



UvA-DARE (Digital Academic Repository)

A Grid-Based Hiv Expert System

Sloot, P.M.A.; Boukhanovsky, A.V.; Keulen, W.; Tirado Ramos, A.; Boucher, C.A.B.

DOI

[10.1007/s10877-005-0673-2](https://doi.org/10.1007/s10877-005-0673-2)

Publication date

2005

Published in

Journal of Clinical Monitoring and Computing

[Link to publication](#)

Citation for published version (APA):

Sloot, P. M. A., Boukhanovsky, A. V., Keulen, W., Tirado Ramos, A., & Boucher, C. A. B. (2005). A Grid-Based Hiv Expert System. *Journal of Clinical Monitoring and Computing*, 19(4-5), 263-278. <https://doi.org/10.1007/s10877-005-0673-2>

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

1 A GRID-BASED HIV EXPERT SYSTEM

2 Peter M.A. Sloot,¹ Alexander V. Boukhanovsky,²

3 Wilco Keulen,³ Alfredo Tirado-Ramos,¹ and

4 Charles A. Boucher⁴

Sloot P MA, Boukhanovsky AV, Keulen W, Tirado-Ramos A, Boucher CA. A grid-based HIV expert system.

J Clin Monit 2005; xxx: 1–16

ABSTRACT. Objectives. This paper addresses Grid-based integration and access of distributed data from infectious disease patient databases, literature on *in-vitro* and *in-vivo* pharmaceutical data, mutation databases, clinical trials, simulations and medical expert knowledge. **Methods.** Multivariate analyses combined with rule-based fuzzy logic are applied to the integrated data to provide ranking of patient-specific drugs. In addition, cellular automata-based simulations are used to predict the drug behaviour over time. Access to and integration of data is done through existing Internet servers and emerging Grid-based frameworks like Globus. Data presentation is done by standalone PC based software, Web-access and PDA roaming WAP access. The experiments were carried out on the DAS, a Dutch Grid testbed. **Results.** The output of the problem-solving environment (PSE) consists of a prediction of the drug sensitivity of the virus, generated by comparing the viral genotype to a relational database which contains a large number of phenotype-genotype pairs. **Conclusions.** Artificial Intelligence and Grid technology is effectively used to abstract knowledge from the data and provide the physicians with adaptive interactive advice on treatment applied to drug resistant HIV. An important aspect of our research is to use a variety of statistical and numerical methods to identify relationships between HIV genetic sequences and antiviral resistance to investigate consistency of results.

KEY WORDS. grid, HIV, PSE, expert system, artificial intelligence, bio-statistics.

1. INTRODUCTION

1.1. Motivation

Forty two million people worldwide have been infected with HIV and 12 million have died, over the last 20 years. Figure 1 shows the pan-epidemic extent of HIV infections.

Effective antiretroviral therapy has lead to sustained HIV viral suppression and immunological recovery in patients who have been infected with the virus. The incidence of AIDS has declined in the Western world with the introduction of effective antiretroviral therapy, though questions on “When to start treatment? What to start with? How to monitor patients?” remain heavily debated. Adherence to antiretroviral treatment remains the cornerstone of effective treatment, and failure to adhere is the strongest predictor of virological failure. Long-term therapy can lead to metabolic complications. Other treatment options are now available, with the recent introduction to clinical practice of fusion inhibitors, second-generation non-nucleoside reverse transcriptase inhibitors, and nucleotide reverse transcriptase inhibitors. The sheer complexity of the disease,

From the ¹Section Computational Science, University of Amsterdam, Kruislaan 403, 1098 SJ Amsterdam, The Netherlands, ²Institute for High Performance Computing and Information Systems, Bering St, 38, St. Petersburg, Russia, ³Virology Education, 69042 Utrecht, The Netherlands, ⁴University Medical Center, University of Utrecht, 3508 GA Utrecht, The Netherlands.

Received—, and in revised form—. Accepted for publication—.

Based on “A Grid-based HIV Expert System”, by P.M.A. Sloot, A.V. Boukhanovsky, W. Keulen, and C.A. Boucher, which appeared in the IEEE/ACM International Symposium on Cluster Computing and the Grid, Cardiff, UK, May 9–12, 2005. ©2005 IEEE.

Address correspondence to Peter M.A. Sloot, Section Computational Science, University of Amsterdam, Kruislaan 403, 1098 SJ Amsterdam, The Netherlands
E-mail: sloot@science.uva.nl

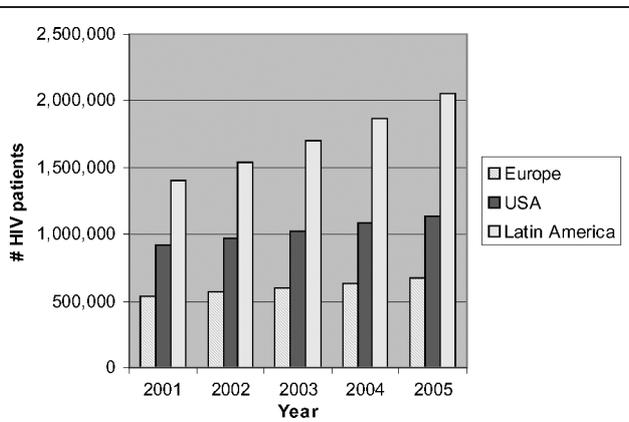


Fig. 1. Worldwide spread of HIV infections, history and near future perspective.

51 the distribution of the data, the required automatic updates
 52 to the knowledgebase and the efficient use and integration
 53 of advanced statistical and numerical techniques necessary
 54 to assist the physician motivated us to explore the novel
 55 possibilities supported by Grid technology.

56 In this position paper we describe ongoing research in
 57 our 3 laboratories (Utrecht, St. Petersburg and Amsterdam)
 58 addressing the development of a Grid based medical deci-
 59 sion support system. The goal of the research is to investi-
 60 gate novel computational methods and techniques that sup-
 61 port the development of a user friendly integrated support
 62 system for physicians. We use emerging Grid-technology
 63 to combine data discovery, data mining, statistical analyses,
 64 numerical simulation and data presentation [1].

65 The paper is organized as follows. Chapter 2 describes
 66 the background of HIV research and a prototypical rule-
 67 based approach to data analyses. In chapter 3 we give an
 68 overview of the two computational techniques we study
 69 to understand the temporal variability of HIV populations
 70 through stochastic modeling and the evolution of HIV
 71 infection and the onset of AIDS through Cellular Automata
 72 (CA) modeling. Chapter 4 describes a first approach to
 73 advanced data presentation through roaming devices such
 74 as Personal Digital Assistants (PDA's).

75 *1.2. Background*

76 *1.2.1. Clinical aspects of HIV*

77 The clinical management of patients infected with Human
 78 Immunodeficiency Virus (HIV) is based on studies on the
 79 pathogenesis of the disease and the results of trials evaluat-
 80 ing the effects of anti-HIV drugs. Retrospective analysis of
 81 large cohorts has identified laboratory markers for disease

82 progression, such as the amount of virus (HIV-RNA) and
 83 the number of T helper cells (CD4 + cells) in blood. In ad-
 84 dition the results of prospective drug trials have generated
 85 data on effectiveness of individual drugs and drug combi-
 86 nations and the effect of drug resistant viruses on therapy
 87 outcome. Currently clinicians are limited in the practical
 88 use of this information because in most cases they are only
 89 provided with statistical relationships between individual
 90 parameters and disease or therapy outcome. Large data sets
 91 have not been analyzed and made available in such a way
 92 that it allows a clinician to use the available data in more
 93 clinical settings. The availability of large databases and the
 94 development of innovative data mining approaches create
 95 the opportunity to develop systems which allow the prac-
 96 ticing clinician to determine the risk profile for disease
 97 development, or the change or success for a given regimen
 98 for his individual patients. Such a system will determine the
 99 rate of success for different drug regimens by taking into
 100 account the effect and interaction of all relevant laboratory
 101 and clinical parameters and by comparing the results for
 102 similar patients available in the database.

103 Currently there are fifteen drugs licensed for treatment of
 104 individuals infected with HIV. These drugs belong to two
 105 classes, one inhibiting the viral enzyme reverse transcrip-
 106 tase and another inhibiting the viral protease. These drugs
 107 are used in combination with therapy to maximally inhibit
 108 viral replication and decrease HIV-RNA to below levels of
 109 detection levels (currently defined as below 50 copies per
 110 ml) in blood. Treatment with drug combinations is suc-
 111 cessful in inhibiting viral replication to undetectable levels
 112 in only 50% of the cases. In the remaining 50% of cases
 113 viruses can be detected with a reduced sensitivity to one
 114 or more drugs from the patients' regimen. The molecular
 115 base for resistance has been, and still is, focus of extensive
 116 research. Over 80 amino acid positions in the viral enzyme
 117 reverse transcriptase (RT) and 40 positions in the protease
 118 enzyme can undergo changes when exposed to selective
 119 drug pressure in vitro or in vivo. For some drugs, at cer-
 120 tain positions, a change towards a specific new amino acid
 121 is seen. At other positions several alternative amino acids
 122 may appear and cause (variable) levels of resistance to one
 123 or more drugs. In theory, therefore, an infinite number
 124 of combinations of amino acid changes could appear and
 125 cause resistance in vivo. Preliminary clinical observations
 126 however show that specific amino acid changes at a limited
 127 number of positions and a limited number of combina-
 128 tions prevail. In addition to changing drug sensitivity some
 129 amino acid changes may also influence the replication po-
 130 tential of HIV. Amino acids selected initially during a failing
 131 regimen cause resistance to the drugs the patient is taking,
 132 but at the same time may decrease the capacity of the virus
 133 to replicate. Changes appearing later do not function to
 134 further increase resistance but merely function to restore

135 the capacity of the virus to replicate (“viral fitness”). Sev-
 136 eral clinical studies have been performed recently to evalu-
 137 ate the clinical benefit of resistance-guided therapy. These
 138 studies show that a better virological response is obtained
 139 in patients who are failing their therapy, when their new
 140 regimen is chosen on the basis of their resistant profile.
 141 In three out of the four studies from last year the results
 142 showed that if new regimens were selected on the basis of
 143 the mutations (viral resistance genotype) the results were
 144 better as compared to standard care approaches. Currently,
 145 the basis for clinical interpretation of the viral genotype is
 146 based on data sets relating mutations to changes in drug sen-
 147 sitivity, and/or data sets directly relating mutations present
 148 in the virus to clinical responses to specific regimens. Ini-
 149 tially, experts compared the observed mutations to lists of
 150 published sequences taken from the literature, and based
 151 on this comparison would select a regimen.

152 1.2.2. Prototype support system

153 Recently, first generation bioinformatics software pro-
 154 grams have been developed to support clinicians. Examples
 155 of such systems are the Virtual Phenotype developed by
 156 Virco NV, and a first generation decision support system
 157 (Retrogram TM) developed by Virology Networks BV in
 158 collaboration with parts of our research team. The out-
 159 put of these programs consists of a prediction of the drug
 160 sensitivity of the virus generated by comparing the viral
 161 genotype to a relational database containing a large num-
 162 ber of phenotype-genotype pairs. The Retrogram decision
 163 software interprets the genotype of a patient by using rules
 164 developed by experts on the basis of the literature, taking
 165 into account the relationship of the genotype and phe-
 166 notype. In addition, it is based on (limited) available data
 167 from clinical studies and on the relationship between the
 168 presence of genotype directly to clinical outcome. It is im-
 169 portant to realise however that these systems focus on bio-
 170 logical relationships and are limited to the role of resistance.
 171 The next step will be to use clinical databases and inves-
 172 tigate the relationship between the viral resistance profile
 173 (mutational profile and/or phenotypic data) and therapy
 174 outcome measures such as amount of virus (HIV-RNA)
 175 and CD4+ cells. A summary of the flow of data is shown
 176 in Figure 2.

177 1.2.3. Data collection

178 Large high quality clinical and patient databases are used
 179 to explore the relationships described above and to de-
 180 velop a first prototype matching system. The Athena co-
 181 hort is a large Dutch observational clinical cohort study

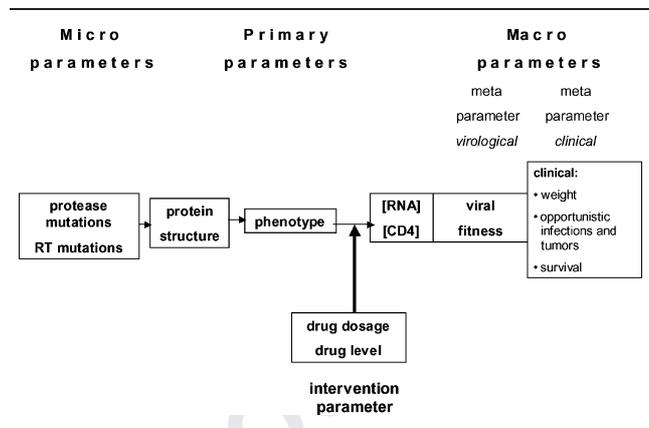


Fig. 2. From molecule to man: Hierarchical data flow model for infectious diseases.

aiming at the surveillance of antiretroviral treatment sup- 182
 ported by the government. The cohort consists of 3000 183
 patients from whom data are centrally collected through a 184
 decentralized data entry system. Within the cohort 600 pa- 185
 tients are studied intensively, whose phenotypic and geno- 186
 typic data, drug levels and CD4+ and HIV-RNA patterns 187
 are collected. Phenotype, genotype, viral fitness and drug 188
 levels as CD4+ and HIV-RNA patterns will be collected 189
 from two large international trials (sponsored by Roche 190
 Pharmaceuticals), evaluating the effect of a new fusion in- 191
 hibitor drug (T20), and representing 1000 patients. The 192
 third database will be from the international multi-center 193
 Great study, sponsored by Virology Networks BV. Within 194
 this study the value of the Retrogram decision support 195
 program is evaluated and similar parameters as described 196
 above will be collected. Within this study 360 patients will 197
 be enrolled. 198

The Viradapt study showed that the virological response 199
 was better in the patient group in which genotype and rule- 200
 based interpretation was used as compared to the standard 201
 of care arm [2]. On the basis of these results, a more elabo- 202
 rate decision support software system (Retrogram version 203
 1.0) was built in collaboration with Virology Networks 204
 B.V. This system ranks the efficacy of the antiretroviral 205
 drugs within each class. The ranking is based on expert 206
 interpretation of two types of data. The software system 207
 estimates the drug sensitivity for the fifteen drugs by in- 208
 terpreting the genotype of a patient by using mutational 209
 algorithms. These mutational algorithms are developed by 210
 a group of experts on the basis of the scientific literature, 211
 taking into account the published data relating genotype to 212
 phenotype. In addition, the ranking is based on data from 213
 clinical studies on the relationship between the presence of 214
 particular mutations and clinical or virological outcome. 215

The Athena cohort is a large Dutch observational clinical 216
 cohort study aiming at the surveillance of antiretroviral 217

218 treatment supported by the Dutch government. The co-
 219 hort consists of 3000 patients from whom clinical, viro-
 220 logical, immunological and data on drug side effects are
 221 centrally collected through a decentralised data entry sys-
 222 tem. Within this cohort 600 patients are studied intensively,
 223 phenotypic and genotypic data, drug levels and CD4+ and
 224 HIV-RNA patterns are collected. From two large interna-
 225 tional trials (sponsored by Roche Pharmaceuticals) evalu-
 226 ating the effect of a new fusion inhibitor drug (T20),
 227 representing 1000 patients from whom also phenotype,
 228 genotype, viral fitness, drug levels as CD4+ and HIV-RNA
 229 patterns will be collected. The third database will be from
 230 the international multi-center Great study sponsored by
 231 Virology Networks BV, within this study the value of the
 232 Retrogram decision support program is evaluated and sim-
 233 ilar parameters a described above will be collected, within
 234 this study 360 patients will be enrolled. Another dataset
 235 will come from the Italian Musa study, in this trial data will
 236 be collected from 450 patients followed over a year. Entry
 237 point to the trial is failing a fist or second regimen, subse-
 238 quently patients will be genotyped and a new regimen will
 239 be selected on the basis of Retrogram 1.4 or the Virtual
 240 Phenotype from Virco (Belgium).

241 Throughout the duration of the project we will collect
 242 additional datasets. These datasets may serve to further re-
 243 fine our models and first version software and may also be
 244 use to perform validation studies.

245 1.2.4. Data analysis

246 The primary goal of the data analysis is to identify pat-
 247 terns of mutations (or naturally occurring polymorphisms)
 248 associated with resistance to antiviral drugs and to predict
 249 the degree of *in-vitro* or *in-vivo* sensitivity to available drugs
 250 from an HIV genetic sequence. The statistical challenges
 251 in doing such analyses arise from the high dimensionality
 252 of these data. A variety of approaches have been de-
 253 veloped to handle this type of data, including clustering,
 254 recursive partitioning, and neural informatics. Neural in-
 255 formatics is used for synthesis of heuristic models received
 256 by methods of knowledge engineering, and results of the
 257 formal multivariate statistical analysis in uniform systems.
 258 Clustering methods have been used to group sequences
 259 that are "near" each other according to some measure of
 260 genetic distance [3]. Once clusters have been identified,
 261 recursive partitioning can be used to determine the im-
 262 portant predictors of drug resistance, as measured by *in-*
 263 *vitro* assays or by patient response to antiviral drugs. Prin-
 264 ciple component analyses can help to identify what are the
 265 most important sources of variability in the HIV genome.
 266 An important aspect of our research is to use a variety of
 267 methods to identify relationships between HIV genetic se-

quences and antiviral resistance to validate the consistency 268
 of results. 269

270 The molecular sequences of the viral enzymes reverse
 271 transcriptase and protease are the micro parameters in the
 272 model. In theory an infinite number of combinations of
 273 mutations could appear and cause (variable) changes in viral
 274 drug sensitivity and viral replication capacity (See also Ta-
 275 ble 1). Clinical datasets however show that specific amino
 276 acid changes at a limited numbers of positions in a lim-
 277 ited number of combinations prevail. HIV-RNA and CD4
 278 are the primary parameters determining disease outcome.
 279 HIV-RNA, the amount of HIV-RNA genomic copies per
 280 ml plasma, has been validated as being highly predictive of
 281 clinical outcome. HIV-RNA and CD4+ cell numbers are
 282 now the standard endpoint in clinical trials for approval of
 283 new antiretroviral drugs. A patient's HIV-RNA may range
 284 between a few hundred to millions of RNA copies per
 285 ml plasma. The CD4+ cell numbers in peripheral blood
 286 range typically between zero and thousand. Whereas the
 287 predictive clinical value of both parameters has been deter-
 288 mined initially in untreated individuals, they have also been
 289 shown to be of predictive value also for patients under an-
 290 tiretroviral therapy. Recently observations have been pub-
 291 lished indicating that in some patients under highly active
 292 antiretroviral therapy (HAART) a disconnect may occur
 293 between the response in HIV-RNA and in CD4 counts.
 294 Typically, in these patients a rise in HIV-RNA as conse-
 295 quence of incomplete inhibition of viral replication under
 296 therapy is not paralleled by a continuous decrease in CD4
 297 counts. This disconnect has been explained by a decrease

Table 1. Parameters for the data analyses. Here the hierarchical ap-
 proach shown in Figure 2 is extended to detail the content of the
 parameters

Micro Parameter	Protease Mutations Reverse Transcriptas Mutations
Primary Parameter	HIV-RNA CD4 Drug Resistance
Macro Parameter	Meta Parameter: Viral Fitness Virological Meta Parameter: Weight Clinical Opportunistic Infections and Tumors Survival
Intervention Parameter	Drug Dosage Bio-availability of Drug/Drug Level

298 in the viral replicative capacity ('viral fitness') which leads
299 to a decrease in capacity to lower CD4 counts.

300 The patient's weight and secondary opportunistic infec-
301 tions and/or malignancies are parameters that determine
302 disease outcome and survival time. Currently there are fif-
303 teen drugs licensed for treatment of individuals infected
304 with HIV: More than ten inhibitors have been developed
305 which inhibit the reverse transcriptase process. These in-
306 hibitors can be classified in two sub-categories that dif-
307 fer in the way they inhibit the RT-enzyme, nucleoside
308 (analogue) RT-inhibitors (NRTI) and the non-nucleoside
309 RT-inhibitors (NNRTI). These compounds inhibit the
310 protease enzyme, which acts much later on in the HIV
311 replication cycle than reverse transcriptase.

312 The protease is responsible for cleaving a long poly-
313 protein into smaller functional proteins. The overall ex-
314 posure to antiretroviral drugs has been shown to be an
315 important factor for the degree of success for a given ther-
316 apy. The overall exposure can be captured by parameters
317 as dosage and bio-availability which will codetermine the
318 drug level within an individual patient. Given the relation-
319 ships between exposure and antiviral efficacy, variability in
320 drug levels (which may be due to differences in patient
321 adherence to their regimens) will contribute to virologi-
322 cal and immunological outcome. Individuals with relatively
323 low exposure are more likely to experience virological fail-
324 ure than those with a high exposure.

325 2. METHODS AND MATERIALS

326 2.1. Modeling the dynamics and temporal variability 327 of HIV-1 populations

328 In addition to rule based and parameter based decision sup-
329 port we developed statistical models and cellular automata
330 based models to study the dynamics of the HIV popula-
331 tions. These 2 numerical models run on Grid-resources.
332 The output is integrated with the medical support system
333 and accessible to the end-user. In this paragraph we briefly
334 outline the two computational methods. Details are be-
335 yond the scope of this paper; we refer to the references
336 provided.

337 2.1.1. A cellular automata model to study the evolution 338 of HIV infection and the onset of AIDS

339 A cellular automata model to study the evolution of HIV
340 infection and the onset of AIDS is developed. The model
341 takes into account the global features of the immune re-
342 sponse to any pathogen, the fast mutation rate of the HIV,

and a fair amount of spatial localization, which may occur 343
in the lymph nodes. The dynamics of the cellular automata 344
requires high throughput computing, which is provided by 345
the resource management of the Grid. In this section, we 346
employ non-uniform Cellular Automata (CA's) to simulate 347
drug treatment of HIV infection, in which each compu- 348
tational domain may contain different CA rules, in con- 349
trast to normal uniform CA models. Ordinary (or par- 350
tial) differential equation models are insufficient to de- 351
scribe the two extreme time scales involved in HIV in- 352
fection (days and decades), as well as the implicit spatial 353
heterogeneity. Zorzenon dos Santos et al. [7] reported a 354
cellular automata approach to simulate three-phase pat- 355
terns of human immunodeficiency virus (HIV) infection 356
consisting of primary response, clinical latency and onset 357
of acquired immunodeficiency syndrome. We developed a 358
non-uniform CA model to study the dynamics of drug 359
therapy of HIV infection, which simulates four-phases 360
(acute, chronic, drug treatment responds and onset of 361
AIDS). Our results indicate that both simulations (with and 362
without treatments) evolve to the same steady state. Three 363
different drug therapies (mono-therapy, combined drug 364
therapy and HAART) can also be simulated in our model. 365
Our model for prediction of the temporal behaviour of the 366
immune system to drug therapy qualitatively corresponds 367
to clinical data. 368

Pseudo Code 1a: HI Model (Adapted from Zorzenon dos 369
Santos R. M., Phys. Rev. Let. 2001). H = healthy cell, 370
A1 and A2 are infected cells at different time steps. 371

Assume: {H, A1(t), A2(t+ τ), D}; 1 time-step = 1
week; Simulation of lymph-node;
Moore neighbourhood and square
lattices used

Rule 1: (a) If it has at least one infected-A1
neighbor, it becomes infected-A1
(b) If it has no infected-A1 neighbor but
does have at least R ($2 < R < 8$)
infected-A2 neighbors, it becomes
infected-A1
(c) Otherwise it stays healthy

Rule 2: An infected-A1 cell becomes infected-A2
after τ time steps

Rule 3: Infected-A2 cells become dead cells

Rule 4: (a) Dead cells can be replaced by healthy
cells with probability *prepl in the next step*.
(b) Each new healthy cell introduced may
be replaced by an infected-A1 with
probability *pinfec*

374 This CA (Pseudo-code 1a) mimics in a simple way the
 375 dynamical properties of a HIV infection; next we intro-
 376 duce drug therapy into the model by modelling a response
 377 function Presp and changing only rule 1.

378 Pseudo Code 1b: Advanced HI Model, taking into
 379 account drug therapy effects.

Rule 1:

(a) If there is one A1 neighbor after the starting of drug
 therapy, $N(0 \leq N \leq 7)$ neighbor healthy cells become
 infected-A1 in the next time steps with probability presp .
 Otherwise, all of eight neighbors become infected-A1.

N represents effectiveness of drugs.

$N = 0$: no replication;

$N = 7$: less effective for the drug.

$\text{Presp}(t - t_s)$ represents certain response function of drug
 effects over the time steps (t). The t_s is the starting of
 treatment.

380

382 The main success of the presented CA model is the ad-
 383 equate modeling of the four-phases of HIV infection with
 384 different time scales into one model. Moreover, we could
 385 also integrate all of the three different therapy procedures.
 386 The simulations show a qualitative correspondence to clin-
 387 ical data. During the phase of drug therapy response, tem-
 388 poral fluctuations for $N > 3$ were observed, this is due to
 389 the relative simple form of the response distribution func-
 390 tion (P_{dis}) applied to the drug effectiveness parameter N
 391 at each time-step. The simulation results indicate that, in
 392 contrast to ODE/PDE, our model supports a more flexible
 393 approach to mimic different therapies through the use of
 394 mapping the parameter space of P_{dis} to clinical data. There-
 395 fore there is ample room to incorporate biologically more
 396 relevant response functions into the model. The data inte-
 397 gration required for the CA, the parametric computation
 398 and the data presentation are supported by the Grid.

399 2.1.2. Multivariate stochastic modeling

400 The modeling of Human Immunodeficiency Virus
 401 (HIV-1) genotype datasets has a goal to identify patterns
 402 of mutations (or naturally occurring polymorphisms) as-
 403 sociated with resistance to antiviral drugs and to predict
 404 the degree of *in-vitro* or *in-vivo* sensitivity to available drugs
 405 from an HIV-1 genetic sequence. The statistical challenges
 406 in doing such analyses arise from the high dimensionality
 407 of these data. Direct application of the well-known genetic
 408 approaches [5] to analysis of HIV-1 genotype results in a lot
 409 of problems. Principal difference is in the fact that, in HIV

DNA analysis, the main scope of interests is the so-called 410
 relevant mutations – a set of mutations, associated with the 411
 drug resistance. These mutations might exist in different 412
 positions over the amino-acid chains. Moreover, the sheer 413
 complexity of the disease and data require the development 414
 of the reliable statistical technique for its analysis and mod- 415
 eling. A multivariate stochastic model for describing the 416
 dynamics of complex non-numerical ensembles, such as 417
 observed in the (HIV) genome, has been developed in [6]. 418
 This model was based on principle component analyses for 419
 numerated variables. Generally speaking, the interpretation 420
 of numerated variables in terms of relevant mutations is not 421
 clear. Below we develop this model directly for the ensemble 422
 of relevant mutations in the RT and protease parts of 423
 the HIV-1 genome. Each element of the ensemble is pre- 424
 sented as the cortege $\Xi_k = \{\xi_j\}_{j=1}^{n_k}$, $k = \overline{1, M}$ with the 425
 variable dimension n_k – the total number of the mutations 426
 in the gene. Each value ξ_k is a literal index and corresponds 427
 the position and new value of the amino acid (e.g., 184 V, 428
 77I, etc.). It allows to associate each mutation with the cat- 429
 egorical random variable $i \in 1 \dots K$, where K is the total 430
 number of possible mutations. Each sub sample of genomes 431
 with a fixed number of mutations $n = \text{const}$ may be con- 432
 sidered as the realizations of a categorical random vector. 433

The representation above is based on the proximity to the 434
 “wild-type” virus and takes into account only the relevant 435
 mutations in a genome. It allows for significant compression 436
 of the DNA representation and simplifies the interpretation 437
 of the results. 438

Principle of the modeling approach. The joint variability of dif- 439
 ferent mutations in the HIV-1 genomes is a complicated 440
 phenomenon. The dimension of the probabilistic charac- 441
 teristics is high, and its analytical investigations and inter- 442
 pretation are hard. Hence, for the studying of HIV-1 pop- 443
 ulations we use a computational statistical approach that 444
 allows to numerically generate an ensemble with the same 445
 probabilistic properties by means of a Monte-Carlo pro- 446
 cedure. This is a well-known powerful method to study 447
 complex system variability. 448

The idea of the stochastic modeling is shown in the 449
 Figure 5. It is based on the evolutionary hypothesis, consid- 450
 ering the group with $n + 1$ mutations as subgroup of group 451
 with n mutations in a previous step. For each gene the tran- 452
 sit from n to $n + 1$ mutation groups is driven by a stochastic 453
 operator $D_{(n+1)}$, which defines the mutations on the $n + 1$ 454
 step, when the mutations on the previous n steps are known. 455
 The initial step of the stochastic procedure begins from the 456
 whole ensemble of wild-type viruses. The number of the 457
 genomes that has been mutated at each step of the stochas- 458
 tic procedure is in accordance with $M_n = \rho_n M$, where ρ_n 459
 are the probabilities of the occurrence of genotypes with n 460
 mutations in a total population of M genes. 461

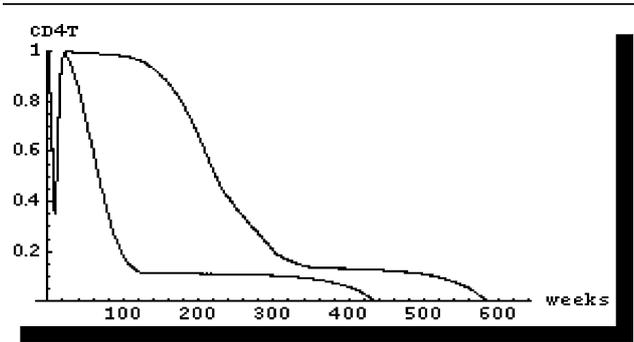


Fig. 3. Temporal behaviour of the CD4 count, with modeled Brownian movement for lymphocytes [8].

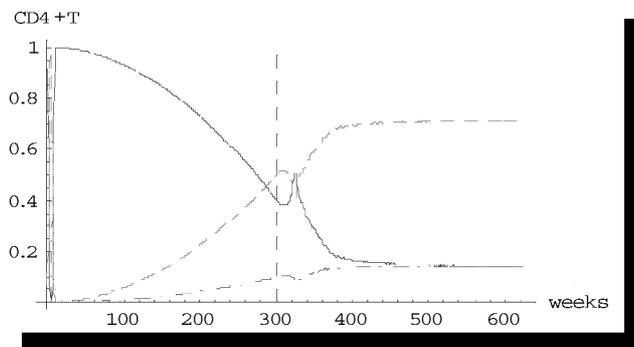


Fig. 4. As in Figure 3, with additionally modeled mono therapy in week 300 [8].

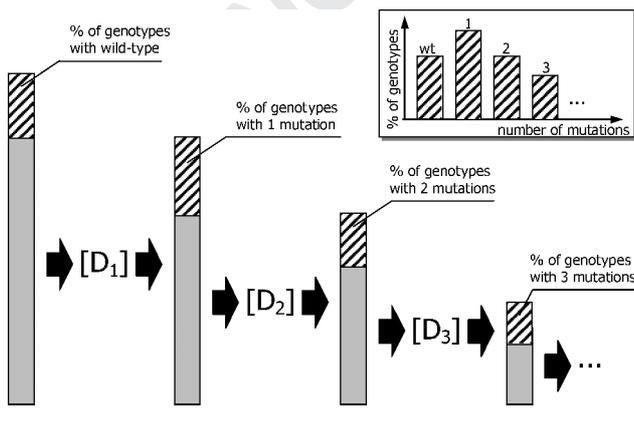


Fig. 5. Principle of the modeling.

462 The stochastic operator D may be considered as a “black
 463 box”. It is formalized in terms of the conditional probabil-
 464 ities of the occurrence of mutation ξ_i , if the mutation ξ_j
 465 arise in the previous step of the generation. For genotypes
 466 with 2 mutations only the values D_{ij} are the conditional
 467 probabilities of the pairs. In this case the matrix $\{D_{ij}\}$ is

the transition Markov probability matrix, containing the 468
 conditional probabilities for simple Markov chains with 469
 the number of these states corresponding to quantity of 470
 the relevant mutations. In more complicate cases, where 471
 $n > 2$, the probability matrix $\{D_{ij}\}$ consists of the con- 472
 ditional probabilities to meet mutation ξ_j in certain gene, 473
 when the mutation ξ_i is present. 474

This approach allows us to reduce the complicated sta- 475
 tistical description of the dataset to a rather simple model, 476
 using only three probabilistic distributions as the initial pa- 477
 rameters of the model: distribution of number n of the 478
 mutations ρ_n ; 479

- distribution $P_{\xi}^{(1)}$ for the relevant mutations in the group 480
 $n = 1$; 481
- transient probability matrix D . 482

All these parameters might be identified on the sample 483
 datasets of the HIV-1 population. 484

Identification of the model. For the identification of parameters 485
 of the model, a large database of HIV-infected patients, col- 486
 lected over several years in USA, is used [4]. These databases 487
 contain genotypes of 43620 patients examined from Au- 488
 gust 9, 1998 to May 5, 2001. We observed 59 different 489
 mutations in the RT genome, including 17 mixed muta- 490
 tions, and 77 different mutations in the protease genome, 491
 including 34 mixed mutations. 492

Distribution ρ_n of number of mutations. The practice of HIV 493
 treatment however, has shown that the variability of the 494
 number of mutations n is high, due to the complexity of 495
 the drug combinations that has been applied. The sample 496
 estimate of distribution ρ_n of the number of mutations in 497
 protease is shown in the Figure 6. It is seen, that the distri- 498
 butions have a clear first peak ($n = 1$), and a shelf (or second 499
 peak), corresponding to $n = 3 \div 5$. Therefore we expect 500
 that there are two groups of genomes in the database, cor- 501
 responding to the low and high number of mutations. The 502
 possible interpretation of the discovered bi-modal distri- 503
 bution is that we have two groups of patients. One group 504
 is the “new” patients who had one or two treatments, thus 505
 their genotype contains relative small numbers of muta- 506
 tions. The second group is the “old” patients, which have 507
 a long treatment history, or new patients, infected through 508
 treated HIV-1 patients [15]. 509

Distributions of the relevant mutations P_{ξ} . Distribution ρ_n al- 510
 lows describe the variability of the groups of the “new” 511
 and “old” patients, only. For a more detailed study of the 512
 virus mutations driving by the certain drugs combinations, 513
 the probabilities of occurrence of the relevant mutations 514
 ξ should be considered. They are estimated by the sample 515

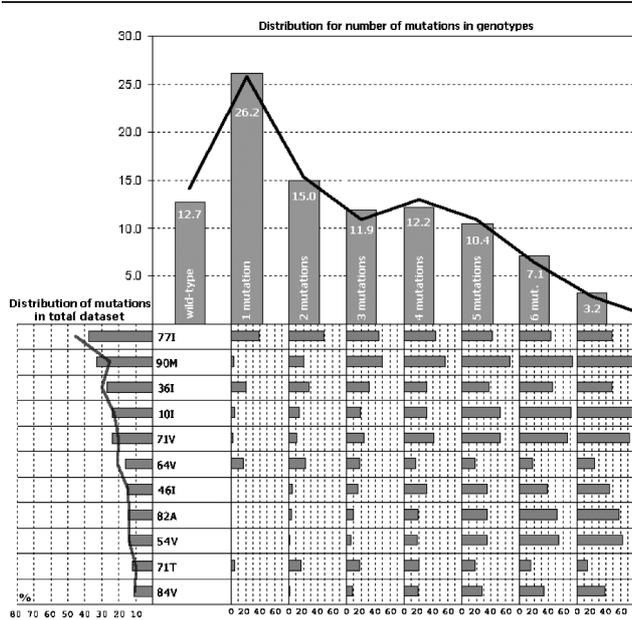


Fig. 6. Statistical description for distribution of mutations in Protease.

516 frequencies:

$$P_{\xi} = \frac{\{\text{Number of genes with mutation } \xi\}}{M} \quad (1)$$

517 Here M is the total number of genomes in the dataset.
 518 Equation (1) describes the marginal impact of each muta-
 519 tion in the total population, without any information about
 520 number and occurrences of other mutations. The proba-
 521 bilities of the most significant relevant mutations ξ_k (in
 522 decreasing order of its probability) are shown in Figure 6.
 523 The marginal estimates of P_{ξ} over the total dataset show
 524 only general impacts of the mutations. For a detailed
 525 analysis of its behavior we also consider the occurrences
 526 $P_{\xi}^{(n)}$ of mutations in the groups of genotypes with exactly
 527 n mutations. These values were computed also by means
 528 of Equation (1), where $M \stackrel{\text{def}}{=} M_n = \rho_n M$ – the number of
 529 genes with n mutations in a database. The sample estimates
 530 of these occurrences are also shown in the Figure 1. It is
 531 clearly seen that the inputs of some mutations are rather dif-
 532 ferent for different n , both for the protease and RT parts of
 533 the genome. E.g., for RT, for $n = 1$, the mutations 184V
 534 and 103N have the main input. The distribution $P_{\xi}^{(1)}$ is the
 535 limit distribution from the procedure shown in Figure 5.

536 From Figure 1 we also observe that the total sum
 537 $\sum_k P_{\xi_k} > 100\%$, excluding case $n = 1$. This demonstrates
 538 that the analysis of the marginal mutations is not enough
 539 for general statistical description of all DNA ensemble vari-
 540 ability, because some positions of DNA may be statistically
 541 dependent [15], especially in relation to viral fitness. Hence,

the joint characteristics of its variability must be taking into 542
 account. 543

Transient probability matrix D . The conditional probability of 544
 the occurrence of mutation ξ_i , if the mutation ξ_j arises 545
 from the previous steps of the generation, is estimated by: 546

$$D_{ij} = \frac{\{\text{Number of genotypes with mutations } \xi_i \text{ and } \xi_j \text{ simult\&neously}\}}{\{\text{Number of genotypes with mutation } \xi_j\}} \quad (2)$$

547
 The dimensionality of the related matrix, obtained from 548
 Equation (2), may be rather high. In order to decrease 549
 the dimensionality we consider the algebraic technique of 550
 orthogonal expansion, applied to transient probability ma- 551
 trices [16]. 552

$$D = \Phi \Lambda^{1/2} \Psi. \quad (3)$$

where Φ are the eigenvectors of matrix DD^T , and Ψ – of 553
 matrix $D^T D$. It allows considering the coefficients $a_k =$ 554
 $\sqrt{\lambda_k}$ as the principal components (PC) [13], and represents 555
 the probability (2) as a series: 556

$$D_{ij} = \sum_k \sqrt{\lambda_k} \phi_{ik} \psi_{jk}. \quad (4)$$

The values λ_k shows the part of the probability, explained 557
 by k -th PC. The sum of the first k -th coefficients λ_k may 558
 be interpreted as a measure of convergence of the series 559
 (4). In Table 2 the values of the first 7 λ_k for the RT and 560
 protease parts of the HIV-1 genome are shown. These data 561
 were obtained for the total database. It can be seen that the 562
 series (4) converges rather fast in both cases: e.g. for the RT 563
 part only the first term of the series explain more 60% of 564
 conditional probability (the first five terms explain 80%). 565

Let us consider the normalized bases $\tilde{\phi}_{ik} = \lambda_k^{0.25}$ 566
 ϕ_{ik} , $\tilde{\psi}_{jk} = \lambda_k^{0.25} \psi_{jk}$. It allows to present the terms in Equa- 567
 tion (4) as the $p_k^{ij} = \tilde{\phi}_{ik} \tilde{\psi}_{jk}$ and interpreted these values 568
 as the independent factor loadings, driving the changes of 569
 the conditional probability D_{ij} over all the mutations ξ_i , ξ_j 570
 in the database. For example, in the Figure 7 the estimates 571

Table 2. Normalized (%) values of the expansion coefficients λ_k in Equation (4)

Part of the genome	# of PC						
	1	2	3	4	5	6	7
RT	61.3	8.2	5.4	2.8	2.1	1.7	1.6
Protease	55.0	6.3	4.5	4.2	3.4	2.7	2.4

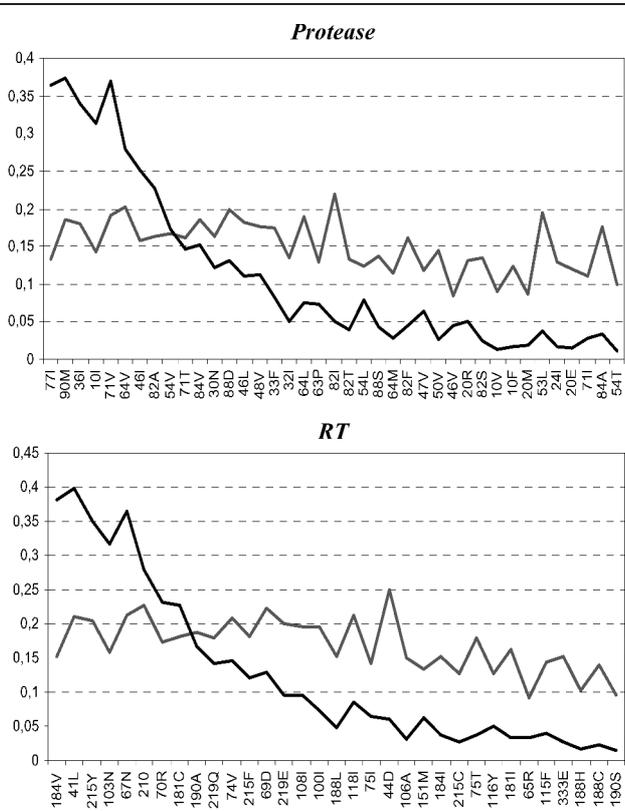


Fig. 7. Orthogonal basic functions of expansion (4) for transient probability matrix.

572 of the first basic functions are shown for RT and protease
 573 parts of the genotype (the input of multiplication of functions
 574 are in the Table 2). It is clearly seen, that the first
 575 term $p_1^{ij} = \tilde{\phi}_{i1}\tilde{\psi}_{j1}$ reflects the total occurrence of the muta-
 576 tions in a genotype (see Figure 6): for the mutations with
 577 the maximal occurrences the input to conditional proba-
 578 bilities of its pairs is also high.

579 *Model validation.* The simulation model is based on the
 580 ρ_n , $P_\xi^{(1)}$, D distributions of the mutations only. No infor-
 581 mation of more complicate mechanisms (distributions of
 582 pairs, triples, etc.) has been used for this identification.

583 The main goal of the verification is the possibility to
 584 reproduce these features of the ensemble through the de-
 585 pendencies formalizing the matrix D . We compared the
 586 total occurrences of all mutations in genotypes, estimated
 587 on the initial and simulated samples, see also Figure 6 (solid
 588 line). It is seen, that the results of the simulation and sample
 589 are rather close.

590 The error of the simulation increases proportionally to
 591 absolute value of the occurrences. Nevertheless, for some
 592 cases the error of the simulation is larger then the boundary

of the confidence interval. This systematic error may be 593
 explained by possible variations in matrix D for groups of 594
 the “old” and “new” patients. 595

Application to forecast of HIV-1 evolution in time. The evolu- 596
 tion of total world populations of HIV-1 and the associ- 597
 ated changing of the related drug resistance levels should 598
 be taken into account. The stochastic models, used to de- 599
 scribe the HIV-1 genotype ensemble in terms of parame- 600
 ters and shown in the Figure 5, can be used for the analysis 601
 of its temporal variability during the observation period 602
 (VIII.1998–V.2001). The temporal variability of the data 603
 may be considered in terms of the samples of the seasons 604
 (3-months periods). The volumes of seasonal samples are 605
 from 1500 till 4500 genotypes; that is enough for obtain- 606
 ing the stable estimations. Only the hypothesis of linear 607
 trends is considered: $\xi(t) = at + b + \delta(t)$, where a is the 608
 most interesting parameter—value of the trend, b is the 609
 shift parameter, and δ is the white noise. In the Table 3 the 610
 integral parameters of trends of the various parameters of 611
 the HIV-1 population (mean value of the parameter, value 612
 of the trend, determination coefficient R^2 and the sample 613
 value of F -criterion) are shown. 614

Trends of single mutations occurrence P_ξ . The database allowed 615
 us to investigate trends in codon frequency in the period 616
 of 1998 till 2001. Results for Protease and RT are shown 617
 in Table 3. The majority of the mutations in the genotype 618
 have a negative trend, only 771 in Protease has significant 619
 positive trend. 620

*Trends of bi-modal distribution for number of mutations in geno- 621
 types ρ_n .* For the decreasing of the data dimensionality and 622
 the statistical discrimination of two groups in the dataset 623
 we consider the model of the mixture of two Bernoulli 624
 distributions: 625

$$\rho_n = p_g C_{m_1}^k q_1^k (1 - q_1)^{m_1 - k} + (1 - p_g) C_{m_2}^k q_2^k (1 - q_2)^{m_2 - k} \quad (5)$$

where p_g is an input of the first group of mutations (and 626
 p_g is an input of the second group, m_1, m_2 —are maximal 627
 numbers of mutations in groups and q_1, q_2 —are probabil- 628
 ities to find each one (arbitrary) mutation in the groups. 629
 The use of Bernoulli distribution logic (based on the rep- 630
 etition of the independent events) is more close to the 631
 description of the mutation process, then the Poisson dis- 632
 tribution, generally applying to description of rare events. 633
 Temporal variability of the parameters (p, q_1, q_2, m_1, m_2) _{t} 634
 of the ρ_n approximation by Equation (5) are shown in 635
 Table 3. In both cases only the parameter p_g (weight 636
 of the left part for group of m_1 mutations) has a clear 637

Table 3. Trend analysis of the parameters of the HIV-1 genotype population (F is compared with Fisher's test $F(1,31,95\%) = 4.14$)

Parameter	Occurrence of mutations, %				p_g , %, Equation (5)	Coefficients $\sqrt{\lambda_k}$, Equation (4)		
	77I	90M	10I	71V		$k = 1$	$k = 2$	$k = 3$
Protease part								
Mean	37.78	32.69	27.97	23.64	48	5.78	1.67	0.83
a (1/month)	0.20	-0.43	-0.72	0.32	0.74	0.13	0.06	0.06
R_2	0.68	0.91	0.61	0.82	0.67	0.80	0.73	0.54
F	16.7	77.6	9.6	47.1	64.0	23.6	26.8	11.8
RT part								
	41L	215Y	103N	67N		$k = 1$	$k = 2$	$k = 3$
Mean	32.86	31.37	30.66	27.21	47	6.65	2.20	2.08
a (1/month)	-0.51	-0.50	-0.32	-0.39	0.49	0.11	0.17	0.07
R^2	0.88	0.93	0.88	0.84	0.75	0.68	0.78	0.71
F	57.4	98.7	59.8	41.8	94.3	21.4	36.1	25.3

638 significant positive trend. For protease value p_g increased
 639 from 39% in Summer, 1998 to 62% in Summer 2001
 640 (with average increment $a = 0.74\%$ per month). Taking
 641 into account trends for separate mutations we observed a
 642 “degradation” of genotypes: the number of patients with
 643 simple genotypes (small number of mutations) is growing
 644 but a number of patients with big count of mutations is
 645 decreased.

646 Trends of transient probabilities D . The analysis of the trends of
 647 parameters for distribution (1) shows that the input of the
 648 first group of mutations with low number n is increased.
 649 Hence, it may be a consequence of the temporal variations
 650 of the interdependencies between different mutations, gov-
 651 erned by the developing of the drug therapy. For the anal-
 652 ysis of these hypothesis, let us consider the trends for the
 653 matrix D , Equation (2). Taking into account the expan-
 654 sions (3, 4), we may reduce the complicate problem for
 655 joint trend analysis for components D_{ij} to the procedure
 656 of trend analysis for independent time series – components
 657 of expansions (4). From the Table 3 it can be seen, that all
 658 the components have a clear positive trends. Taking into
 659 account the shape of first bases functions, see Figure 7, it is
 660 clear, that generally the joint probabilities D_{ij} of the mu-
 661 tations is increased also; moreover, the power of increasing
 662 corresponds to the total occurrences of the mutation in the
 663 ensemble.

664 The discrimination of the groups of “old” and “new”
 665 patients in terms of bi-modal distribution (5) allow to fore-
 666 cast the growth of the total number of HIV-infected people
 667 in time:

$$N(t) = N_{\text{patients}}^{\text{new}}(t) + N_{\text{patients}}^{\text{old}}(\varepsilon t), \varepsilon \ll 1. \quad (6)$$

Here ε – is the slow time parameter, which shows the rapid
 668 increasing of the new patients group in comparison with
 669 the old patients. The part of “new” patients of the sample
 670 is p_g (old patients – $(1 - p_g)$) from (5). Hence, the growth
 671 curve is:
 672

$$N(t) = N_{\text{patients}}^{\text{old}}(0) \left[1 + \frac{p_g(t)}{1 - p_g(t)} \right], \quad (7)$$

where $p_g(t) = p_0 + a_g t$ – is the linear trend with the pa-
 673 rameters from Table 3, and $N_{\text{patients}}^{\text{old}}(0)$ is the initial value of
 674 “old” (treated) patients on the beginning of the forecast.
 675

In Figure 8 the “crucial” forecast of the HIV-1 popula-
 676 tion growth are shown. It is based on the fact that altogether
 677

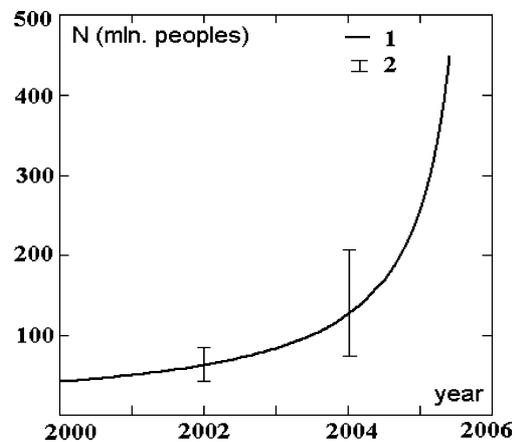


Fig. 8. Qualitative forecast of HIV-1 population grows. 1 – mean value (7), 2 – 90% confidence interval.

678 42 million people worldwide have been infected with HIV
 679 at the beginning of XXI century, and 12 million have died
 680 over the last 20 years. Moreover, not taken into account
 681 is the arising of new drugs and different prophylactic and
 682 social preventive activities for restriction of HIV-1 infec-
 683 tion. Really, this result is qualitative only; for quantita-
 684 tive conclusions the more sophisticated research should be
 685 done.

686 **3. RESULTS**

687 *3.1. Data presentation: Roaming PDA access*

688 *3.1.1. User Scenario*

689 RetroGram™ (www.retrogram.com) is a unique HIV-
 690 genotype expert based interpretation software program,
 691 which weighs the effect of specified genotype changes on
 692 clinical drug activity. It accepts a list of substitutions to the
 693 protease and reverse transcriptase genes with respect to the
 694 NL4-3 reference strain. This is accomplished by running
 695 a “simulation”, which applies some hundred rules relat-
 696 ing substitutions on the HIV genome to knowledge of
 697 effects on drug response. The latter comes from over hun-
 698 dreds of references from the clinical literature. The rules are
 699 checked against the reported substitutions, and each drug is
 700 evaluated for its suitability. In a later stage we added Web-
 701 access where a Web interface is used to submit the input
 702 and take out the output. We want to make the simulations
 703 wireless-accessible. Developing a wireless Internet version
 704 from scratch will not be cost-efficient and causes maintain-
 705 ability problems. For example, the rules mentioned above
 706 are often changed and these changes have to be reflected in
 707 both versions. Furthermore, for privacy and security rea-
 708 sons the developer is not granted access to the source code
 709 of the “simulation”. Thus, it is much more convenient to
 710 have wireless access to the Web-based interface. In this case
 711 the “simulation” take places in a unique server and privacy
 712 and security are guaranteed. A typical user scenario is de-
 713 scribed below and the associated graphical representation
 714 of the Retrogram Web access is given in Figure 9.

715 After the user has successfully logged in, the *Patient Detail*
 716 page is displayed (Figure 10). The form, taking place in
 717 this page is used to enter the personal data of the patient.
 718 Two fields are required in the form, *Patient ID* and *Data of*
 719 *Sample*.

720 According to the information taken from the laboratory
 721 the user enters the laboratory test results (i.e. Protease or
 722 RT substitutions) for the patient in the *Laboratory Informa-*
 723 *tion* page. Next a script invoked on the server does the
 724 following:

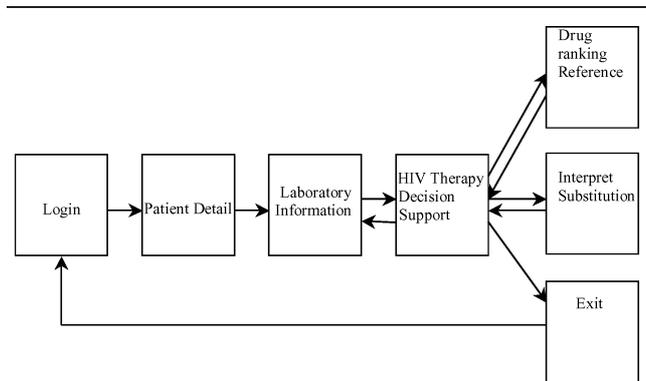


Fig. 9. Web-based Retrogram use case sequence.

Script 1: Server validation script

725

Validate inputs:
 Validate Protease or RT substitutions if they conform
 to certain rules.
 A single substitution should be represented by an
 integer (for position in the gene) and a letter (for the
 amino acid). The position in the gene is in the range
 from 1 to 99 for Protease position and from 1 to 599
 for RT position. The amino acid code is one of the
 following codes: A C D E F G H I K L M N P Q R
 S T U V W Y.
 Submit the inputs to the “simulation” program and
 take back the drugs ranking result.
 Show the Drugs ranking result in the ‘*HIV Therapy*
decision support’ screen:
 After applying certain rules on the laboratory test
 result return to the final drugs ranking or drug’s level
 of suitability indication as follows:
 A (green): This drug can be used
 B (yellow): Consider use if no class A drug available
 C (amber): Consider use if no class A or B drug
 available
 D (red): Consider use if no class A, B or C drug
 available
 U (grey): Unranked, insufficient data available

726

In the ‘*HIV Therapy decision support*’ screen, clicking on
 any drug name in the ranking lists will display a list of avail-
 able references from the scientific literature supporting the
 particular ranking for that drug. In the ‘*HIV Therapy*
decision support’ screen, clicking on the ‘*Interpret substitution*’
 button will show classification of the patient’s substitutions
 into *relevant, natural or additional*.

732
733

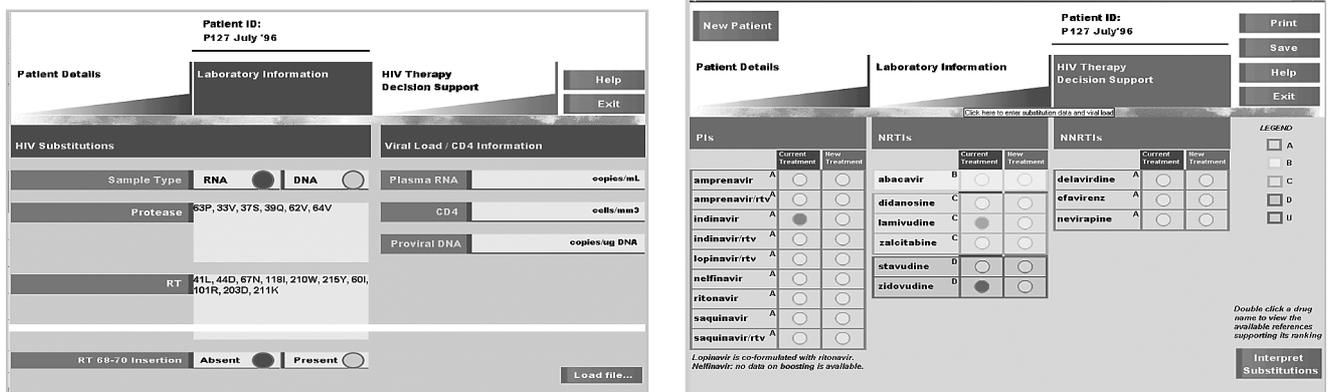


Fig. 10. Web Retrogram: user enters patient substitutions (left), drug ranking results (right).

734 3.1.2. Roaming, wireless access

735 In the designing phase of wireless versions of the application
 736 the constraints of the mobile devices should be considered.
 737 At the same time we have tried to maintain the same level
 738 of usability and readability as in the original Web version.
 739 This is accomplished by maintaining the same structure as
 740 that in the Web but with some modifications. For example,
 741 the Patient detail form has many fields and putting them
 742 in one screen would cause problems in the usability of
 743 the program (it's supposed that the mobile device has a
 744 resolution comparable to a normal PDA, i.e., something
 745 around 160 × 160 pixels). Thus we use three screens for
 746 Patient Detail data. The Patient Detail Web page has 2
 747 required fields. We put them in the first screen after the
 748 'login' screen. In this way, if the user is not interested in
 749 entering optional data, she can directly go to the Laboratory
 750 Information.

751 *Proxy method Implementation.* A Proxy method is imple-
 752 mented for accessing the web-based software from mobile
 753 devices. The Proxy server takes places between the remote
 754 server (the Retrogram server) and the mobile device. A
 755 *mininavigator* script developed in the Proxy is responsible
 756 for the following:

- 757 • Take the patient data from the mobile user (i.e. patient
- 758 detail, laboratory information)
- 759 • Create an HTTP communication with the remote
- 760 server,
- 761 • Submit data to the remote server. These data are basically
- 762 the input for the Retrogram 'simulation'.
- 763 • Take the result from the remote server (HTML code
- 764 generated from retrogram.asp script),
- 765 • Parse HTML code and retrieve only relevant informa-
- 766 tion (i.e. drug ranking, error messages, drug references

etc.). It uses this relevant information to build wireless 767
 pages (i.e. WML page in case of WAP or Web-clipping 768
 page). 769

- Send the wireless pages to the mobile device. 770

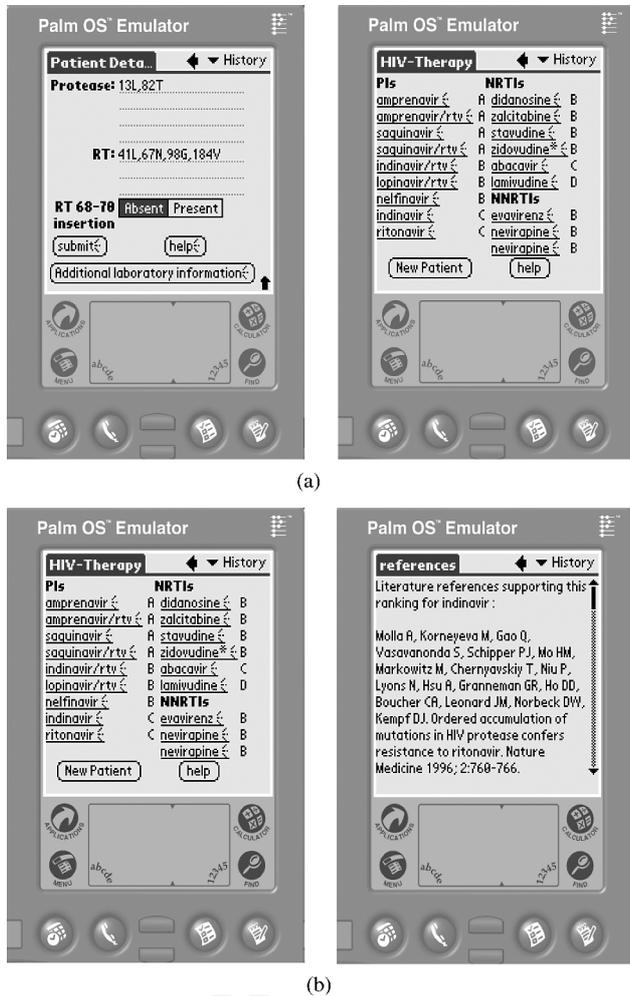
The Proxy is implemented using PHP: Hypertext Pre- 771
 processor as a server-site scripting language [9–11] running 772
 on the Apache Web server [12]. 773

Two versions are developed using the Proxy method: 774
 WAP version and web clipping. If a user wants to enter the 775
 'patient details' fields, he has to move from one screen to 776
 the other and come back again. The fields already filled in 777
 the previous screens should not be lost. Thus maintaining 778
 the client's state is necessary. In the WAP case we simply 779
 use cookies but in web clipping cookies are supported only 780
 in PALM OS 4.0 version or higher. For this reason the 781
 "hidden field" method is used this is another method used 782
 for maintaining state in the Internet. The following figures 783
 are the user interfaces that have been captured. They track 784
 the user's path through the running of the application, as 785
 shown in Figures 11(a) and 11(b), where the user enters 786
 the patient's details and accesses ranking results. 787

J2ME Implementation. The same user interface is applied in 788
 the J2ME implementation. There are two main differences 789
 between the J2ME implementation and the Proxy one: 790

1. J2ME enables the device to communicate directly to 791
 the Retrogram server without an intermediate Proxy 792
2. In J2ME the client's interface is contained within the 793
 device. In the Proxy method, every time the interface 794
 should be changed, the Proxy is responsible for gener- 795
 ating a new page. 796

The following illustrates the necessary steps one should 797
 take in order to fetch an HTML page generated from a 798



(a)

(b)

Fig. 11. (a) User corrects the input and submit again (left), drug ranking results (right). (b) Users clicks to the drug 'indinavir' (left), references supporting this ranking (right).

799 script in the remote host. Specifically this is an example
 800 illustrating how the user can login to a script in the Ret-
 801 rogram server and extract the cookie from the header re-
 802 sponse:

- 803 1. Open an HTTP connection
- 804 2. Open an input stream
- 805 3. Make an HTTP POST request
- 806 4. Extract the cookie from the header response
- 807 5. Close the connection

808 In the J2ME implementation of Retrogram the entire
 809 client's interface takes places in the device. The connec-
 810 tion to the server is established in the following cases: user
 811 login, with connection with the server is necessary in order
 812 to validate the user and/or password. The user submits the

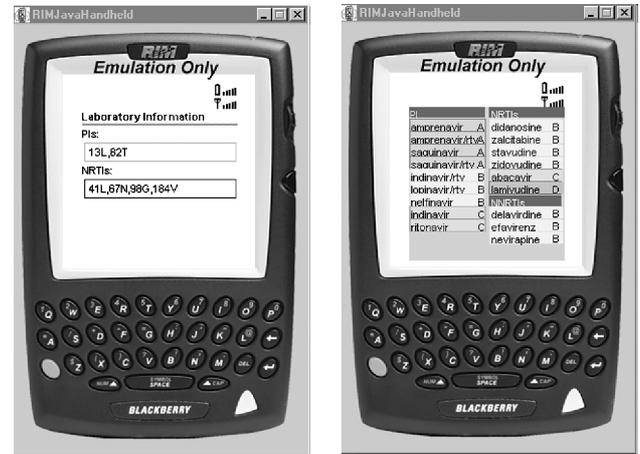


Fig. 12. J2ME method; user enters patient's substitutions (left), drug ranking results (right).

813 username and password, and the application judges them
 814 for their correctness by scanning the HTML response from
 815 the Retrogram server. The user submits the patient's labo-
 816 ratory information data. The application should connect
 817 to the server in order to submit the data, take the result
 818 (HTML format) and extract the drugs ranking. Next the
 819 user looks for the references that suggest a certain drug
 820 ranking. The database with all the references exists in the
 821 Retrogram server, therefore the connection is necessary.
 822 The application submits to a Retrogram script the cookie
 823 and the name of the drug. The drug references are given
 824 back from the server in HTML format. The application
 825 should clean up the HTML tags and show the references
 826 as plain text. Finally the user looks for classification of the
 827 patient's substitutions. This classification is part of the Ret-
 828 rogram 'simulation' and thus the connection to the server
 829 is still necessary. In Figure 12 we illustrate the process of
 830 taking the drugs ranking using the J2ME method.

831 Currently we have the J2ME version in use for different
 832 users to study the usability and extendibility. More details
 833 on the implementation can be found in reference [13].

3.2. Virtual laboratory infrastructure 834

3.2.1. A virtual organization for retrogram-centered workflow 835

836 Grid technology is a major cornerstone of today's com-
 837 putational science and engineering, with its basic unit of
 838 Grid organization called the Virtual Organization (VO).
 839 A VO is a set of Grid entities, such as individuals, appli-
 840 cations, services or resources, which are related to each
 841 other by some level of trust. In the most basic example,

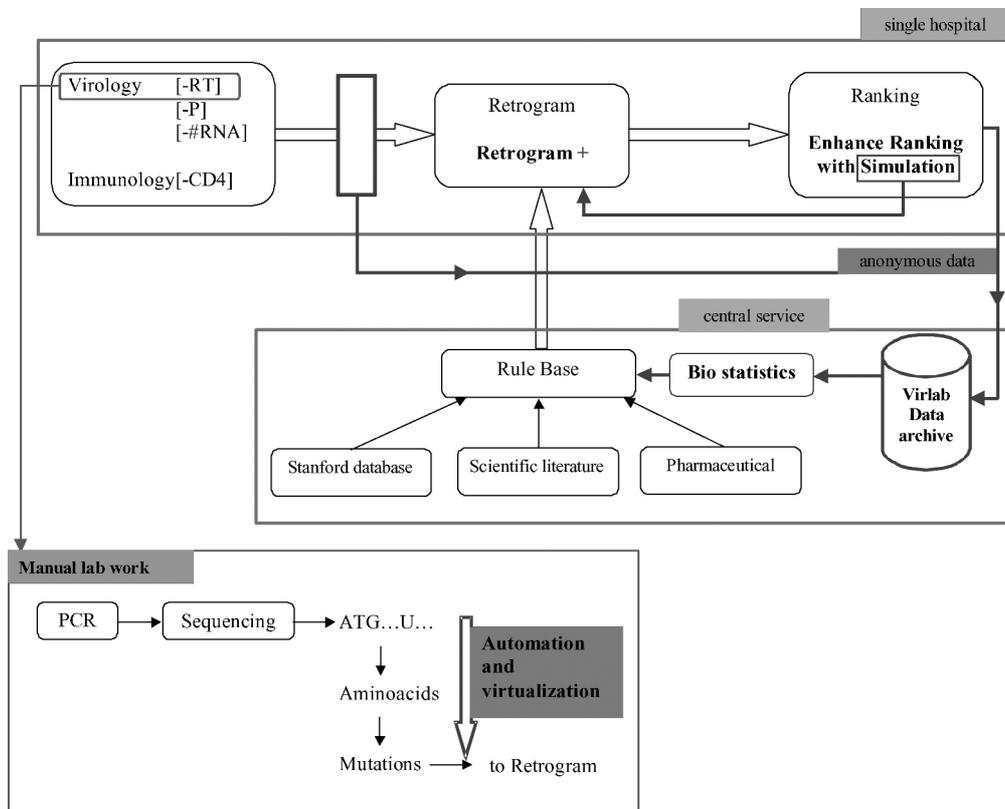


Fig. 13. A Retrogram-centered workflow.

842 service providers would only allow access to the mem- 865
 843 bers of the same VO. We are currently building a dis- 866
 844 tributed Grid-based overall decision support infrastruc- 867
 845 ture to support the Retrogram-centered workflow shown in 868
 846 Figure 13. 869

847 This VO will offer a Grid virtual laboratory that will 870
 848 assist users in the interpretation the genotype of a patient 871
 849 by using rules developed by experts on the basis of the liter- 872
 850 ature, taking into account the relationship between the 873
 851 genotype and phenotype. The workflow is based on highly 874
 852 distributed available data from clinical studies and on the 875
 853 relationship between the presence of genotype and the clin- 876
 854 ical outcome. In order to cover the fast temporal and spatial 877
 855 scales required to infer information from a molecular (ge- 878
 856 nomic) level up to patient medical data multi-scale methods 879
 857 are applied, where simulation, statistical analysis and data 880
 858 mining are combined and used to enhance the rule-based 881
 859 decision. In this scenario, information sources are widely 882
 860 distributed, and the data processing requirements are highly 883
 861 variable, both in the type of resources required and the pro- 884
 862 cessing demands. Experiment design, integration of infor-
 863 mation from various sources, as well as transparent schedul-
 864 ing and execution of experiments will be supported by this

865 support system based on distributed Grid middleware. The 866
 867 DAS2 testbed (Netherlands) will initially provide the addi- 868
 869 tional computational power for our compute intensive jobs. 869
 870 We will reuse Grid middleware from successful European 870
 871 projects such as CrossGrid (www.crossGrid.org) and VL-e 871
 872 (www.vl-e.nl) to provide basic Grid services of data man- 872
 873 agement, resource management, and information services 873
 874 on top of Globus. For transparent use of this infrastructure 874
 875 we will build a presentation layer that will provide a user-
 876 friendly interface to both medical doctors and scientists.

4. DISCUSSION

4.1. Conclusions and future work

877 In this paper we discussed an integrative approach to bio- 877
 878 medicine at large and to infectious diseases in particular. 878
 879 We showed how in the understanding of processes ‘from 879
 880 molecule to man’ Grid technology can play a crucial role. 880
 881 In order to cover the fast time and spatial scales required to 881
 882 infer information from a molecular (genomic) level up to 882

883 patient medical data, we need to apply multi-scale meth-
 884 ods where simulation, statistical analysis, data-mining is
 885 combined in an efficient way. Moreover the required in-
 886 tegrative approach asks for distributed data collection (e.g.
 887 HIV mutation databases, patient data, literature reports etc.)
 888 and a virtual organization (physicians, hospital administra-
 889 tion, computational resources etc.). Also the access to and
 890 use of large-scale computation (both high performance as
 891 well as distributed) is essential since many of the compu-
 892 tations involved require near real-time response and are
 893 to complex to run on a personal computer or PDA. Fi-
 894 nally data presentation is crucial in order to lower the
 895 barrier of actual usage by the physicians, here the Grid
 896 technology (server-client approach) can play an important
 897 role.

898 Although many of the aspects discussed in this paper
 899 have proven to work in concept, the complete integration
 900 of the systems and the evaluation of day-to-day use is
 901 still under development [17]. In addition each of the
 902 underlying methods (Rule-based, statistical and CA based
 903 models) remain topics of further studies. We will set up a
 904 use-base with the system described running under various
 905 European Grid testbeds. The first testbed we will use is
 906 the so-called DAS2, and eventually the CrossGrid testbed,
 907 which supports specific features for interactive computa-
 908 tion, an essential ingredient for a medical decision support
 909 system.

910 The authors gratefully acknowledge Fan Chen and Ferdinand
 911 Alimadhi for assistance in implementing the CA models and
 912 the roaming PDA access. The Dutch Virtual Laboratory on e-
 913 science project supported parts of the research presented here:
 914 <http://www.VL-e.nl>.

915 GLOSSARY

916 Grid: Distributed architecture for solving computational
 917 problems by making use of the resources from the mem-
 918 bers of a virtual organization, treating them as a virtual
 919 cluster.

920 CA: Cellular Automata, a discrete model studied in com-
 921 putational theory and mathematics, which consists of
 922 regular grid of cells, each in one of a finite number of
 923 states.

924 Decision Support System: Computer-based system that
 925 helps in the process of decision-making.

926 Web Interface: User interfaces for information available via
 927 the web.

928 Proxy: Computer service which allows clients to make in-
 929 direct network connections to other services.

HTTP: Hyper Text Transfer Protocol, a request/response
 930 protocol for transferring information on the Web. 931

HTML: Hyper Text Markup Language, a markup language
 932 designed for the creation of web pages. 933

WML: Wireless Markup Language, a markup language
 934 used in mobile phones. 935

J2ME: Java 2 Platform Micro Edition, a collection of Java
 936 interfaces for embedded consumer appliances such as
 937 cellular phones. 938

DAS2: Distributed ASCI Super Computer 2, a wide-area
 939 distributed computer connecting 5 Dutch Universities. 940

REFERENCES

1. Zhao Z, Belleman RG, van Albada GD, Sloot PMA. AG-IVE: An Agent-Based Solution to Constructing Interactive Simulation Systems, in Series Lecture Notes in Computer Science, April 2002; 2329: 693–703. 941–945
2. Durant J, Clevenbergh P, Halfon P, Delguidice P, Porsin S, Simonet P, Montagne N, Dohin E, Schapiro JM, Boucher C, Dellamonica P. Improving HIV therapy with drug resistance genotyping: The Viradapt Study. *Lancet* 1999; 353: 2195–2199. 946–950
3. Sevin AD, DeGruttola, Nijhuis M, Schapiro JM, Foulkes AS, Para MF, Boucher CAB. Methods for Investigation of the Relationship between Drug-Susceptibility Phenotype and Human Immunodeficiency Virus Type 1 Genotype with Applications to AIDS Clinical Trials Groupw 333. *The Journal of Infectious Diseases* 2000; 182: 59–67. 951–956
4. The Genotype database is obtained from a large service testing laboratory from the US. It contains the resistance profiles of the Protease and Reverse Transcriptase genes of the HIV-1 virus obtained from plasma samples of HIV-1 infected patients. No clinical background information on medication or drug history is available. 957–961
5. Mathematical Methods for DNA Sequences. In: Waterman MS, eds. CRC Press Inc., Boca Raton, Florida, 1999. 963–964
6. Kiryukhin I, Saskov K, Boukhanovsky AV, Keulen W, Boucher, CA, Sloot PMA. Stochastic modeling of temporal variability of HIV-1 population. In: Sloot PMA, Abrahamson D, Bogdanov AV, Dongarra JJ, Zomaya AY, Gorbachev YE, eds. Computational Science – ICCS 2003, Melbourne, Australia and St. Petersburg, Russia, Proceedings Part I, in series Lecture Notes in Computer Science, vol. 2657, pp. 125–135. Springer Verlag, June 2003. ISBN 3-540-40194-6. 965–972
7. Zorzenon dos Santos RM, Coutinho S. Dynamics of HIV infection: A cellular automata approach. *Phys Rev Lett* 2001; 87(16): 168102–1–4. 973–975
8. Sloot PMA, Chen F, Boucher CA. Cellular automata model of drug therapy for HIV infection. In: Bandini S, Chopard B, Tomassini M, eds. 5th International Conference on Cellular Automata for Research and Industry, ACRI 2002, Geneva, Switzerland, October 9–11, 2002. Proceedings, in series Lecture Notes in Computer Science, vol. 2493, pp. 282–293. October 2002. 976–981
9. PHP: Hypertext Preprocessor: <http://www.php.net>. 982–983

- 984 10. The resource for PHP developers: <http://www.phpbuilder.com>
- 985 11. Zend Technologies – PHP tools for the development, pro- 996
- 986 tection and scalability of PHP applications – PHP for Linux, 997
- 987 Unix and Apache, Encoder, Accelerator Studio, Debugger: 998
- 988 <http://www.zend.com>. 999
- 989 12. The Apache Software Foundation: <http://www.apache.org>.
- 990 13. Alimadhi F. Mobile Internet: Wireless access to Web- 1000
- 991 based interfaces of legacy simulations, MSc thesis, Uni- 1001
- 992 versity of Amsterdam, The Netherlands, September 2002: 1002
- 993 <http://www.science.uva.nl/research/pacs/papers/master.html>. 1003
- 994 14. Cross-Grid: Grid technology of Interactive Distributed Com- 1004
- 995 putation: <http://www.eu-crossGrid.org/>. 1005
15. Little SJ, Holte S, Routy JP, Daar ES, Markowitz M, Collier AC, 1006
- Koup RA, Mellors JW, Connick E, Conway B, Kilby M, Wang 996
- L, Whitcomb JM, Hellmann NS, Richman DD. Antiretroviral- 997
- drug resistance among patients recently infected with HIV. *N* 998
- Engl J Med* 2002; 8;347(6): 385–394. 999
16. Karlin S. *A First Course in Stochastic Processes*. Academic Press. 1000
- NY-London, 1968. 1001
17. Sloot PMA, Boucher CA, Kiryukhin I, Saskov K, 1002
- Boukhanovsky AV. A grid-based problem-solving envi- 1003
- ronment for biomedicine. In: Nørager S, ed. *Proceedings of* 1004
- the First European HealthGrid Conference*, January, 16th–17th, 1005
- 2003, pp. 300–323. Commission of the European Commu- 1006
- nities, Information Society Directorate-General, Brussels, 1007
- Belgium, 2003. 1008

UNCORRECTED PROOF

1009 **Query**

1010 Q1. Au: Pls. provide dates.

UNCORRECTED PROOF