2011; 22(1):58-68.
5
6

7
8 83. Couturier R. Designing Scientific Applications on GPUs. Belfort (France): CRC/Taylor &
9 Francis; 2014.
10
11

12
13 84. TOP500 list. Available from: <https://www.top500.org>
14
15

16
17 85. Amdahl GM. Validity of the single processor approach to achieving large scale computing
18 capabilities. In Proceedings of spring joint computer conference. Atlantic City, New Jersey.
19 ACM; 1967. p. 483-485.
20
21

22
23 86. Gustafson JL. Reevaluating Amdahl's law. Commun ACM.1988; 31(5):532-533.
24
25

26
27 87. Kasam V, Salzemann J, Botha M, et al. WISDOM-II: screening against multiple targets
28 implicated in malaria using computational grid infrastructures. Malar J. 2009 May 1;8:88.
29
30

31 **** This article describes the WISDOM drug discovery challenge**
32
33

34 88. Breton V, Jacq N, Kasam V, et al. Grid-added value to address malaria. IEEE Trans Inf
35 Technol Biomed. 2008; 12(2):173-81.
36
37

38
39 89. Pharmaceutical composition for preventing and treating Malaria, containing compounds that
40 inhibit Plasmepsin II activity, and method of treating Malaria using the same. Available from:
41 <https://patents.google.com/patent/WO2009131384A2/en>
42
43
44

45
46 90. Pettifer S, Ison J, Kalas M, et al. The EMBRACE web service collection. Nucleic Acids Res.
47 2010; 38(Web Server issue):W683-8.
48
49

50
51 91. Milanesi L. BIOINFOGRID: Bioinformatics simulation and modelling based on grid. Model
52 Simulat Sci. 2007: 178-186.
53
54
55
56
57
58
59
60

- 1
2
3 92. Chiappori F, Mattiazzi L, Milanese L, et al. A novel molecular dynamics approach to
4 evaluate the effect of phosphorylation on multimeric protein interface: the alphaB-Crystallin case
5 study. *BMC Bioinformatics*. 2016;17 Suppl 4:57.
6
7
8
9
10 93. Dove MT, Sullivan LA, Walker AM, et al. Molecular dynamics in a grid computing
11 environment: experiences using DL_POLY_3 within the e Minerals escience project. *Molecular*
12 *Simulation*. 2006; 32(12-13):945-952.
13
14
15
16
17 94. Harvey MJ, De Fabritiis G. AceCloud: Molecular Dynamics Simulations in the Cloud. *J*
18 *Chem Inf Model*. 2015; 55(5):909-14.
19
20
21
22 95. Addison E, Keinan S. Using Quantum Molecular Design & Cloud Computing to Improve the
23 Accuracy & Success Probability of Drug Discovery. *Drug Development & Delivery*. 2016;
24 16(2):62-66.
25
26
27
28
29 96. Jiang W, Phillips JC, Huang L, et al. Generalized Scalable Multiple Copy Algorithms for
30 Molecular Dynamics Simulations in NAMD. *Comput Phys Commun*. 2014; 185(3):908-916.
31
32
33
34 97. Merelli I, Morra G, D'Agostino D, et al. High performance workflow implementation for
35 protein surface characterization using grid technology. *BMC Bioinformatics*. 2005; 6 Suppl
36 4:S19.
37
38
39
40
41 98. Hull D, Wolstencroft K, Stevens R, et al. Taverna: a tool for building and running workflows
42 of services. *Nucleic Acids Res*. 2006; 34(Web Server issue):W729-32.
43
44
45
46 99. Oinn T, Li P, Kell DB, et al. Taverna/my Grid: aligning a workflow system with the life
47 sciences community. In *Workflows for e-Science*. Springer, London. 2007: 300-19.
48
49
50
51 100. Stevens R, McEntire R, Goble C, et al. myGrid and the drug discovery process. *Drug*
52 *Discov Today*. 2004; 2(4):140-148.
53
54
55
56
57
58
59
60

1
2
3 101. Buyya R, Branson K, Giddy J, et al. The Virtual Laboratory: A toolset to enable distributed
4 molecular modelling for drug design on the world-wide grid. *Concurr Comput Pract Exp*. 2003;
5 15(1):1–25.
6
7

8
9
10 102. Zhang W, Du X, Ma F, et al. DDGrid: Harness the Full Power of Supercomputing Systems.
11 In *Grid and Cooperative Computing Workshops, 2006. GCCW'06. Fifth International*
12 *Conference*. 2006, Oct: 1-4. IEEE.
13
14

15
16
17 103. Wang Q, Ye Y, Yu K, et al. A Graphical Workflow Modeler for Docking Process in Drug
18 *Discovery*. In *Grid and Cloud Computing: Concepts, Methodologies, Tools and Applications*.
19 2012: 1408-22. IGI Global.
20
21

22
23
24 104. Koday MT, Nelson J, Chevalier A, et al. A Computationally Designed Hemagglutinin
25 *Stem-Binding Protein Provides In Vivo Protection from Influenza Independent of a Host*
26 *Immune Response*. *PLoS Pathog*. 2016; 12(2):1–23.
27
28

29
30
31 105. Fleishman SJ, Whitehead TA, Ekiert DC, et al. Computational design of protein targeting
32 *the conserved stem region of influenza hemagglutinin*. *Science*. 2011; 332(6031):816–821.
33

34 *** The paper describes a computational method for designing proteins. The method is**
35 **applied to the design of different proteins.**
36
37

38
39 106. Gordon SR, Stanley EJ, Wolf S, et al. Computational design of an α -gliadin peptidase. *J Am*
40 *Chem Soc*. 2012; 134(50):20513–20520.
41
42

43
44 107. Wolf C, Siegel JB, Tinberg C, et al. Engineering of Kuma030: A Gliadin Peptidase That
45 *Rapidly Degrades Immunogenic Gliadin Peptides in Gastric Conditions*. *J Am Chem Soc*. 2015;
46 137(40):13106–13113.
47
48

49
50
51 108. Das R, Qian B, Raman S, et al. Structure prediction for CASP7 targets using extensive
52 *all-atom refinement with Rosetta@ home*. *Proteins*. 2007; 69(S8): 118-28.
53
54

- 1
2
3 109. Lyskov S, Gray JJ. The RosettaDock server for local protein-protein docking. *Nucleic*
4 *Acids Res.* 2008; 36(Web Server issue):W233-8.
5
6
7
8 110. Kim DE, Chivian D, Baker D. Protein structure prediction and analysis using the Robetta
9 server. *Nucleic Acids Res.* 2004; 32(Supp 2):W526-W531.
10
11
12
13 111. Lensink MF, Mendez R, Wodak SJ. Docking and scoring protein complexes: CAPRI 3rd
14 Edition. *Proteins.* 2007; 69(4):704-18.
15
16
17
18 112. Bhardwaj G, Mulligan VK, Bahl CD, et al. Accurate de novo design of hyperstable
19 constrained peptides. *Nature.* 2016; 538(7625):329–335.
20
21
22
23 113. Hosseinzadeh P, Bhardwaj G, Mulligan VK, et al. Comprehensive computational design of
24 ordered peptide macrocycles. *Science.* 2017; 358(6369):1461–1466.
25
26
27
28 114. Jayachandran G, Vishal V, Pande VS. Using massively parallel simulation and Markovian
29 models to study protein folding: examining the dynamics of the villin headpiece. *J Chem Phys.*
30 2006; 124(16):164902.
31
32
33
34
35 115. Most powerful distributed computing network [Internet]. 2007 Sep 16. Available from:
36 <http://guinnessworldrecords.com>.
37
38
39
40 116. Beberg A, Ensign D, Jayachandran G, et al. Folding@home: Lessons From Eight Years of
41 Volunteer Distributed Computing. In: *Parallel & Distributed Processing, IEEE International*
42 *Symposium*: 1–8. 2009.
43
44
45
46 * **Authors describe the larger distributed platform worldwide for protein structural**
47 **analysis**
48
49
50
51 117. Kondo D, Javadi B, Malecot P, et al. Cost-benefit analysis of cloud computing versus
52 desktop grids. In *Parallel & Distributed Processing; 2009 May 23-29. IPDPS 2009. IEEE*
53 *International Symposium.* 2009. p. 1-12.
54
55
56
57
58
59
60

- 1
2
3
4
5 118. Bertis V, Bolze R, Desprez F, et al. Large scale execution of a bioinformatic application on
6 a volunteer grid. *Parallel and Distributed Processing*; 2008 Apr 14-18; Miami (FL); IPDPS 2008.
7 IEEE International Symposium. 2008. p. 1-8.
8
9
10
11
12 119. Merelli I, Pérez-Sánchez H, Gesing S, et al. Managing, analysing, and integrating big data
13 in medical bioinformatics: open problems and future perspectives. *Biomed Res Int*. 2014;
14 2014:134023.
15
16
17
18 120. McKenna A, Hanna M, Banks E, et al. The Genome Analysis Toolkit: a MapReduce
19 framework for analyzing next-generation DNA sequencing data. *Genome Res*. 2010; 20(9):1297-
20 303.
21
22
23
24
25 121. Zhou W, Li R, Yuan S, et al. MetaSpark: a Spark-based distributed processing tool to
26 recruit metagenomic reads to reference genomes. *Bioinformatics*. 2017; 33(7):1090-1092.
27
28
29
30 122. Niu B, Zhu Z, Fu L, et al. FR-HIT, a very fast program to recruit metagenomic reads to
31 homologous reference genomes. *Bioinformatics*. 2011; 27(12):1704-5.
32
33
34
35 123. Harnie D, Saey M, Vapirev AE, et al. Scaling machine learning for target prediction in drug
36 discovery using Apache Spark. *Future Generat Comput Syst*. 2017; 67:409-17.
37
38
39
40 124. Buch I, Harvey MJ, Giorgino T, et al. High-throughput all-atom molecular dynamics
41 simulations using distributed computing. *J Chem Inf Model*. 2010; 50(3):397-403.
42
43
44
45 125. Harvey MJ, Giupponi G, Fabritiis GD. ACEMD: accelerating biomolecular dynamics in the
46 microsecond time scale. *J Chem Theory Comput*. 2009; 5(6): 1632-1639.
47
48
49
50 126. Hess B, Kutzner C, Van der spoel D, et al. GROMACS 4: algorithms for highly efficient,
51 load-balanced, and scalable molecular simulation. *J Chem Theory Comput*. 2008; 4(3): 435-447.
52
53
54
55
56
57
58
59
60

- 1
2
3 127. Marin I, Goga N, Goga M. Benchmarking MD Systems Simulations on the Graphics
4 Processing Unit and Multi-core Systems. IEEE International Symposium on Systems
5 Engineering (ISSE); 2016 Oct 3-5; Edinburgh (UK); 2016.
6
7
8
9
10 128. Hung CL, Chen CC. Computational approaches for drug discovery. Drug Dev Res. 2014;
11 75(6):412-8.
12
13
14
15 129. Vogt M, Bajorath J. Chemoinformatics: a view of the field and current trends in method
16 development. Bioorg Med Chem. 2012; 20(18):5317-23.
17
18
19
20 130. Ma C, Wang L, Xie XQ. GPU accelerated chemical similarity calculation for compound
21 library comparison. J Chem Inf Model. 2011; 51(7):1521-7.
22
23
24
25 131. Malhat MG, El-Sisi AB. Parallel Ward Clustering for Chemical Compounds Using
26 OpenCL. Tenth International Conference on Computer Engineering & Systems; 2015 Dec 23-24;
27 Cairo (Egypt); ICCES; 2016.
28
29
30
31
32 132. Lo YT, Wang HW, Pai TW, et al. Protein–ligand binding region prediction (PLB-SAVE)
33 based on geometric features and CUDA acceleration. BMC Bioinformatics. 2013; 14 Suppl 4:S4.
34
35
36
37 133. Imbernón B, Cecilia JM, Pérez-Sánchez H, et al. METADOCK: A parallel metaheuristic
38 schema for virtual screening methods. Int J High Perform Comput Appl. 2017.
39
40
41
42 134. Yuriev E, Holien J, Ramsland PA. Improvements, trends, and new ideas in molecular
43 docking: 2012-2013 in review. J Mol Recognit. 2015; 28(10):581-604.
44
45
46
47 135. Lionta E, Spyrou G, Vassilatis DK, et al. Structure-based virtual screening for drug
48 discovery: principles, applications and recent advances. Curr Top Med Chem. 2014;
49 14(16):1923-38.
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 136. McIntosh-Smith S, Price J, Sessions RB, et al. High performance in silico virtual drug
4 screening on many-core processors. *Int J High Perform Comput Appl*. 2015; 29(2):119-134.
5
6
7
8 137. Fang Y, Ding Y, Feinstein WP, et al. GeauxDock: Accelerating Structure-Based Virtual
9 Screening with Heterogeneous Computing. *PLoS One*. 2016; 11(7):e0158898.
10
11
12
13 138. Feinstein W, Brylinski M. Structure-Based Drug Discovery Accelerated by Many-Core
14 Devices. *Curr Drug Targets*. 2016; 17(14):1595-1609.
15
16
17
18 139. Krige S, Mackey M, McIntosh-Smith S, et al. Porting a Commercial Application to
19 OpenCL. *Proceedings of the International Workshop on OpenCL 2013 & 2014 - IWOCL '14*;
20 2014.
21
22
23
24
25 140. Harvey MJ, Giupponi G, Fabritiis GD. ACEMD: Accelerating Biomolecular Dynamics in
26 the Microsecond Time Scale. *J Chem Theory Comput*. 2009; 5(6):1632-9.
27
28
29
30 141. Sukhwani B, Herbordt MC. Fast Binding Site Mapping Using GPUs and CUDA. 2010
31 IEEE International Symposium on Parallel & Distributed Processing, Workshops and Phd Forum
32 (IPDPSW); 2010; 04.
33
34
35
36
37 142. McArt DG, Bankhead P, Dunne PD, et al. cudaMap: a GPU accelerated program for gene
38 expression connectivity mapping. *BMC Bioinformatics*. 2013; 14:305.
39
40
41
42 143. Hu X, Liu Q, Zhang Z, et al. SHEsisEpi, a GPU-enhanced genome-wide SNP-SNP
43 interaction scanning algorithm, efficiently reveals the risk genetic epistasis in bipolar disorder.
44 *Cell Res*. 2010; 20(7):854-7.
45
46
47
48
49 144. Irwin JJ, Shoichet BK, Mysinger MM, et al. Automated docking screens: a feasibility study.
50 *J Med Chem*. 2009; 52(18):5712-5720.
51
52
53
54
55
56
57
58
59
60

- 1
2
3 145. Tsai TY, Chang KW, Yu-Chian Chen C. iScreen: world's first cloud-computing web server
4 for virtual screening and de novo drug design based on TCM database@Taiwan. J Comput
5 Aided Mol Des. 2011; 25(6):525-531.
6
7
8
9
10 146. Grosdidier A, Zoete V, Michielin O. SwissDock, a protein-small molecule docking web
11 service based on EADock DSS. Nucleic Acids Res. 2011; 39(Web Server issue):W270-7.
12
13
14
15 147. Prakhov ND, Chernorudskiy AL, Gainullin MR. VSDocker: a tool for parallel high-
16 throughput virtual screening using AutoDock on Windows-based computer clusters.
17 Bioinformatics. 2010; 26(10):1374-1375.
18
19
20
21
22 148. Sánchez-Linares I, Pérez-Sánchez H, Cecilia JM, et al. High-Throughput parallel blind
23 Virtual Screening using BINDSURF. BMC Bioinformatics. 2012; 13 Suppl 14:S13.
24
25
26
27 149. Zoete V, Daina A, Bovigny C, et al. SwissSimilarity: A Web Tool for Low to Ultra High
28 Throughput Ligand-Based Virtual Screening. J Chem Inf Model. 2016; 56(8):1399-404.
29
30
31
32 150. Shang J, Dai X, Li Y, et al. HybridSim-VS: a web server for large-scale ligand-based virtual
33 screening using hybrid similarity recognition techniques. Bioinformatics. 2017; 33(21):3480-
34 3481.
35
36
37
38
39 151. Li H, Leung KS, Wong MH, et al. USR-VS: a web server for large-scale prospective virtual
40 screening using ultrafast shape recognition techniques. Nucleic Acids Res. 2016; 44(W1):W436-
41 W441.
42
43
44
45
46 152. Wang X, Chen H, Yang F, et al. iDrug: a web-accessible and interactive drug discovery and
47 design platform. J Cheminform. 2014; 6(28).
48
49
50
51 153. Gong J, Cai C, Liu X, et al. ChemMapper: a versatile web server for exploring
52 pharmacology and chemical structure association based on molecular 3D similarity method.
53 Bioinformatics. 2013; 29(14):1827-1829.
54
55
56
57
58
59
60

1
2
3
4
5 154. Koes DR, Camacho CJ. ZINCPharmer: pharmacophore search of the ZINC database.
6 Nucleic Acids Res. 2012; 40(W1):W409-W414.
7
8
9

10 155. Dong J, Yao ZJ, Wen M, et al. BioTriangle: a web-accessible platform for generating
11 various molecular representations for chemicals, proteins, DNAs/RNAs and their interactions. J
12 Cheminform. 2016; 8:34.
13
14
15

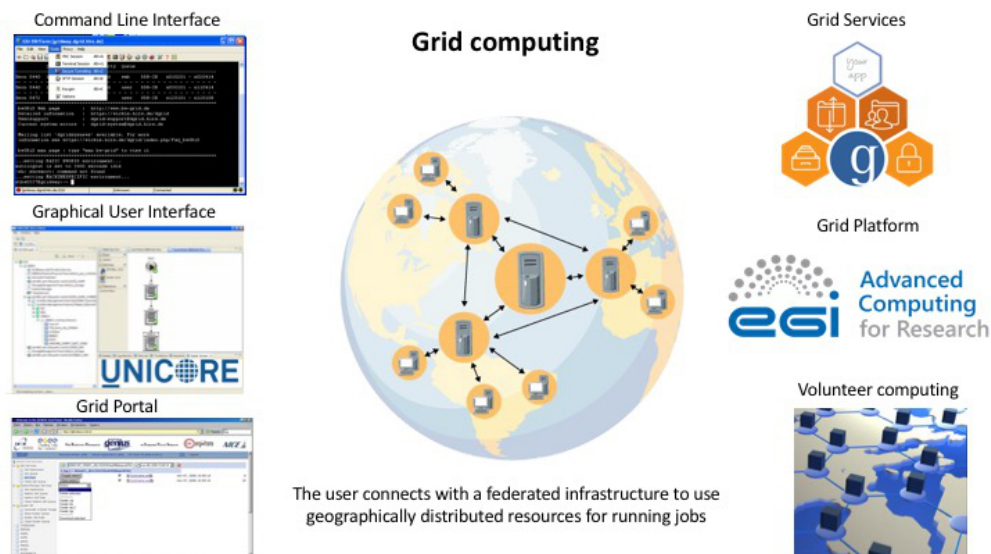
16 156. Click2Drug. Swiss Institute of Bioinformatics. Available from: <http://www.click2drug.org>
17
18
19

20 157. Novartis Case Study. Amazon Web Services [Internet]. Available from:
21 <https://aws.amazon.com/it/solutions/case-studies/novartis/>.
22
23
24

25 158. List M. Using Docker Compose for the Simple Deployment of an Integrated Drug Target
26 Screening Platform. J Integr Bioinform. 2017; 14(2).
27
28
29

30 159. Merelli I, Fornari F, Tordini F, et al. Exploiting Docker Containers over Grid Computing
31 for a comprehensive Study of Chromatin Conformation in different cell types. Journal of Parallel
32 and Distributed Computing, in press.
33
34
35
36

37 160. AstraZeneca taps AI for drug discovery in deal with Berg. Reuters. 2018. Available from:
38 [https://www.reuters.com/article/us-astrazeneca-ai-berg/astrazeneca-taps-ai-for-drug-discovery-](https://www.reuters.com/article/us-astrazeneca-ai-berg/astrazeneca-taps-ai-for-drug-discovery-in-deal-with-berg-idUSKCN1B81G1)
39 [in-deal-with-berg-idUSKCN1B81G1](https://www.reuters.com/article/us-astrazeneca-ai-berg/astrazeneca-taps-ai-for-drug-discovery-in-deal-with-berg-idUSKCN1B81G1)
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

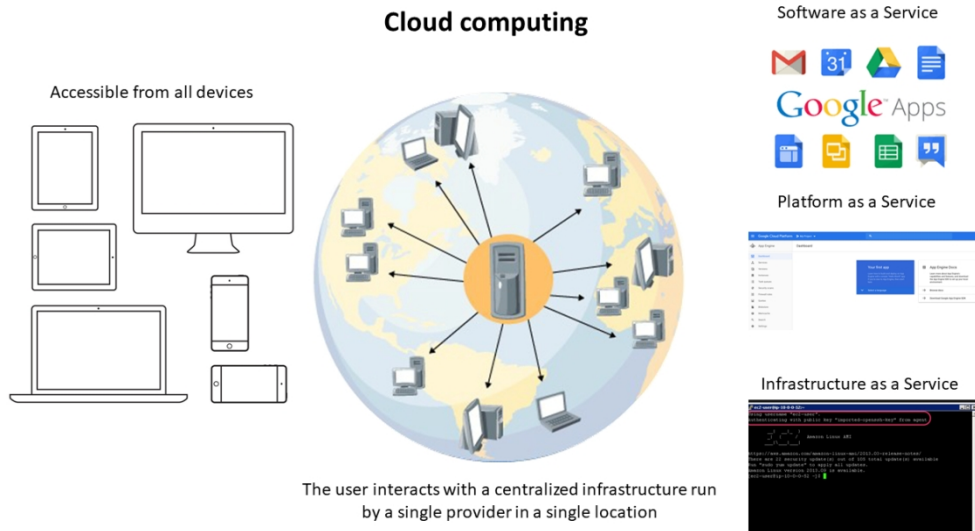


25 The Grid computing paradigm. On the left side, three examples of user interface are shown (Command Line
26 Interface, Graphical User Interface, Grid Portal). The central panel shows the common federated and
27 geographical dispersed infrastructure typical of grid computing. On the right side, the three general types of
28 grid infrastructures are reported (Grid Services, Grid Platform, Volunteer Computing).

29 338x190mm (54 x 54 DPI)

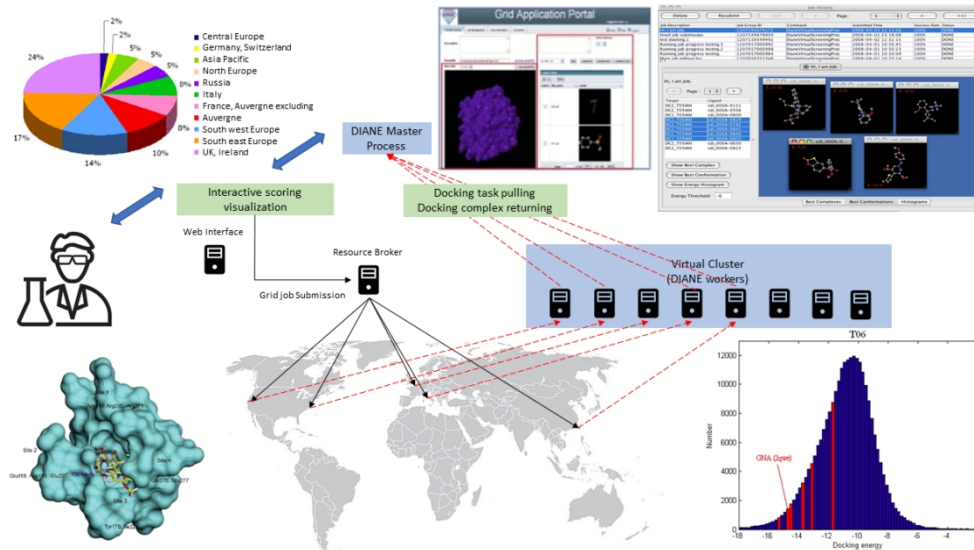
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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



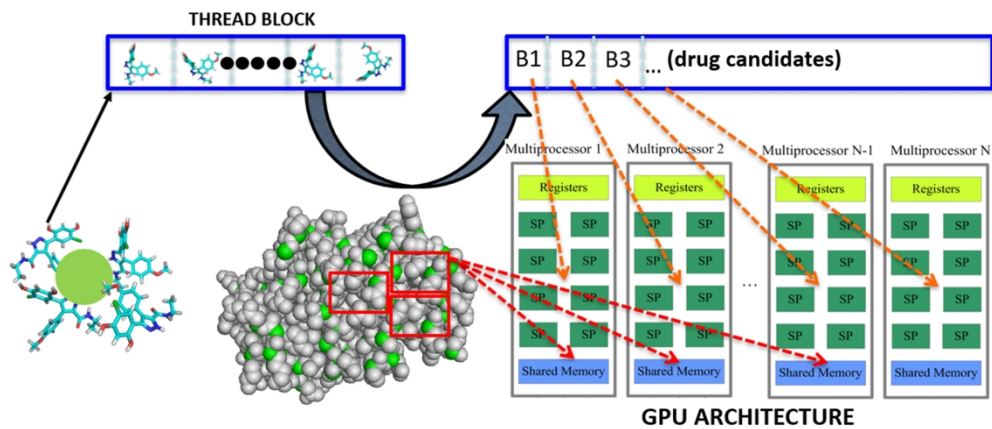
The Cloud computing paradigm. On the left side, the figure shows how cloud infrastructures are easily accessible from any kind of device. The central panel shows the common centralized and single-entry point infrastructure typical of Cloud Computing. On the right side, three popular examples of Cloud platforms are shown (Software as a Service, Platform as a Service, Infrastructure as a Service).

108x60mm (300 x 300 DPI)



The architecture of the grid-based infrastructure adopted for WISDOM and the Avian flu virtual screenings. DIANE was used as job dispatcher and a Grid Application Portal was developed to visually inspect achieved results. In the top left corner, the distribution of the jobs among the different geographical regions of the world. In the bottom left corner, the structure of the Neuraminidase that was target of the screening. In the top right corner, the job monitoring application used in the virtual screening. In the bottom right corner, the enrichment curve of one of the leading compounds identified with the virtual screening.

338x190mm (96 x 96 DPI)



20
21
22
23

Data distribution model on GPU in METADOCK. Drug candidates are grouped in blocks of threads and mapped to the multiprocessors for their computation. The shared memory, in the lower part, stores parts of the receptor molecule for reuse by the candidates that integrate the same block of threads.

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

440x186mm (300 x 300 DPI)