

achieve high SVR rates, while patients with a partial EVR (pEVR: 2 log drop in HCV RNA but still detectable at week 12) have lower rates of SVR.⁵ Since PEG-IFN RBV combination therapy is costly and accompanied by potential adverse effects, the ability to predict the possibility of RVR or cEVR before therapy and identifying curable patients may significantly influence the selection of patients for therapy. Moreover, identification of baseline predictors of poor response is particularly important to establish a rationale for identifying therapeutic targets to improve the efficacy of antiviral therapy.

Data mining is a method of predictive analysis which explores tremendous volumes of data to discover hidden patterns and relationships in highly complex datasets and enables the development of predictive models. The classification and regression tree (CART) analysis is a core component of the decision tree tool for data mining and predictive modeling,⁶ is deployed to decision makers in various fields of business, and currently is being used in the area of biomedicine.^{7–13} The results of CART analysis are presented as a decision tree, which is intuitive and facilitates the allocation of patients into subgroups by following the flow-chart form.¹⁴ CART has been shown to be competitive with other traditional statistical techniques such as logistic regression analysis.¹⁵

In the present study, we used the CART analysis to explore baseline predictors of response to PEG-IFN plus RBV therapy among clinical, biochemical, virological and histological pretreatment variables and to define a pre-treatment algorithm to discriminate chronic hepatitis C patients who are likely to respond to PEG-IFN plus RBV therapy.

MATERIALS AND METHODS

Patients

A TOTAL OF 419 chronic hepatitis C patients were treated with PEG-IFN alpha-2b and RBV at Musashino Red Cross Hospital between December 2001 and December 2007. Among them, 400 patients who fulfilled the following inclusion criteria were enrolled in the present study. (i) infection by genotype 1b (ii) HCV RNA higher than 100 KIU/mL by quantitative PCR (Cobas Amplicor HCV Monitor, Roche Diagnostic systems, CA) which is usually used for the definition of high viral load in Japan (iii) lack of co-infection with hepatitis B virus or human immunodeficiency virus (iv) lack of other causes of liver disease such as autoimmune hepatitis, primary biliary cirrhosis, or alcohol intake of more than 20 g per day, and (v) having completed at

least 12 weeks of therapy with an early virological response that could be evaluated. Patients received PEG-IFN alpha-2b (1.5 microgram/kg) subcutaneously every week and were administered a weight adjusted dose of RBV (600 mg for <60 kg, 800 mg for 60–80 kg, and 1000 mg for >80 kg) which is the recommended dosage in Japan. Data from two third of patients (269 patients) were used for the model building set and the remaining one third of patients (131 patients) were used as a validation set. Consent in writing was obtained from each patient and the study protocol conformed to the ethical guidelines of the 1975 Declaration of Helsinki and was approved by the institutional review committee.

Laboratory tests

Blood samples were obtained before therapy, and at least once every month during therapy and analyzed for hematologic tests, blood chemistries, and HCV RNA. In the present study, RVR and cEVR was defined as undetectable HCV RNA by qualitative PCR with a lower detection limit of 50 IU/mL (Amplicor, Roche Diagnostic systems, CA) at week 4 and 12, respectively. SVR was defined as undetectable HCV RNA at week 24 after the completion of therapy.

Histological examination

For all patients, liver biopsy specimens were obtained before therapy and were evaluated independently by three pathologists who were blinded to the clinical details. If there was a disagreement, the scores assigned by the majority of pathologists were used for the analysis. Fibrosis and activity were scored according to the METAVIR scoring system.¹⁶ Fibrosis was staged on a scale of 0–4: F0 (no fibrosis), F1 (mild fibrosis: portal fibrosis without septa), F2 (moderate fibrosis: few septa), F3 (severe fibrosis: numerous septa without cirrhosis) and F4 (cirrhosis). Activity of necroinflammation was graded on a scale of 0–3: A0 (no activity), A1 (mild activity), A2 (moderate activity) and A3 (severe activity). Percentage of steatosis was quantified by determining the average proportion of hepatocytes affected by steatosis and graded on a scale of 0–3: grade 0 (no steatosis), grade 1 (0–9%), grade 2 (10–29%), and grade 3 (over 30%) as we reported previously.¹⁷

Database for analysis

A pretreatment database of 72 variables was created containing histological findings (grade of fibrosis, activity, and steatosis), laboratory tests including the quantity of HCV RNA by Cobas Amplicor, and clinical information (age, gender, body weight, and body mass index).

The baseline characteristics and test results are listed in Table 1. The overall rate of RVR/cEVR was 43% in the model building set and 48% in the validation set. There were no significant differences in the clinical backgrounds between these two groups. Hepatitis C viral mutations, such as mutations in interferon-sensitivity determining region or core amino acid residues 70 and 91, were not included in the present analysis. The dataset of laboratory tests was based on the digitized records in this hospital. Continuous data was split into categorized data by increment of 10; For example, age was categorized into <30, 30–39, 40–49, 50–59, 60–69, and ≥70.

Statistical analysis

Based on this database, the recursive partitioning analysis algorithm referred to as CART was implemented to define meaningful subgroups of patients with respect to the possibility of achieving RVR/cEVR. The CART belongs to a family of nonparametric regression methods based on binary recursive partitioning of data. The software automatically explore the data to search for optimal split variables, builds a decision tree structure and finally classifies all subjects into particular subgroups that are homogeneous with respect to the outcome of interest.¹⁸ During the CART analysis, first, the entire study population, and thereafter, all newly defined subgroups, were investigated at every step of the analysis to determine which variable at what cut-off point yielded the most significant division into two prognostic subgroups that were as homogeneous as possible with respect to estimates of RVR/cEVR possibilities. This algorithm uses the impurity function (Gini criterion function) for splitting.¹⁹ A restriction was imposed on the tree construction such that terminal subgroups resulting from any given split must have at least 20 patients. The CART procedure stopped when either no additional significant variable was detected or when the sample size was below 20. The resulting final subgroups were most homogeneous with respect to the probability of achieving RVR/cEVR. For this analysis, data mining software Clementine version 12.0 (SPSS Inc, Chicago, IL) was utilized. SPSS 15.0 (SPSS Inc, Chicago, IL) was used for logistic regression analysis.

RESULTS

Factors associated with RVR/cEVR by standard statistical analysis

WE FIRST ANALYZED 72 variables by univariate and multivariate logistic regression analysis to find factors associated with RVR/cEVR (Table 2).

Patients with RVR/cEVR were significantly younger than those without. Among histological findings, grade of steatosis and stage of fibrosis was significantly lower in RVR/cEVR. Among hematologic tests, hemoglobin and hematocrit was significantly higher in RVR/cEVR. Among blood chemistry tests, creatinine and low-density lipoprotein cholesterol (LDL-C) was significantly higher and gamma-glutamyltransferase (GGT), low-density-lipoprotein cholesterol (LDL-C), and blood sugar were significantly lower in RVR/cEVR. The level of HCV RNA was significantly lower in RVR/cEVR. There were no significant differences in other tests.

Multivariate logistic regression analysis was performed on age, fibrosis stage, steatosis, HCV RNA, creatinine, hemoglobin, GGT, LDL-C, and blood sugar: hematocrit was not included since it is closely associated with hemoglobin. On multivariate analysis, age, grade of steatosis, level of HCV RNA, creatinine, hemoglobin, GGT, and LDL-cholesterol remained significant whereas stage of fibrosis, hemoglobin and blood sugar were not.

The CART analysis

The CART analysis was carried out on the model building set of 269 patients using the same variables as logistic regression analysis. Figure 1 shows the resulting decision tree. The CART analysis automatically selected five predictive variables to produce a total of seven subgroups of patients. The grade of steatosis was selected as the variable of initial split with an optimal cut-off of 30%. The possibility of achieving RVR/cEVR was only 18% for patients with hepatic steatosis of 30% or more compared to 47% for patients with hepatic steatosis of less than 30%. Among patients with hepatic steatosis of less than 30%, the level of serum LDL-C, with an optimal cut-off of 100 mg/dL, was selected as the variable of second split. Patients with higher LDL-C level had the higher probability of RVR/cEVR (57% vs. 32%). Among patients with LDL-C of less than 100 mg/dL, age, with an optimal cut-off of 60, was selected as the third variable of split. Younger patients had the higher probability of RVR/cEVR (49% vs. 15%). Among patients younger than 60, the blood sugar, with an optimal cut-off of 120 mg/dL, was selected as the fourth variable of split. Patients with lower blood sugar level had the higher probability of RVR/cEVR (71% vs. 31%). Among patients with hepatic steatosis of less than 30% and LDL-C of 100 mg/dL or more, age, with an optimal cut-off of 50, was selected as the third variable of split, younger being the predictor of higher RVR/cEVR probability (77% vs. 50%). Among patients older than 50,

Table 1 Clinical characteristics of patients

	Model set n = 269	Validation set n = 131	P-value
Sex (M/F)	127/142	55/76	0.325
Age (years)	57.7 ± 10.1	57.6 ± 10.0	0.932
Body weight (kg)	59.6 ± 11.0	57.5 ± 9.5	0.094
Body mass index (kg/m ²)	23.2 ± 3.1	23.3 ± 3.8	0.934
Total protein (g/dL)	7.6 ± 0.5	7.7 ± 0.6	0.558
Albumin (g/dL)	4.2 ± 0.3	4.2 ± 0.3	0.349
Globulin (g/dL)	3.4 ± 0.5	3.4 ± 0.6	0.989
Aspartate aminotransferase (IU/L)	58.1 ± 43.1	55.8 ± 37.5	0.601
Alanine aminotransferase (IU/L)	70.9 ± 49.2	66.4 ± 52.6	0.462
Gamma-glutamyltransferase (IU/L)	49.6 ± 44.0	45.2 ± 34.4	0.33
Lactate dehydrogenase (IU/L)	289.3 ± 112.3	301.5 ± 109.3	0.417
Total bilirubin (mg/dL)	0.71 ± 0.28	0.69 ± 0.23	0.317
Direct bilirubin (mg/dL)	0.23 ± 0.12	0.25 ± 0.10	0.147
Indirect bilirubin (mg/dL)	0.48 ± 0.21	0.44 ± 0.16	0.064
Alkaline phosphatase (IU/L)	290.9 ± 107.6	292.5 ± 107.6	0.917
Leucine aminopeptidase (IU/L)	64.3 ± 14.3	65.5 ± 12.3	0.543
Thymol turbidity test (KU)	7.1 ± 3.4	8.0 ± 3.7	0.062
Zinc sulfate turbidity test (KU)	15.4 ± 4.9	16.3 ± 5.4	0.188
Choline esterase (IU/L)	318.1 ± 81.7	321.1 ± 78.1	0.798
Ammonia (microg/dL)	39.7 ± 20.2	45.0 ± 15.6	0.668
Blood sugar (mg/dL)	125.9 ± 41.1	117.4 ± 47.9	0.081
Glycohemoglobin (%)	5.6 ± 1.6	5.4 ± 1.2	0.797
Total cholesterol (mg/dL)	170.8 ± 33.9	175.6 ± 36.8	0.170
Low-density-lipoprotein-cholesterol (mg/dL)	96.5 ± 25.2	100.9 ± 28.5	0.153
High-density-lipoprotein-cholesterol (mg/dL)	54.2 ± 15.9	55.2 ± 17.4	0.612
Triglyceride (mg/dL)	108.5 ± 47.8	102.8 ± 46.4	0.306
Creatinine (mg/dL)	0.72 ± 0.15	0.74 ± 0.17	0.236
Urea nitrogen (mg/dL)	14.1 ± 3.4	14.9 ± 3.9	0.123
Uric acid (mg/dL)	5.3 ± 1.2	5.2 ± 1.2	0.715
Sodium (mEq/L)	142.2 ± 2.0	142.4 ± 2.0	0.471
Potassium (mEq/L)	4.3 ± 0.3	4.3 ± 0.4	0.578
Chloride (mEq/L)	104.0 ± 2.2	104.0 ± 2.6	0.905
Calcium (mg/dL)	9.1 ± 0.4	9.2 ± 0.4	0.479
Phosphorus (mg/dL)	3.5 ± 0.5	3.5 ± 0.6	0.814
Magnesium (mg/dL)	2.2 ± 0.2	2.3 ± 0.3	0.390
Amylase (IU/L)	178.7 ± 125.8	175.1 ± 133.1	0.118
Creatine kinase (IU/L)	114.9 ± 147.6	119.3 ± 73.7	0.849
Iron (microg/dL)	104.7 ± 53.2	109 ± 37	0.726
Ferritin (ng/mL)	111.3 ± 103.3	59.7 ± 118.5	0.405
C-reactive peptide (mg/dL)	0.2 ± 1.1	0.1 ± 0.1	0.586
Immunoglobulin G (mg/dL)	1849 ± 426	1988 ± 525	0.129
Immunoglobulin M (mg/dL)	141 ± 69	205 ± 106	0.200
Immunoglobulin A (mg/dL)	323 ± 675	291 ± 81	0.784
Triiodothyronine (pg/mL)	2.3 ± 0.3	2.2 ± 0.3	0.358
Thyroxin (ng/dL)	0.9 ± 0.1	0.9 ± 0.1	0.872
Thyroid stimulating hormone (micro IU/mL)	1.8 ± 1.4	1.7 ± 0.7	0.939
White blood cell count (/microl)	5243 ± 1591	5286 ± 1101	0.843
Segmented neutrophils (%)	55.4 ± 10.8	57.0 ± 10.0	0.297
Band neutrophils (%)	1.5 ± 1.6	0.5 ± 0.6	0.250
Eosinophils (%)	2.9 ± 2.3	2.4 ± 1.4	0.127

Table 1 Continued

	Model set n = 269	Validation set n = 131	P-value
Basophiles (%)	0.6 ± 0.4	0.6 ± 0.3	0.727
Lymphocytes (%)	34.6 ± 9.6	34.0 ± 9.3	0.682
Monocytes (%)	6.6 ± 2.2	6.2 ± 2.6	0.149
Red blood cell count (10 ⁴ /microl)	458 ± 43	455 ± 47	0.643
Hemoglobin (g/dL)	14.4 ± 1.5	14.5 ± 1.5	0.618
Hematcrit (%)	42.7 ± 4.0	42.9 ± 4.4	0.717
Reticulocytes (%)	1.4 ± 0.4	1.4 ± 0.4	0.762
Mean corpuscular volume (fL)	93.3 ± 4.5	93.8 ± 5.41	0.466
Mean corpuscular hemoglobin concentration (pg)	31.5 ± 1.9	31.7 ± 2.3	0.583
Mean corpuscular hemoglobin concentration (g/dL)	33.8 ± 0.9	33.7 ± 1.3	0.910
Platelets (10 ⁴ /microl)	16.8 ± 5.4	16.3 ± 4.5	0.480
Prothrombin time (s)	11.7 ± 1.2	11.7 ± 0.9	0.762
Prothrombin time (activity %)	104.6 ± 14.4	102.6 ± 14.8	0.363
Prothrombin time (international normalized ratio)	1.0 ± 0.1	1.0 ± 0.1	0.387
Thrombin time (%)	97.2 ± 31.3	109 ± 31.5	0.231
Activated partial thromboplastin time (s)	29.7 ± 4.4	29.1 ± 2.7	0.260
Hepaplastin test (%)	97.8 ± 20.3	95.4 ± 19.4	0.523
Fibrinogen (%)	237 ± 44	225 ± 45	0.069
Hepatitis C virus RNA (<850/≥850 KIU/mL)	130/139	70/61	0.394
Histological grade of			
Activity (A1/A2/A3)	138/107/24	62/55/14	0.714
Fibrosis (F1/F2/F3/F4)	135/74/57/3	58/40/27/6	0.131
Steatosis (0%/1–9%/10–29%/30%≧)	89/109/37/34	49/45/21/16	0.643
Hepatitis C virus RNA negative at week 12 (yes/no)	116/153	63/68	0.349

the level of GGT, with an optimal cutoff of 40 U/L, were then selected as the fourth level of split, low levels being the predictor of higher RVR/cEVR probability (60% vs. 35%).

All five factors selected as significant variables in the CART analysis were also significantly associated with RVR/cEVR by univariate analysis (Table 2). In addition, steatosis, LDL-C, age and GGT were also independently

Table 2 Factors associated with rapid or complete early virological response by univariate and multivariate logistic regression analysis

Parameter	Category	Univariate			Multivariate		
		Odds	95% CI	P-value	Odds	95% CI	P-value
Age (years)	<50 vs. ≥50	2.65	1.51–4.65	<0.001	2.03	1.04–3.97	0.039
Fibrosis stage	F1-2 vs. F3-4	2.47	1.31–4.66	0.005	1.77	0.85–3.68	0.120
Steatosis (%)	<30 vs. ≥30	4.11	1.64–10.29	0.003	2.88	1.07–7.79	0.037
Hepatitis C virus RNA (KIU/mL)	<850 vs. ≥850	1.97	1.21–3.22	0.007	1.93	1.09–3.43	0.025
Creatinine (mg/dL)	≥0.8 vs. <0.8	3.30	1.96–5.56	<0.001	3.54	1.88–6.67	<0.001
Hemoglobin (g/dL)	≥14.5 vs. <14.5	1.76	1.08–2.87	0.023	1.38	0.74–2.57	0.320
Hematcrit (%)	≥43 vs. <43	1.75	1.07–2.84	0.003			
Gamma-glutamyltransferase (IU/L)	<40 vs. ≥40	2.06	1.26–3.37	0.004	2.45	1.32–4.56	0.005
Low-density-lipid cholesterol (mg/dL)	≥100 vs. <100	2.71	1.61–4.55	<0.001	2.21	1.21–4.06	0.010
Blood sugar (mg/dL)	<120 vs. ≥120	2.00	1.02–3.95	0.045	1.42	0.64–3.13	0.390

CI, confidence interval.

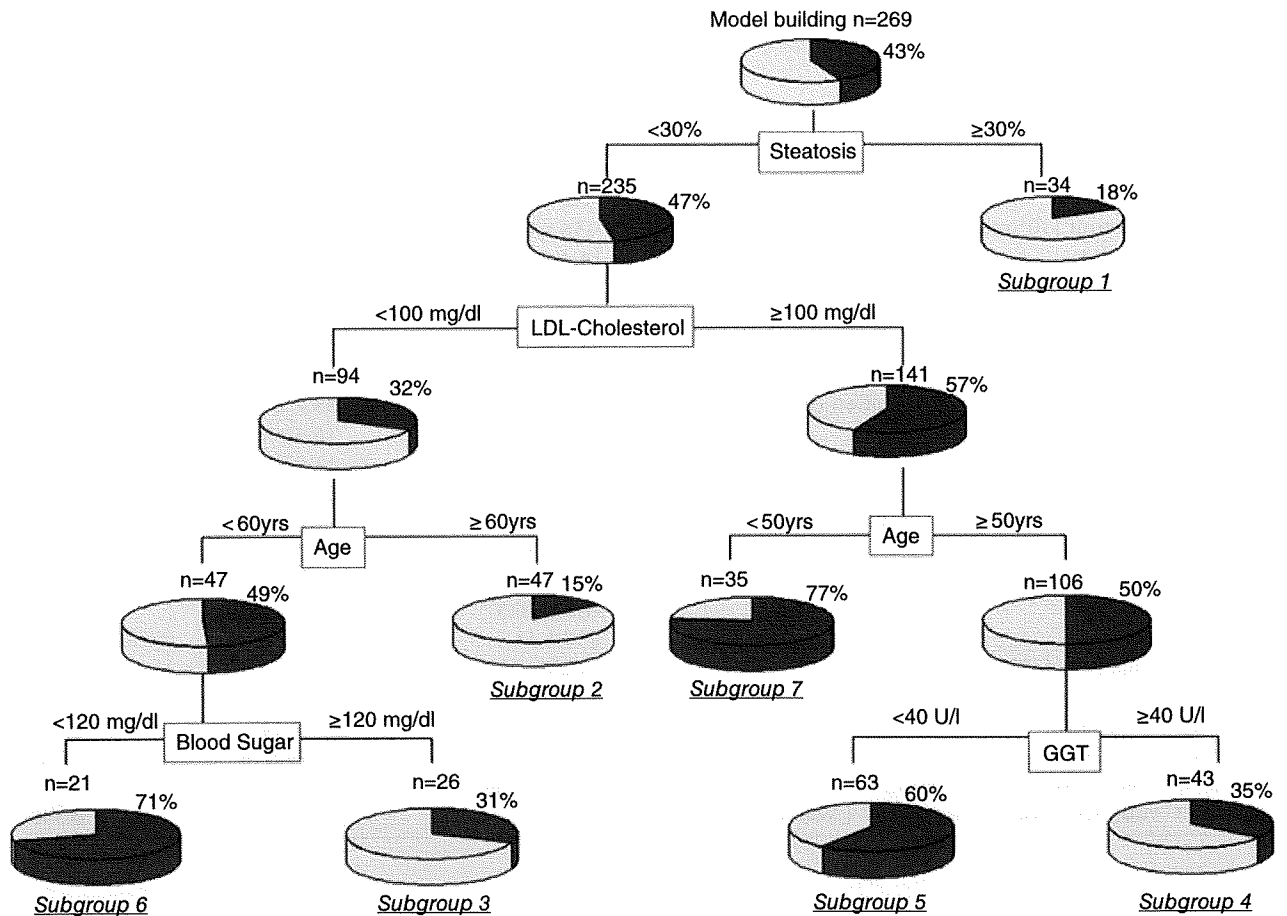


Figure 1 Classification and regression tree analysis. Boxes indicate the factors used for splitting and the cut-off value for the split. Pie charts indicate the rate of RVR/cEVR for each group of patients after splitting. Terminal subgroups of patients discriminated by the analysis are numbered from one to seven. GGT, gamma-glutamyltransferase; LDL, low-density-lipoprotein.

associated with RVR/cEVR by multivariate logistic regression analysis while blood sugar was not (Table 2). On the other hand, HCVRNA and creatinine which were significantly associated with RVR/cEVR by multivariate analysis were not selected as significant variables in CART analysis.

The probabilities of RVR/cEVR for the seven subgroups derived by this process were highly variable. The subgroup whose hepatic steatosis was less than 30%, serum LDL-C was 100 mg/dL or more and of an age less than 50 years (subgroup 7) showed the highest probability of RVR/cEVR (77%), while the subgroup whose hepatic steatosis more than 30% (subgroup 1) and the subgroup whose hepatic steatosis was less than 30% but serum LDL-C was less than 100 mg/dL and of an age

greater than 60 years (subgroup 2) showed the lowest probability of RVR/cEVR (18% and 15%, respectively).

Validation of the CART analysis

The results of the CART analysis were validated with a validation dataset of 131 cases which is independent of the model building dataset. Each patient in the validation set was allocated to subgroups 1-7 using the flow-chart form of the CART tree. The rates of RVR/cEVR were 20% for subgroups 1 and 2, 29% for subgroups 3, 38% for subgroup 4, 59% for subgroup 5, 71% for subgroup 6, and 85% for subgroups 7. The rates of RVR/cEVR for each subgroup of patients were closely correlated between the model building dataset and the validation dataset (Fig. 2).

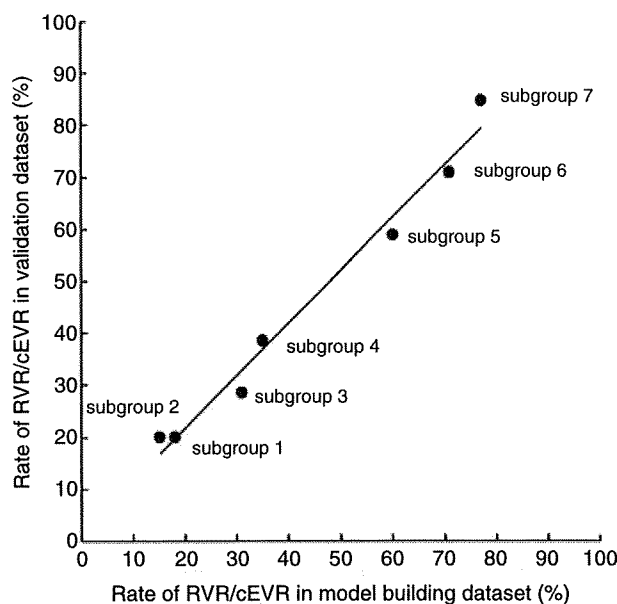


Figure 2 Validation of the classification and regression tree (CART) analysis: Subgroup stratified comparison of the rate of rapid or complete early virological response (RVR/cEVR) between the model building and validation datasets. Each patient in the validation set was allocated to subgroups 1–7 by following the flow-chart form of the CART tree and the rates of RVR/cEVR were calculated. The rate of RVR/cEVR in each subgroup was plotted. The x-axis represents the rate of RVR/cEVR in the model building datasets and the y-axis represents the rate of RVR/cEVR in the validation datasets. The rates of achieving RVR/cEVR in each subgroup of patients closely correlated between the model building dataset and the validation dataset ($r^2 = 0.987$).

Construction of 3 groups according to the probability of RVR/cEVR

If the seven subgroups were reconstructed into three groups according to their rate of RVR/cEVR, the rate of RVR/cEVR was 16% for low probability group (subgroup 1 and 2), 46% for intermediate probability group (subgroup 3, 4, and 5) and 75% for high probability group (subgroup 6 and 7; $P < 0.0001$).

Effect of adherence

Adherence of PEG-IFN and RBV was not included as a variable of analysis since the present study aimed to develop a pre-treatment model for the prediction of response. To analyze the possible effect of adherence on the result of CART analysis, three groups of patients divided by CART (low, intermediate and high probability group) were further stratified according to adherence

of PEG-IFN and RBV. Poor adherence was defined as taking less than 80% planned dose of PEG-IFN or RBV at 12 weeks, and good adherence was defined as taking more than 80% planned dose of both PEG-IFN and RBV at 12 weeks. The result is shown in Figure 3. Among patients with good adherence, the rate of RVR/cEVR was 19% for low probability group, 52% for intermediate probability group and 77% for high probability group. Among poor adherence group, the rate of RVR/cEVR was 13% for low probability group, 41% for intermediate probability group and 73% for high probability group. Collectively, even after adjustment for adherence, 3 groups of patients divided by CART analysis still had low, intermediate and high probability of achieving RVR/cEVR, respectively.

DISCUSSION

IN THE PRESENT study, we performed the CART analysis and built a simple decision tree model for the pre-treatment prediction of response to PEG-IFN plus

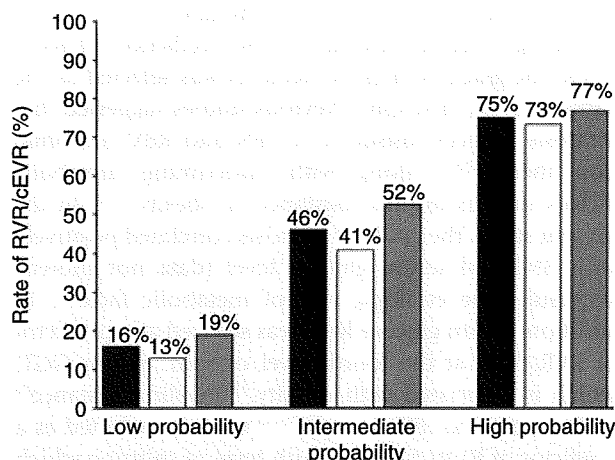


Figure 3 The rate of rapid or complete early virological response (RVR/cEVR) between the classification and regression tree (CART) groups stratified by adherence. The three groups of patients divided by CART (low, intermediate and high probability group) were further stratified according to adherence of peg-interferon (PEG-IFN) plus ribavirin (RBV). Black, white and gray boxes in the bar chart indicate total patients, patients with poor adherence (taking less than 80% planned dose of PEG-IFN or RBV at 12 weeks), and good adherence (taking more than 80% planned dose of both PEG-IFN and RBV at 12 weeks), respectively. Even after adjustment for adherence, 3 groups of patients divided by CART analysis still had low, intermediate and high probability of achieving RVR/cEVR, respectively.

RBV therapy. The analysis highlighted five host variables relevant to response: steatosis, LDL-C, age, blood sugar and GGT. Classification of patients based on these variables identified subgroups of patients with high probabilities of achieving RVR/cEVR among difficult to treat chronic hepatitis C patients. The reproducibility of the model was confirmed by the independent validation datasets. According to the result of the CART, patients were categorized into 3 groups: the rate of RVR/cEVR was 16% for low probability group, 46% for intermediate probability group and 75% for high probability group. The result of the CART analysis could be readily applicable to clinical practice because patients could be allocated to specific subgroups with a defined rate of response simply by following the flow-chart form. Although an early disappearance of serum HCV RNA is the prerequisite for achieving SVR, no reliable baseline predictors of response to PEG-IFN plus RBV therapy are established to date. Thus, this model may have the potential to support decisions in patient selection for PEG-IFN plus RBV therapy or to tailor treatment strategies for individual patients. Moreover, our result may provide a rationale for treating metabolic factors to improve the efficacy of antiviral therapy.

Among variables relevant to the prediction of RVR/cEVR, the grade of hepatic steatosis was selected as the variable of the first split. Previous studies suggested that steatosis induces resistance to IFN and RBV combination therapy^{20,21} along with underlining metabolic factors such as insulin resistance or obesity.^{21–24} In the present study, the grade of steatosis correlated positively with BMI and serum glucose level (data not shown) suggesting the etiologic role of metabolic factors. In addition, serum glucose level was selected as a predictor of RVR/cEVR at the fourth level of split. Serum GGT, which is associated with obesity,²⁵ insulin resistance²⁶ and response to IFN therapy,^{27–30} was also selected as a predictor of RVR/cEVR at fourth level of splitting which may emphasize the importance of metabolic factors in therapeutic resistance. These findings raise the possibility that treatment of these metabolic factors may improve the virological response to the PEG-IFN plus RBV therapy. This hypothesis should be examined by a prospective study.

We and others have reported that steatosis, obesity and insulin resistance are associated with the progression of fibrosis,^{17,31–33} which can interfere indirectly with the effect of IFN on hepatocytes. Other possible mechanisms of resistance by steatosis or metabolic factors include dysregulation of adipocytokines³⁴ or oxidative stress which may inhibit intracellular IFN signaling

pathway.³⁵ Despite these findings, the precise mechanism of resistance is not established and further investigation is needed.

Another factor relevant in the prediction of RVR/cEVR was LDL-C. LDL-C was selected as the second factor for splitting by CART, and was an independent predictor of RVR/cEVR by logistic regression analysis. LDL-C recently has attracted attention as a novel predictor of response to IFN or PEG-IFN plus RBV.^{30,36,37} Since *in vitro* study showed that LDL-C receptor acts as a receptor for HCV and LDL-C competitively inhibit the binding of HCV,³⁸ high level of serum LDL-C may inhibit HCV entry to hepatocytes and attenuate replication. LDL-C and its receptor may be a future therapeutic target.

Not all factors selected as significant variables in the CART analysis were also significantly associated with response by standard statistical analysis: blood sugar was associated with response by univariate analysis but not by multivariate logistic regression analysis. On the other hand, HCV RNA and creatinine which were significantly associated with RVR/cEVR by multivariate analysis were not selected as significant variables in CART analysis. These differences may indicate both the unique feature and the limitations of the CART analysis. To note, blood sugar was significantly associated with RVR/cEVR within specialized subgroups of patients defined by the CART analysis: in subgroup of patients with steatosis <30%, LDL-C <100 mg/dL and younger than 60, which indicate the unique feature of the CART analysis that it could visualize significant predictors that specifically apply to selected patients. The limitation is that not all significant factors may be adopted in the decision tree since we applied the rule to stop CART procedure when the sample size was below 20. This rule was applied to avoid the generation of over-fit model which may lack universality. Therefore, it is possible that HCV RNA or creatinine may become a significant variable in the CART analysis if larger number of patients were included in the analysis. Stage of fibrosis was significantly associated with response to therapy by univariate analysis but not by multivariate analysis and not selected as a significant variable in the CART analysis. The possible reason is that advanced fibrosis is associated with older age as a confounding factor.

CART analyses are gaining acceptance in medical research in addition to biomedical field. Recent publications include the prediction of aggressive prostate cancer,⁸ diabetic vascular complications,¹⁹ prognosis of melanoma,^{7,39} response to preoperative radiotherapy for rectal tumor,⁹ prognostic groups in colorectal carcinoma,¹² and outcome after liver failure.¹¹ An advantage

of CART over traditional regression models is that it can identify prognostic subgroups that are useful in clinical practice. Because the results of CART analysis are presented as a decision tree, which is intuitive, they can be readily interpreted by medical professionals without any specific knowledge of statistics. The most important consideration is that five variables used in the decision tree were clinical parameters that are readily available by the usual work-up of patients before therapy. Especially, glucose, GGT and LDL-C are simple biochemical markers that are easily measured at a low cost. Using this model, we can rapidly develop an estimate of the response before treatment, which may facilitate clinical decision making.

In conclusion, we built a pre-treatment model for the prediction of virological response in PEG-IFN plus RBV therapy. Because this decision tree model was made up of simple host factors such as steatosis, LDL-C, age, blood sugar and GGT, it can be easily applied to clinical practice. This model may have the potential to support decisions in patient selection for PEG-IFN plus RBV therapy based on the possibility of response against a potential risk of adverse events or costs, and may provide a rationale for treating metabolic factors to improve the efficacy of antiviral therapy.

ACKNOWLEDGEMENTS

THIS STUDY WAS supported by a grant-in-aid from Ministry of Health, Labor and Welfare, Japan. There exist no conflicts of interest.

REFERENCES

- 1 Strader DB, Wright T, Thomas DL, Seeff LB. Diagnosis, management, and treatment of hepatitis C. *Hepatology* 2004; 39: 1147–71.
- 2 Fried MW, Shiffman ML, Reddy KR *et al*. Peginterferon alfa-2a plus ribavirin for chronic hepatitis C virus infection. *N Engl J Med* 2002; 347: 975–82.
- 3 Manns MP, McHutchison JG, Gordon SC *et al*. Peginterferon alfa-2b plus ribavirin compared with interferon alfa-2b plus ribavirin for initial treatment of chronic hepatitis C: a randomised trial. *Lancet* 2001; 358: 958–65.
- 4 Davis GL, Wong JB, McHutchison JG, Manns MP, Harvey J, Albrecht J. Early virologic response to treatment with peginterferon alfa-2b plus ribavirin in patients with chronic hepatitis C. *Hepatology* 2003; 38: 645–52.
- 5 Lee SS, Ferenci P. Optimizing outcomes in patients with hepatitis C virus genotype 1 or 4. *Antivir Ther* 2008; 13 (Suppl 1): 9–16.
- 6 Breiman L, Friedman RA, Olshen CJ, Stone CM. *Classification and Regression Trees*. Calif: Wadsworth, 1980.
- 7 Averbook BJ, Fu P, Rao JS, Mansour EG. A long-term analysis of 1018 patients with melanoma by classic Cox regression and tree-structured survival analysis at a major referral center: Implications on the future of cancer staging. *Surg* 2002; 132: 589–602.
- 8 Garzotto M, Beer TM, Hudson RG *et al*. Improved detection of prostate cancer using classification and regression tree analysis. *J Clin Oncol* 2005; 23: 4322–9.
- 9 Zlobec I, Steele R, Nigam N, Compton CC. A predictive model of rectal tumor response to preoperative radiotherapy using classification and regression tree methods. *Clin Cancer Res* 2005; 11: 5440–3.
- 10 Jin H, Lu Y, Harris ST *et al*. Classification algorithms for hip fracture prediction based on recursive partitioning methods. *Med Decis Making* 2004; 24: 386–98.
- 11 Baquerizo A, Anselmo D, Shackleton C *et al*. Phosphorus an early predictive factor in patients with acute liver failure. *Transplantation* 2003; 75: 2007–14.
- 12 Valera VA, Walter BA, Yokoyama N *et al*. Prognostic groups in colorectal carcinoma patients based on tumor cell proliferation and classification and regression tree (CART) survival analysis. *Ann Surg Oncol* 2007; 14: 34–40.
- 13 Martin MA, Meyricke R, O'Neill T, Roberts S. Mastectomy or breast conserving surgery? Factors affecting type of surgical treatment for breast cancer – a classification tree approach. *BMC Cancer* 2006; 6: 98.
- 14 LeBlanc M, Crowley J. A review of tree-based prognostic models. *Cancer Treat Res* 1995; 75: 113–24.
- 15 Costanza MC, Paccaud F. Binary classification of dyslipidemia from the waist-to-hip ratio and body mass index: a comparison of linear, logistic, and CART models. *BMC Med Res Methodol* 2004; 4: 7.
- 16 Bedossa P, Poynard T. An algorithm for the grading of activity in chronic hepatitis C. The METAVIR Cooperative Study Group. *Hepatology* 1996; 24: 289–93.
- 17 Kurosaki M, Matsunaga K, Hirayama I *et al*. The presence of steatosis and elevation of alanine aminotransferase levels are associated with fibrosis progression in chronic hepatitis C with non-response to interferon therapy. *J Hepatol* 2008; 48: 736–42.
- 18 Segal MR, Bloch DA. A comparison of estimated proportional hazards models and regression trees. *Stat Med* 1989; 8: 539–50.
- 19 Miyaki K, Takei I, Watanabe K, Nakashima H, Omae K. Novel statistical classification model of type 2 diabetes mellitus patients for tailor-made prevention using data mining algorithm. *J Epidemiol* 2002; 12: 243–8.
- 20 Akuta N, Suzuki F, Tsubota A *et al*. Efficacy of interferon monotherapy to 394 consecutive naive cases infected with hepatitis C virus genotype 2a in Japan: therapy efficacy as consequence of tripartite interaction of viral, host and interferon treatment-related factors. *J Hepatol* 2002; 37: 831–6.

- 21 Poynard T, Ratziu V, McHutchison J *et al.* Effect of treatment with peginterferon or interferon alfa-2b and ribavirin on steatosis in patients infected with hepatitis C. *Hepatology* 2003; 38: 75-85.
- 22 Bressler BL, Guindi M, Tomlinson G, Heathcote J. High body mass index is an independent risk factor for non-response to antiviral treatment in chronic hepatitis C. *Hepatology* 2003; 38: 639-44.
- 23 Romero-Gomez M, Del Mar Vilorio M, Andrade RJ *et al.* Insulin resistance impairs sustained response rate to peginterferon plus ribavirin in chronic hepatitis C patients. *Gastroenterology* 2005; 128: 636-41.
- 24 Konishi I, Horiike N, Hiasa Y *et al.* Diabetes mellitus reduces the therapeutic effectiveness of interferon-alpha2b plus ribavirin therapy in patients with chronic hepatitis C. *Hepatol Res* 2007; 37: 331-6.
- 25 Marchesini G, Avagnina S, Barantani EG *et al.* Aminotransferase and gamma-glutamyltranspeptidase levels in obesity are associated with insulin resistance and the metabolic syndrome. *J Endocrinol Invest* 2005; 28: 333-9.
- 26 Fraser A, Ebrahim S, Smith GD, Lawlor DA. A comparison of associations of alanine aminotransferase and gamma-glutamyltransferase with fasting glucose, fasting insulin, and glycated hemoglobin in women with and without diabetes. *Hepatology* 2007; 46: 158-65.
- 27 Mazzella G, Salzetta A, Casanova S *et al.* Treatment of chronic sporadic-type non-A, non-B hepatitis with lymphoblastoid interferon: gamma GT levels predictive for response. *Dig Dis Sci* 1994; 39: 866-70.
- 28 Villela-Nogueira CA, Perez RM, de Segadas Soares JA, Coelho HS. Gamma-glutamyl transferase (GGT) as an independent predictive factor of sustained virologic response in patients with hepatitis C treated with interferon-alpha and ribavirin. *J Clin Gastroenterol* 2005; 39: 728-30.
- 29 Berg T, Sarrazin C, Herrmann E *et al.* Prediction of treatment outcome in patients with chronic hepatitis C: significance of baseline parameters and viral dynamics during therapy. *Hepatology* 2003; 37: 600-9.
- 30 Akuta N, Suzuki F, Kawamura Y *et al.* Predictive factors of early and sustained responses to peginterferon plus ribavirin combination therapy in Japanese patients infected with hepatitis C virus genotype 1b: amino acid substitutions in the core region and low-density lipoprotein cholesterol levels. *J Hepatol* 2007; 46: 403-10.
- 31 Adinolfi LE, Gambardella M, Andreana A, Tripodi MF, Utili R, Ruggiero G. Steatosis accelerates the progression of liver damage of chronic hepatitis C patients and correlates with specific HCV genotype and visceral obesity. *Hepatology* 2001; 33: 1358-64.
- 32 Ortiz V, Berenguer M, Rayon JM, Carrasco D, Berenguer J. Contribution of obesity to hepatitis C-related fibrosis progression. *Am J Gastroenterol* 2002; 97: 2408-14.
- 33 Muzzi A, Leandro G, Rubbia-Brandt L *et al.* Insulin resistance is associated with liver fibrosis in non-diabetic chronic hepatitis C patients. *J Hepatol* 2005; 42: 41-6.
- 34 Charlton MR, Pockros PJ, Harrison SA. Impact of obesity on treatment of chronic hepatitis C. *Hepatology* 2006; 43: 1177-86.
- 35 Di Bona D, Cippitelli M, Fionda C *et al.* Oxidative stress inhibits IFN-alpha-induced antiviral gene expression by blocking the JAK-STAT pathway. *J Hepatol* 2006; 45: 271-9.
- 36 Minuk GY, Weinstein S, Kaita KD. Serum cholesterol and low-density lipoprotein cholesterol levels as predictors of response to interferon therapy for chronic hepatitis C. *Ann Intern Med* 2000; 132: 761-2.
- 37 Gopal K, Johnson TC, Gopal S *et al.* Correlation between beta-lipoprotein levels and outcome of hepatitis C treatment. *Hepatology* 2006; 44: 335-40.
- 38 Agnello V, Abel G, Elfahal M, Knight GB, Zhang QX. Hepatitis C virus and other flaviviridae viruses enter cells via low density lipoprotein receptor. *Proc Natl Acad Sci USA* 1999; 96: 12766-71.
- 39 Leiter U, Buettner PG, Eigentler TK, Garbe C. Prognostic factors of thin cutaneous melanoma: an analysis of the central malignant melanoma registry of the german dermatological society. *J Clin Oncol* 2004; 22: 3660-7.



Reproducibility and usability of chronic virus infection model using agent-based simulation; comparing with a mathematical model

Jun Itakura^{a,*}, Masayuki Kurosaki^a, Yoshie Itakura^a, Sinya Maekawa^b, Yasuhiro Asahina^a, Namiki Izumi^a, Nobuyuki Enomoto^b

^a Division of Gastroenterology and Hepatology, Musashino Red Cross Hospital, 1-26-1 Kyonan-cho, Musashino-shi, Tokyo 180-8610, Japan

^b First Department of Internal Medicine, Faculty of Medicine, University of Yamanashi, 1110, Shimogatou, Chuou-shi, Yamanashi 409-3898, Japan

ARTICLE INFO

Article history:

Received 30 June 2009

Received in revised form 27 August 2009

Accepted 6 September 2009

Keywords:

Agent-based model

Virus infectious disease

ABSTRACT

We created agent-based models that visually simulate conditions of chronic viral infections using two software. The results from two models were consistent, when they have same parameters during the actual simulation. The simulation results comprise a transient phase and an equilibrium phase, and unlike the mathematical model, virus count transit smoothly to the equilibrium phase without overshooting which correlates with actual biology in vivo of certain viruses. We investigated the effects caused by varying all the parameters included in concept; increasing virus lifespan, uninfected cell lifespan, uninfected cell regeneration rate, virus production count from infected cells, and infection rate had positive effects to the virus count during the equilibrium period, whereas increasing the latent period, the lifespan-shortening ratio for infected cells, and the cell cycle speed had negative effects. Virus count at the start did not influence the equilibrium conditions, but it influenced the infection development rate. The space size had no intrinsic effect on the equilibrium period, but virus count maximized when the virus moving speed was twice the space size. These agent-based simulation models reproducibly provide a visual representation of the disease, and enable a simulation that encompasses parameters those are difficult to account for in a mathematical model.

© 2009 Elsevier Ireland Ltd. All rights reserved.

1. Introduction

All viruses need hosts as a basis for their life. When a virus enters the host body, it invades cells and uses both its own enzymes and those of the host cells to replicate. Host cells infected by viruses launch a self-defense system known as the innate immune system (See and Wark, 2008; Nanche, 2009), which inhibits viral replication and uses the human leukocyte antigen system and cytokines to elicit an immune response. Immune cells that have received signals from host cells activate other immune cells, neutralize viruses in the serum by means of antibodies, and prevent the virus from replicating and proliferating by destroying or curing host cells. Viral infection is a disorder based on the interactions between viruses and cells.

The power relationship between these agents changes along with the progression of the disease. In the very early stages of infection, as the host defense mechanisms are immature, the virus has the ability to overwhelm the host cells, actively replicate, and proliferate. Subsequently, as the capacity of the immune system improves, the speed of viral proliferation drops and the virus count reaches a peak. Infected host cells begin to be disrupted by the immune system or virus particles, and symptoms appear as a result. If the immune system is stronger than the virus, then the viral counts decline, and, in transient viral disorders, the virus is finally eliminated and the host recovers. In chronic viral disorders, however, the power relationship between the virus and host cells reaches equilibrium, and a long-term power balance is maintained with the virus count reaching a plateau.

Mathematical models have been proposed to study the dynamics of such viral disorders, and are regarded as being of value in understanding this phenomenon (Ho et al., 1995; Nowak et al., 1996; Neumann et al., 1998). However, these models are difficult to understand for clinicians, and their applicability is somewhat limited in everyday practice. In clinical research, measurements of viral dynamics in patients for short duration have been made for human

Abbreviations: HIV, human immunodeficiency virus; HBV, hepatitis B virus; HCV, hepatitis C virus.

* Corresponding author. Tel.: +81 422 32 3111; fax: +81 422 32 9551.

E-mail address: jitakura@musashino.jrc.or.jp (J. Itakura).

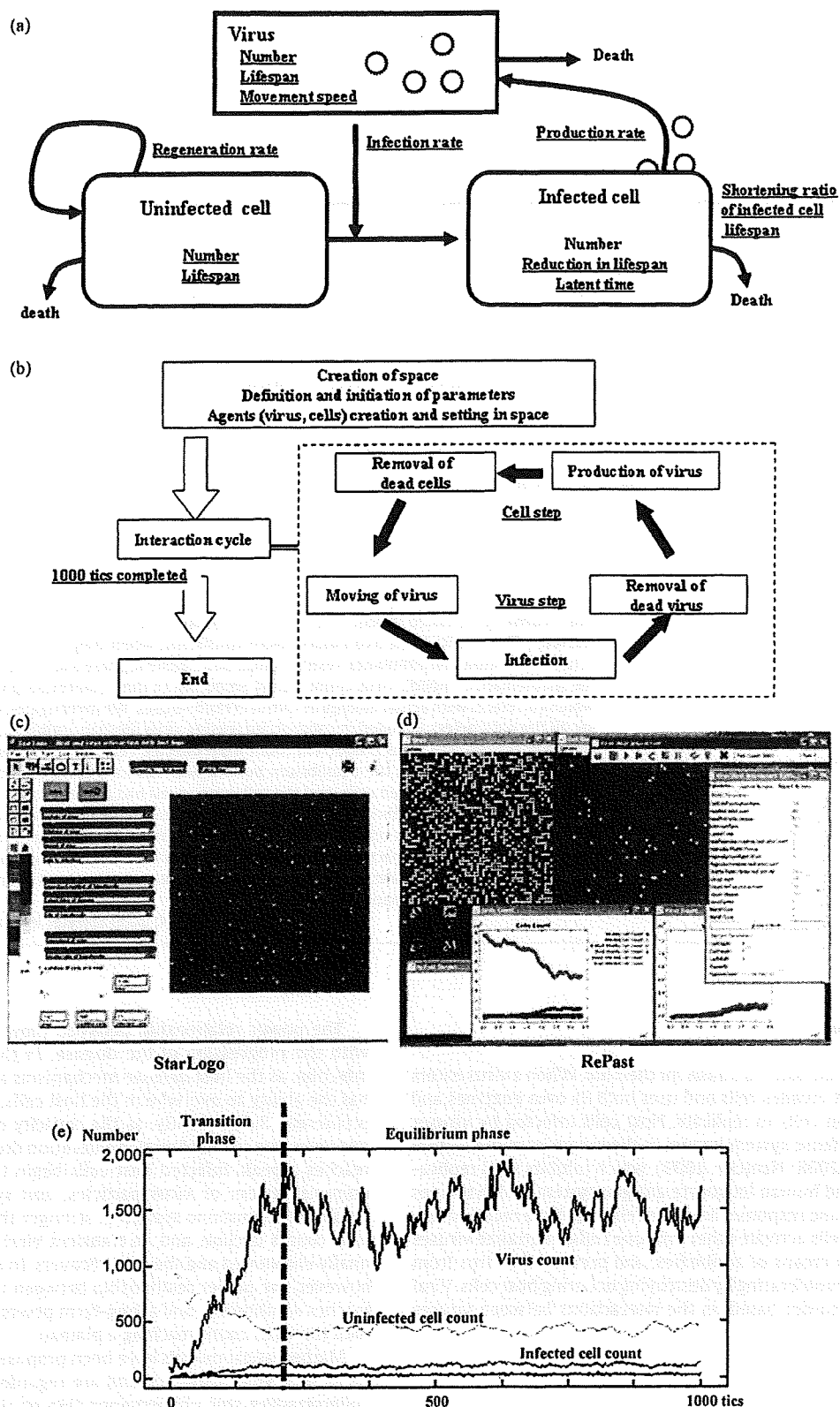


Fig. 1. Simulation design and an example of simulation results. (a) Model concept. Viruses, uninfected cells, and infected cells were treated as agents, and parameters were set for each of these and for interactions between agents (underlined). (b) Flowchart of the program. After preparing the simulation, we entered the interaction cycle, in which virus steps (such as movement) and cell steps were repeated. One cycle was counted as 1 tic, and the simulation concluded after 1000 tics. (c and d) Simulation screen using (c) StarLogo and (d) RePast. Yellow circles are viruses, green squares are uninfected cells, and orange and red indicate infected cells, with orange indicating the latent period. In StarLogo, all the agents are shown on the same screen, but in RePast, viruses and cells are shown in separate windows. (e) Example of a simulation chart in StarLogo. After the start of simulation the virus count and infected cell count increase while the uninfected cell count decreases, with equilibrium state reached after a certain number of tics.

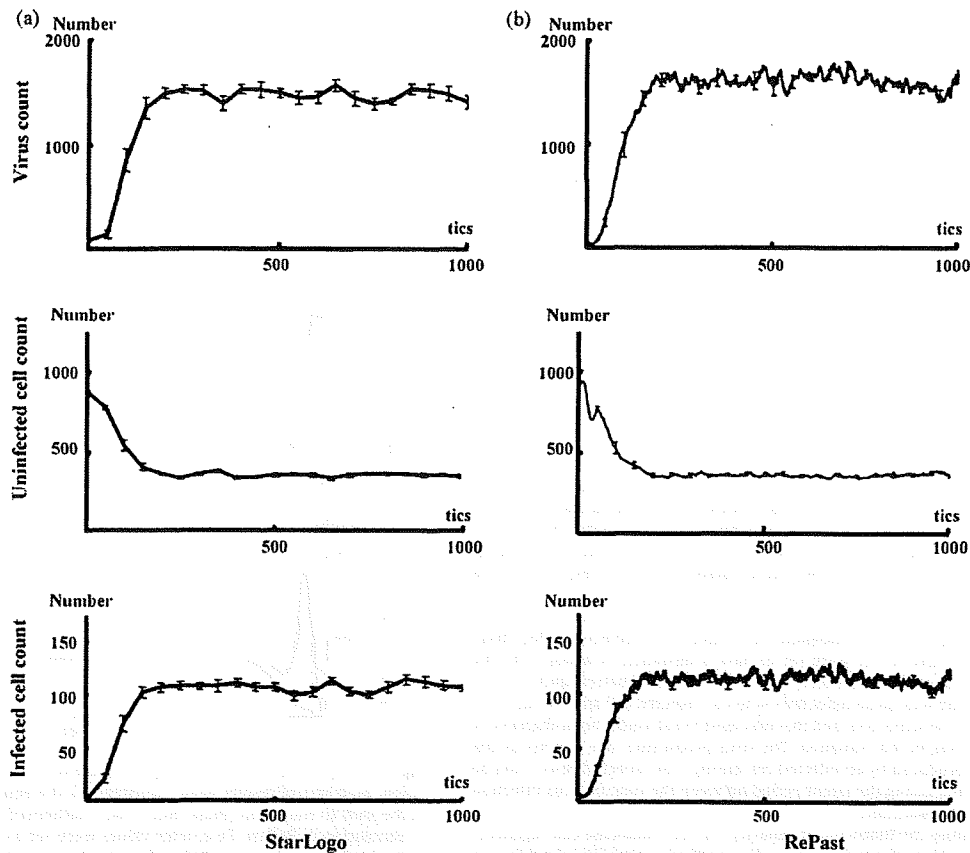


Fig. 2. Comparison of simulation results in (a) StarLogo and (b) RePast. The results were consistent when the parameters were made consistent. (Virus count [average \pm SD]: StarLogo 1458.03 ± 173.1 , RePast 1462.71 ± 178.8 , $p=0.94$. Uninfected cell count: 364.24 ± 30.4 , 368.11 ± 33.4 , $p=0.83$. Infected cell count: 105.73 ± 13.0 , 107.74 ± 13.0 , $p=0.24$. Unpaired Student's *t*-test.) Parameter values were set as follows: initial virus count, 100; uninfected cell count, 880; infected cell count, 0; virus speed of movement, 5 grids/tic; infection rate, 10%; uninfected cell regeneration rate, 1%; latent period, 3 tics; and virus reproduction rate, 5/cells/tic. The following parameter settings were taken from actual measurements: virus lifespan, 4.5 tics; uninfected cell lifespan, 49.8 tics; and infected cell lifespan, 6.7 tics.

immunodeficiency virus (HIV) (Ho et al., 1995), hepatitis B virus (HBV) (Nowak et al., 1996) and hepatitis C virus (HCV) (Neumann et al., 1998), and research is also underway on a range of models based on animal experiments and cell culture systems. As chronic viral disorders persist over long periods of time complete follow-up of viral dynamics is difficult. Furthermore, limitations of items that can be measured, such as the difficulty of measuring whole numbers of host cells, make it extremely difficult to investigate their consistency in mathematical models.

The recent ascend of dynamic-models owes much to advances in computers. Computer performance has improved markedly in recent years, not only in terms of their calculating capacity but also with regard to image displays, and models that offer a visual representation of viral disorders are now being reported (Gilbert and Bankes, 2002; Duca et al., 2007; Shapiro et al., 2008; Castiglione et al., 2007). One advantage of such visual models is that by providing a visual representation, they make understanding the disease status easy. Another benefit is that they enable parameters to be identified that are hidden as background noise in mathematical models. However, these models have some problems; it is difficult to prove the reproducibility of the simulation results derived from different languages or libraries, difficult to prove the validity of the model's concepts, and difficult to prove that the simulation results accurately reflect the reality. In this study, we created agent-based computer models that visually simulate the conditions of chronic viral infections using two software. The reproducibility of two agent-based computer models and the differences between agent-based models and the mathematical model were analyzed.

This agent-based model enabled us to investigate how each parameter included in the concept affects the conditions of chronic viral infections.

2. Methods

2.1. Selection of Software

In this study, we used two different types of softwares: StarLogo version 2.0 (<http://education.mit.edu/starlogo/>) supplied by MIT Media Laboratory and Recursive Porous Agent Simulation Toolkit (RePast-3.0, <http://repast.sourceforge.net/>) supplied by the Argonne National Laboratory. StarLogo uses Logo, one of the simplest programming languages, and has a fixed graphical user interface. RePast is a library that uses Java, another programming language, which also has a fixed graphical user interface.

Logo is an assembly language, and StarLogo carries out processing sequentially. Java is an object-oriented language, and RePast has a faster processing speed than StarLogo. In addition, StarLogo has a number of stipulations to simplify simulations, such as parameters can only be set up to five decimal places and the simulation space is also fixed as 51×51 square grids. RePast, on the other hand, has fewer such restrictions. Thus, it offers a higher degree of freedom in program settings than StarLogo. Taking simulation space as an example, in spite of the restrictions imposed by the underlying operating system's image display system, any number of grids can be set and a hexagonal grid could also be chosen rather than a square one. However, users must stipulate and set all parameters themselves. This means that they must first declare the shape of the grid and the number of grids they will use to fill the simulation space. Java is also more difficult to learn than Logo, and debugging and correcting the program is also more difficult. Thus, it is difficult to judge whether or not the results agree with the planned simulation.

In effect, these two different types of softwares are polar opposites. It is simple to start a simulation in StarLogo, but producing results takes time and it is difficult to carry out more complex simulations. In RePast it is difficult to compose the program and judge whether or not the planned study has actually been achieved, but the

simulation itself takes only a short time to complete and there are lesser restrictions in the construction of a simulation model.

2.2. Concept for Modeling

We applied the basic virus–host interaction mathematical model to the agent-based simulation system with slight modifications. The mathematical model was used to describe the dynamics of HIV (Ho et al., 1995), HBV (Nowak et al., 1996), and HCV (Neumann et al., 1998) and the only agents involved were host cells and viruses, without the inclusion of immune cells. The effects of the immune system are expressed by varying parameters such as lifespan of host cells and viruses.

Fig. 1a illustrates the study concept. Viruses have the ability to infect healthy host cells (uninfected cells) and the infected cells produce new viruses. Both cells and viruses have definite lifespans, and the lifespan of infected cells is usually shorter than that of uninfected cells. Uninfected cells automatically regenerate within the space, whereas infected cells only arise due to infection of uninfected cells. Viruses also lack the ability to regenerate themselves and are only produced from infected cells.

2.3. Parameter Settings

In the present study, as the StarLogo settings are circumscribed, we limited the simulation space to 51 × 51 square grids. However, we made an exception here while investigating the effects of size of space on the simulation results. The numbers of viruses, uninfected cells, and infected cells could only be set before the start of the simulation. As described in the later, our simulation ran in cycles, with 1 cycle defined as 1 tic.

In mathematical simulation models, the death rate is required as a parameter. However, in our program we set lifespans for viruses and uninfected cells. These lifespans were not uniform, but were set to have a deviation of about 10%. The lifespan of cells was shortened by infection with ratio decided beforehand.

The infection ratio was meaningful only when an infected cell and a virus coincidentally occupied the same grid, and this was used to calculate the probability of the infection occurring in that situation. The virus production rate was set as the number of viruses produced by an infected cell during 1 tic. Infected cells could be set as a parameter indicating the latent period between the time of virus infection and the time of virus replication.

In order to emulate the tissue repair capacity, we set uninfected cell regeneration rate such that grids without any cells had a specified probability of producing uninfected cells on top of themselves. As a result, the more the cell count declined within a space the more regenerated uninfected cells were produced, whereas the number of regenerated cells declined as cell count increased.

The number of grids through which a virus could move in 1 tic was set as the speed of movement, and the direction of movement was set within a range of 90° toward the top of the simulation space. The program used a circulatory method of movement; when a virus arrived at the top of the space, it was translocated to the bottom, and moved again toward the top. Cells were fixed on the grid.

2.4. Simulation Flowchart

Fig. 1b shows a flowchart of the program. First, the simulation space was produced, after which each parameter was defined and the initial settings were made. Next the agents – viruses and uninfected and infected cells – were produced. The simulation cycle was as follows. Viruses moved to a new grid, and if an uninfected cell was present, this was infected with a probability based on the infection rate. The lifespan of the virus decreased, and viruses that had completed their lifespan and those that had caused an infection were removed from the space. Infected cells produced new viruses, the lifespans of both uninfected and infected cells decreased. Then, cells that had completed their lifespan were eliminated and a new cycle began. The program was set such that the simulation ended after this cycle had repeated 1000 times. This meant that one simulation was complete after 1000 tics.

2.5. Data Collection

The RePast model was programmed such that data for each tic was saved automatically as a text file at the end of the simulation. This text file could be opened by a database software. The StarLogo model was programmed to stop the simulation and collect data after every 50 tics.

2.6. Mathematical Model

In order to compare the results of this agent-based simulation, we used a viral infection mathematical model, which we improved as follows.

$$\frac{dT}{dt} = s[2601 - (T + I)] - dT - bVT \tag{1}$$

$$\frac{dI}{dt} = bVT - dI \tag{2}$$

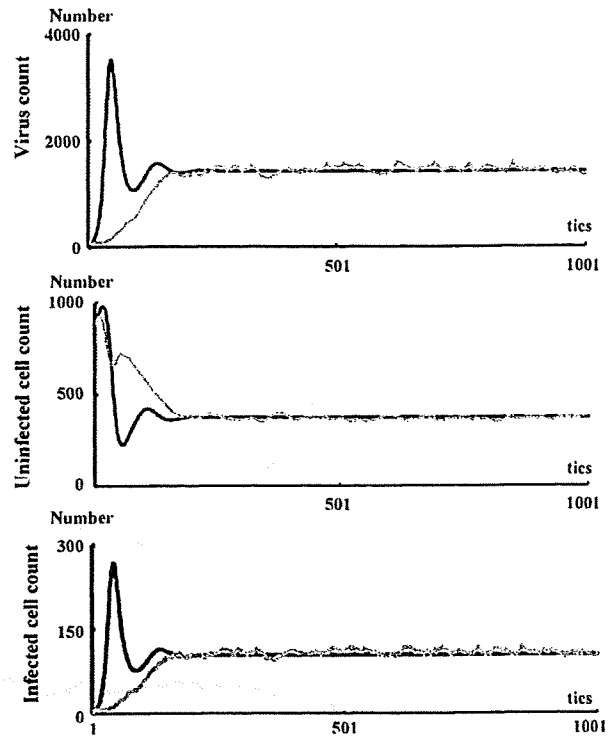


Fig. 3. Comparison of results of agent-based simulation and mathematical simulation. Both sets of results were consistent for the equilibrium phase, but differed in the shift in transition phase. Black line: mathematical model; grey line: results of simulation in RePast. Parameter values were set as follows: initial virus count, 100; uninfected cell count, 880; infected cell count, 0; virus speed of movement, 5 grids/tic; infection rate, 10%; uninfected cell regeneration rate, 1%; latent period, 3 tics; virus reproduction rate, 5/cells/tic; virus lifespan, 10 tics; uninfected cell lifespan, 50 tics; and cell lifespan-shortening ratio as a result of infection, 69%.

$$\frac{dV}{dt} = pI - cV \tag{3}$$

where, T is the uninfected cell count, I is the infected cell count, and V is the virus count. Uninfected cells are supplied to the space with a probability $s[2601 - (T + I)]$, as the number of grids in this agent-based simulation model was 2601 (51 × 51). The death rate of uninfected cells is d , the death rate of infected cells is δ , and the death rate of viruses is c . The infection rate is indicated by β . Viruses are released from infected cells at a probability p .

2.7. Statistical Analysis

Statistical analyses were performed by statistical tests using the program StatView 5.0 (SAS Institute Inc.). All tests of significance were two-tailed, with p values of <0.05 considered to be significant.

3. Results

3.1. Reproducibility of Chronic Viral Infection Disease Models Using Agent-based Simulation Methods

We constructed the chronic viral infection model with StarLogo library. Fig. 1c shows the simulation screen, and Fig. 1e shows one sample result. Immediately after the start of the simulation, the virus count temporarily dropped in accordance with the onset of an infection. Subsequently, the virus count started to increase with an increase in the infected cells and a decrease in the uninfected cells. After a certain number of tics (around 300 in this example), although the virus count, infected cell count, and uninfected cell count had some fluctuation, an equilibrium state was reached. We use the following descriptive terms in this paper: the transient phase is the period during which virus growth peaks, and the equilibrium phase is the period during which an equilibrium state is

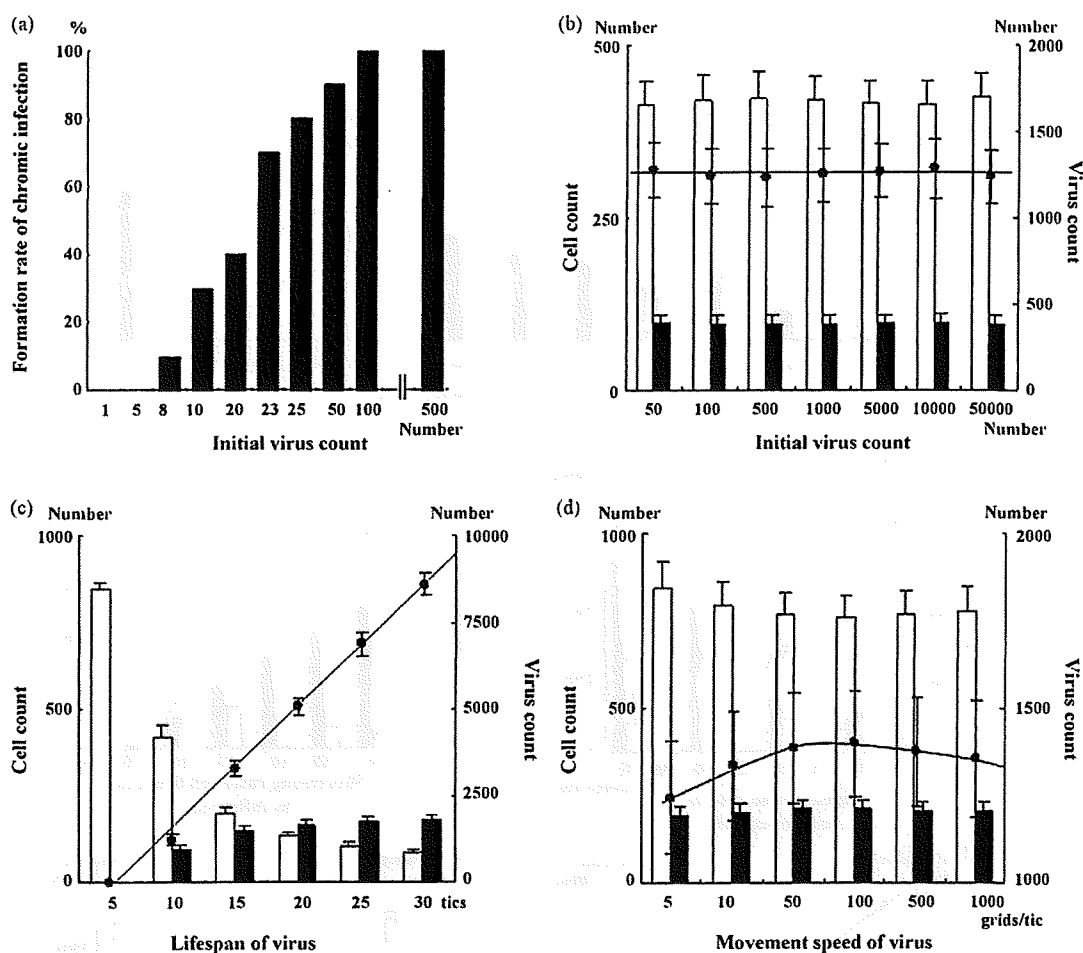


Fig. 4. Effects of changes in viral parameters. (a) The higher the initial virus count, the greater is the increase in the rate of formation of chronic infection, but (b) there was no effect on the conditions in the equilibrium phase. (c) Extending the virus lifespan increased the virus count. (d) Increasing the speed of virus movement to 100 grids/tic increased the virus count, but increasing it to 500 grids/tic had the opposite effect, with a slight declining trend. (a) Black bars: number of infections produced; (b–d) black circles: virus count; line: virus count approximation curve; white bars: uninfected cell count; black bars: infected cell count.

established. When the simulation was performed multiple times, the features described above were maintained, and the average values for virus, infected cell, and uninfected cell counts during the equilibrium state were all consistent.

Fig. 1d shows the simulation screen of the RePast. When we attempted setting all the initial parameters to the same values as those in the StarLogo, the results were not consistent. When we recalculated the parameters from the simulation results, in RePast, the parameters were largely maintained at the levels of the settings, but in StarLogo, the lifespans of both cell types became shorter than the settings while the simulation was in progress. We made the results of both simulations consistent by using the same parameters during the actual simulation (Fig. 2a and b).

3.2. Comparison Between Agent-based Simulation Models and Mathematical Simulation Model

We investigated whether the results of a chronic viral infection disease model produced by RePast would be consistent with the results of a mathematical model. For the mathematical model, we carried out an approximate integration using a four-dimensional Runge–Kutta method to ensure that the uninfected cell count and infected cell count would be in the same class. Parameters were always fixed as constant between simulations. The simulation results were consistent for the equilibrium

phase, but transitions in virus count during the transient phase varied, with a shift to equilibrium state following two overshoots in the mathematical model, but a monotonic increase following a logistic curve in the agent-based model (Fig. 3). In the mathematical model, when the equilibrium condition was calculated with $dT/dt = dI/dt = dV/dt = 0$, the equilibrium-phase virus count, uninfected cell count, and infected cell count were very similar to those of the agent-based model (virus count: mathematical model 371.8/space, agent-based model 371.1 ± 32.4 /space [average \pm SD]; uninfected cell count: mathematical model 1605/space, agent-based model 1454 ± 194 /space; infected cell count: mathematical model 115.9/space, agent-based model 108.3 ± 14.2 /space).

3.3. Usability of the Models; Effect of Changing Parameters

We investigated the changes in the equilibrium phase brought about by changing each parameter. All the investigations below were carried out by using RePast, and we used the average values from ten simulations.

3.4. Viral Parameters

The lower the virus counts at the beginning of the simulation, the lower the probability of a chronic infection (Fig. 4a). However, the initial virus count did not have any effect on the equilibrium

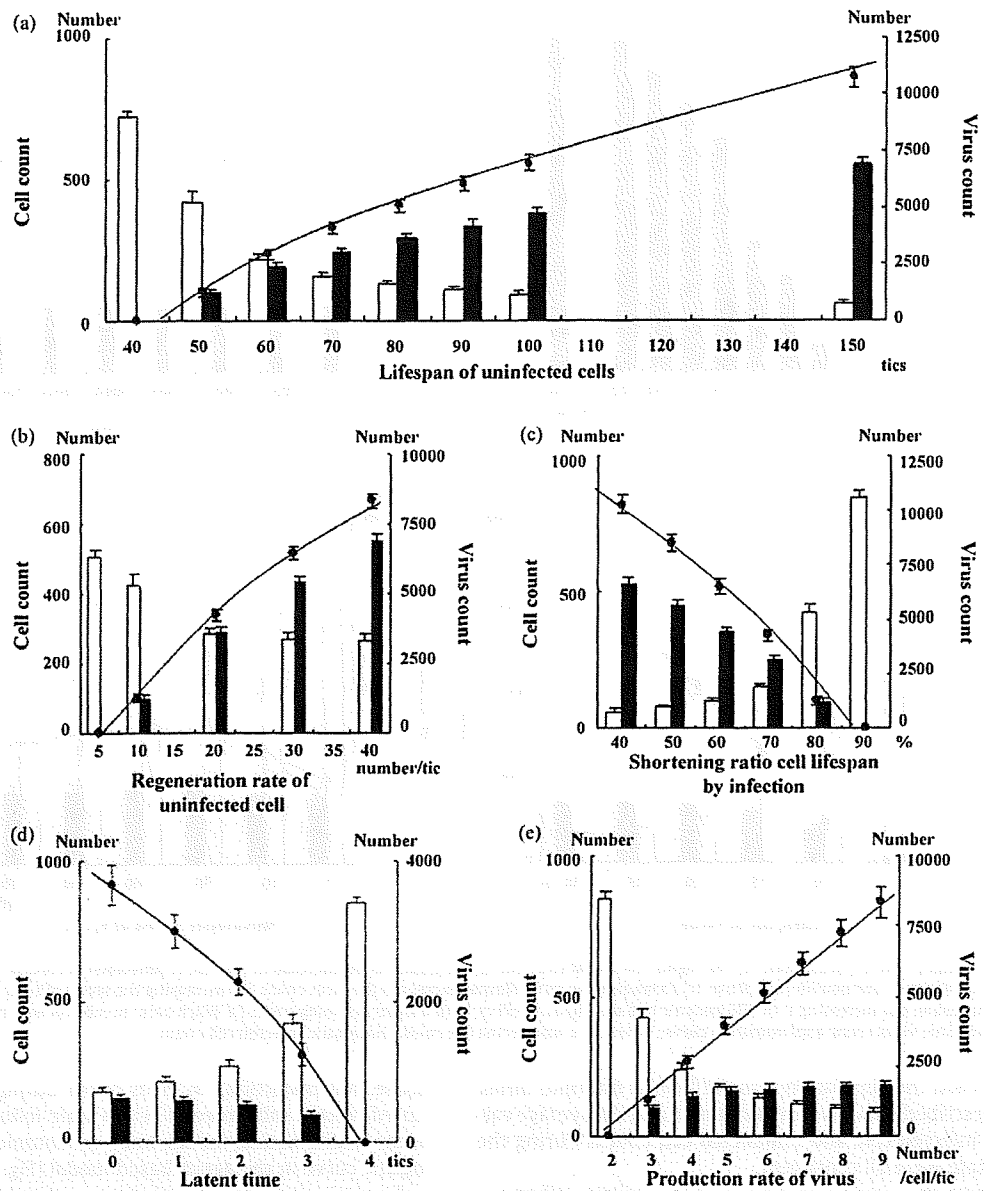


Fig. 5. Effects of changes in cell parameters. (a) Extending the uninfected cell lifespan and (b) increasing the uninfected cell regeneration rate increased the virus count. (c) Raising the lifespan-shortening ratio as a result of infection shortened the lifespan of infected cells, thereby decreasing the virus count. (d) Extending the latent period shortened the period of virus production from infected cells, thereby decreasing the virus count. (e) Increasing the virus production count resulted in a linear increase in equilibrium-phase virus count. Black circles: virus count; line: virus count approximation curve; white bars: uninfected cell count; black bars: infected cell count.

phase itself (Fig. 4b). Extending the lifespan of viruses resulted in a linear increase in equilibrium-phase virus count (Fig. 4c). Although the infected cell count increased, the rate of increase gradually declined. Changing the speed of viral movement resulted in the equilibrium-phase virus count to eventually decline after 100 grids/tic was reached, allowing movement over an area twice the size of the simulation space (Fig. 4d).

3.5. Uninfected Cell Parameters

Extending the lifespan of uninfected cells led to an increased virus count during the equilibrium phase (Fig. 5a). Increasing the uninfected cell regeneration rate also contributed to increased equilibrium-phase virus count (Fig. 5b). In both the cases, the

increases in virus count and infected cell count were not linear, but showed a tendency for the rate of increase to decline gradually.

3.6. Infected Cell Parameters

We carried out an investigation of the effects of variation in the lifespan-shortening ratio on the virus count on the assumption that cell lifespan is shortened by infection. When this ratio was increased, the virus count decreased (Fig. 5c). An extended latent period was also related to a decreased virus count (Fig. 5d). However, the virus production from infected cells led to a linear increase in the virus count (Fig. 5e).

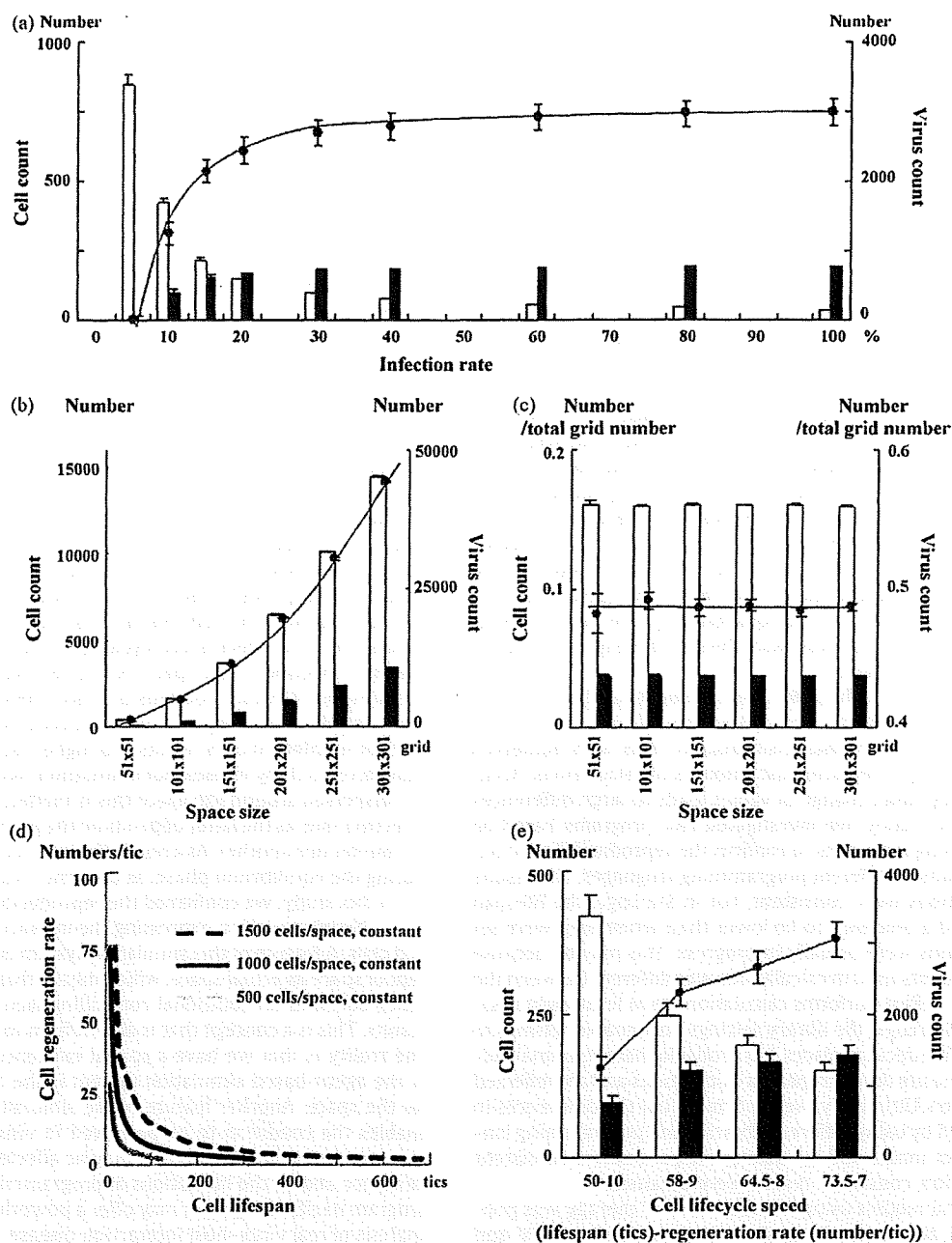


Fig. 6. (a) Increasing the infection rate increased the virus count in equilibrium periods, but the virus count did not change at infection rates of 30% or more. (b) The size of the simulation space increased not only virus count but also the cell count; however, (c) when virus and cell counts were divided by the total number of grids in the space, they were constant for all space sizes. (d) Changing the lifespan and regeneration rate of uninfected cells in opposite directions at the same time makes it possible to change only the cell cycle speed without altering the uninfected cell count. (e) When the cell cycle speed was reduced, the virus count increased toward the right of the graph. This may be because the effect of extending the lifespan of cells exceeds that of reducing their regeneration rate. (a–c and e) Black circles: virus count; line: virus count approximation curve; white bars: uninfected cell count; black bars: infected cell count.

3.7. Infection Rate and Space Size

Increasing the infection rate caused an increase in the virus count, but the change was minimal at an infection rate of 30% or more. The same results were seen for infected cell count, but a decrease in uninfected cell count resulted in a tendency for the infection rate to decrease by up to 60% (Fig. 6a).

The larger the space, higher the increase in both virus and cell counts (Fig. 6b). This increase was proportional to space size, how-

ever, when virus and cell counts were divided by the total number of grids in the space they were all constant (Fig. 6c).

3.8. Cell Cycle Speeds

Running a simulation with the initial virus count set to zero enables only the equilibrium condition for uninfected cells to be simulated. Changing the lifespan and regeneration rate of uninfected cells in opposite directions at the same time makes it possible

to change the cell cycle speed without altering the uninfected cell count (Fig. 6d). We used this technique to investigate how changing the cell cycle speed affected the equilibrium phase. Fig. 6e shows the results. Cell lifespan increases while the cell cycle speed declines. The equilibrium virus count increased in accordance with slower cell cycle speeds.

4. Discussion

In this study, we investigated the models using two agent-based simulation methods to program a simple virus–host chronic infection model. The same model written in two different programming language systems displayed the same results. The transient phase was unlike that seen in a mathematical simulation with no overshoot in virus count, but rather a smooth transition to the equilibrium phase. The virus count at the start of the simulation only had effect on the rate of infection development. Increases in virus lifespan, uninfected cell lifespan, uninfected cell regeneration rate, virus production count from infected cells, and infection rate all led to increased equilibrium-phase virus count. Rises in the infected cell lifespan-shortening ratio, latent period, and cell cycle speed decreased the equilibrium-phase virus count. The size of the space itself had no innate effect on the equilibrium phase, but a speed of movement of the virus that was twice the size of the space produced the maximum virus count.

Reproducibility is the basis for all scientific study, but there are many problems to prove it in computer simulations, such as programming bugs. As agent-based simulation deals with numerous agents individually, it requires vast amounts of calculations. Accumulation of very small change of values leads to large differences of results. In this study, we investigated two programs based on two programming languages to confirm the reproducibility of our simulation results in different programming languages. The results of two simulations were consistent, but in StarLogo, the lifespan parameters had a tendency to be lower than when they were set while simulations were actually in progress. This may be because the number of digits used in calculations was different between the two programs. RePast performs calculations to at least eight decimal places. In StarLogo, the library settings only enable settings to be made up to five decimal places. It is probable that these small differences accumulate during repeated calculations and are reflected in the simulation. Ultimately, we confirmed that the differences in results obtained by using different libraries and programming languages were not innate and by making the parameters consistent during simulation, consistent results were obtained.

Mathematical models using formulae for HIV therapy was published in 1994, the method has since been applied to HBV and HCV (Ho et al., 1995; Nowak et al., 1996; Neumann et al., 1998), and they were thought to be good reflections of the reality. In the mathematical model, viruses and cells are conceived as individuals in the concept itself, but both of them are perceived *en masse* when calculations are performed. However a feature of the agent-based simulation is that it deals with individual viruses and cells as separate agents. By moving each agent individually, it probes the factors influencing overall shifts from the micro viewpoint. When the space is viewed as a whole, it is possible to watch on the screen the collective movement of groups of agents. Recently, models that provide a visual representation of Epstein-Barr virus and HIV infection have been reported, both of which are useful for an instinctive and intuitive understanding (Duca et al., 2007; Shapiro et al., 2008; Castiglione et al., 2007).

In agent-based simulation model, virus count transit smoothly to the equilibrium phase. On the other hand, virus counts overshoot during transient phase in mathematical model. We think this difference is derived from technicality of different model-

ing. The difference in concepts between mathematical models and agent-based models is the space. The mathematical model has no space in concept, but agents move across the space in the agent-based model. In agent-based models, the densities of virus and cells change overtime especially in the transition phase because of the limited space. These changes of the densities of virus and cells lead to the dynamic change of the encounter rate of viruses and cells. The mathematical model does not make such concept of the density; the encounter rate is constant. This may be the reason for the difference between two models in the transition phase. Since no overshoot of virus counts in transient phase had been reported from *in vivo* studies of hepatitis C virus and simian immunodeficiency virus (Dahari et al., 2005; Nowak et al., 1997), agent-based model correlates with actual biology *in vivo* at least for these viruses. The increase of initial virus count at the start of simulation correlates with higher encounter rate of viruses and cells which make the linear increasing of infection forming rate. Mathematical model can only express the infection formation rate as “infected or not”.

The importance of viral passing speed in the agent-based model is also explained by the “space”. Although the virus actually moves through the blood stream in our body and virus could not decide their moving speeds by themselves, there is most appropriate speed for virus to meet the cells on the simulation space by the highest probability. The effect of cell cycle speed should be mentioned by another affection of the space. A fast cell cycle speed means that the lifespan of uninfected cells is short. Then fast cell cycle speed leads to the short lifespan of infected cells. A higher regeneration rate for uninfected cells results in a higher rate of infection among uninfected cells by viruses, but in situations where viruses and cells are dispersed around the space this is ineffective in increasing the infection rate, as the latter depends on the probability that they will encounter one another. As a result, the infected cell count decreases during the equilibrium phase, as does the virus count.

In this study, we confirmed the reproducibility and usability of agent-based models in expressing the interaction between viruses and cells. A feature of this simulation system is that it uses the concept of space as actual space, which means that the existence of the space becomes an additional controlling factor on the simulation results. This is a concept that is absent from mathematical models. The reality is that we have a spatial existence, and an advantage of the agent-based simulation system is the fact that it accounts for the space. Another feature of the simulation system is that it enables the condition to be perceived in visual terms, making it easy to understand. However it may be affected by computer performance and by the limitations of programming languages or the program itself, this system may offer a powerful tool for the future analysis of real virus–host interaction disease.

Conflict of interest

No conflicts of interest exist for all authors.

References

- Castiglione, F., Pappalardo, F., Bernaschi, M., Motta, S., 2007. Optimization of HAART with genetic algorithms and agent-based models of HIV infection. *Bioinformatics* 23, 3350–3355, doi:10.1093/bioinformatics/btm408.
- Dahari, H., Major, M., Zhang, X., Mihalik, K., Rice, C.M., Perelson, A.S., Feinstone, S.M., Neumann, A.U., 2005. Mathematical modeling of primary hepatitis C infection: noncytolytic clearance and early blockage of virion production. *Gastroenterology* 128, 1056–1066, doi:10.1053/j.gastro.2005.01.049.
- Duca, K.A., Shapiro, M., Delgado-Eckert, E., Hadinoto, V., Jarrar, A.S., Laubenbacher, R., Lee, K., Luzuriaga, K., Polys, N.F., Thorley-Lawson, D.A., 2007. A virtual look at Epstein-Barr virus infection: biological interpretations. *PLoS Pathog.* 3, 1388–1400, doi:10.1371/journal.ppat.0030137.
- Gilbert, N., Bankes, S., 2002. Platforms and methods for agent-based modelling. *Proc. Natl. Acad. Sci. U.S.A.* 99 (Suppl. 3), 7197–7198.

- Ho, D.D., Neumann, A.U., Perelson, A.S., Chen, W., Leonard, J.M., Markowitz, M., 1995. Rapid turnover of plasma virions and CD4 lymphocytes in HIV-1 infection. *Nature* 373, 123–126, doi:10.1038/373123a0.
- Naniche, D., 2009. Human immunology of measles virus infection. *Curr. Top. Microbiol. Immunol.* 330, 151–171.
- Neumann, A.U., Lam, N.P., Dahari, H., Gretch, D.R., Wiley, T.E., Layden, T.J., Perelson, A.S., 1998. Hepatitis C viral dynamics in vivo and the antiviral efficacy of interferon-alpha therapy. *Science* 282, 103–107, doi:10.1126/science.282.5386.103.
- Nowak, M.A., Bonhoeffer, S., Hill, A.M., Boehme, R., Thomas, H.C., McDade, H., 1996. Viral dynamics in hepatitis B virus infection. *Proc. Natl. Acad. Sci. U.S.A.* 93, 4398–4402.
- Nowak, M.A., Lloyd, A.L., Vasquez, G.M., Wiltout, T.A., Wahl, L.M., Bischofberger, N., Williams, J., Kinter, A., Fauci, A.S., Hirsch, V.M., Lifson, J.D., 1997. Viral dynamics of primary viremia and antiretroviral therapy in simian immunodeficiency virus infection. *J. Virol.* 71, 7518–7525.
- Shapiro, M., Duca, K.A., Lee, K., Delgado-Eckert, E., Hawkins, J., Jarrah, A.S., Laubacher, R., Polys, N.F., Hadinoto, V., Thorley-Lawson, D.A., 2008. A virtual look at Epstein-Barr virus infection: simulation mechanism. *J. Theor. Biol.* 252, 633–648, doi:10.1016/j.jtbi.2008.01.032.
- See, H., Wark, P., 2008. Innate immune response to viral infection of the lungs. *Paediatr. Respir. Rev.* 9, 243–250, doi:10.1016/j.prrv.2008.04.001.

Genome-wide association of *IL28B* with response to pegylated interferon- α and ribavirin therapy for chronic hepatitis C

Yasuhiro Tanaka^{1,18}, Nao Nishida^{2,18}, Masaya Sugiyama¹, Masayuki Kurosaki³, Kentaro Matsuura¹, Naoya Sakamoto⁴, Mina Nakagawa⁴, Masaaki Korenaga⁵, Keisuke Hino⁵, Shuhei Hige⁶, Yoshito Ito⁷, Eiji Mita⁸, Eiji Tanaka⁹, Satoshi Mochida¹⁰, Yoshikazu Murawaki¹¹, Masao Honda¹², Akito Sakai¹², Yoichi Hiasa¹³, Shuhei Nishiguchi¹⁴, Asako Koike¹⁵, Isao Sakaida¹⁶, Masatoshi Imamura¹⁷, Kiyooki Ito¹⁷, Koji Yano¹⁷, Naohiko Masaki¹⁷, Fuminaka Sugauchi¹, Namiki Izumi³, Katsushi Tokunaga² & Masashi Mizokami^{1,17}

The recommended treatment for patients with chronic hepatitis C, pegylated interferon- α (PEG-IFN- α) plus ribavirin (RBV), does not provide sustained virologic response (SVR) in all patients. We report a genome-wide association study (GWAS) to null virological response (NVR) in the treatment of patients with hepatitis C virus (HCV) genotype 1 within a Japanese population. We found two SNPs near the gene *IL28B* on chromosome 19 to be strongly associated with NVR (rs12980275, $P = 1.93 \times 10^{-13}$, and rs8099917, 3.11×10^{-15}). We replicated these associations in an independent cohort (combined P values, 2.84×10^{-27} (OR = 17.7; 95% CI = 10.0–31.3) and 2.68×10^{-32} (OR = 27.1; 95% CI = 14.6–50.3), respectively). Compared to NVR, these SNPs were also associated with SVR (rs12980275, $P = 3.99 \times 10^{-24}$, and rs8099917, $P = 1.11 \times 10^{-27}$). In further fine mapping of the region, seven SNPs (rs8105790, rs11881222, rs8103142, rs28416813, rs4803219, rs8099917 and rs7248668) located in the *IL28B* region showed the most significant associations ($P = 5.52 \times 10^{-28}$ – 2.68×10^{-32} ; OR = 22.3–27.1). Real-time quantitative PCR assays in peripheral blood mononuclear cells showed lower *IL28B* expression levels in individuals carrying the minor alleles ($P = 0.015$).

Hepatitis C is a global health problem that affects a significant proportion of the world's population. The World Health Organization

estimated that in 1999, there were 170 million HCV carriers worldwide, with 3–4 million new cases appearing each year. HCV infection affects more than 4 million people in the United States, where it represents the leading cause of cirrhosis and hepatocellular carcinoma as well as the leading cause of liver transplantation¹. The American Gastroenterological Association estimated that drugs are the largest direct costs of hepatitis C¹.

The most effective current standard of care in patients with chronic hepatitis C, a combination of PEG-IFN- α with ribavirin, does not produce SVR in all patients treated. Large-scale studies on 48-week-long PEG-IFN- α /RBV treatment in the United States and Europe showed that 42–52% of patients with HCV genotype 1 achieved SVR^{2–4}, and similar results were found in Japan. However, older patients (greater than 50 years of age) had a significantly lower rate of SVR due to poor adherence resulting from adverse events and laboratory-detectable abnormalities such as neutropenia and thrombocytopenia^{5,6}. Specifically, various well-described side effects (such as a flu-like syndrome, hematologic abnormalities and adverse neuropsychiatric events) often necessitate dose reduction, and 10–14% of patients require premature withdrawal from interferon-based therapy⁷. To avoid these side effects in patients who will not be helped by the treatment, as well as to reduce the substantial cost of PEG-IFN- α /RBV treatment, it would be useful to be able to predict an individual's response before or early in treatment. Several viral factors, such as genotype 1, high baseline viral load, viral

¹Department of Clinical Molecular Informative Medicine, Nagoya City University Graduate School of Medical Sciences, Nagoya, Japan. ²Department of Human Genetics, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan. ³Division of Gastroenterology and Hepatology, Musashino Red Cross Hospital, Tokyo, Japan. ⁴Department of Gastroenterology and Hepatology, Tokyo Medical and Dental University, Tokyo, Japan. ⁵Division of Hepatology and Pancreatology, Kawasaki Medical College, 577 Matsushima, Kurashiki, Japan. ⁶Department of Internal Medicine, Hokkaido University Graduate School of Medicine, Sapporo, Japan. ⁷Molecular Gastroenterology and Hepatology, Kyoto Prefectural University of Medicine, Kyoto, Japan. ⁸National Hospital Organization Osaka National Hospital, Osaka, Japan. ⁹Department of Medicine, Shinshu University School of Medicine, Matsumoto, Japan. ¹⁰Division of Gastroenterology and Hepatology, Internal Medicine, Saitama Medical University, Saitama, Japan. ¹¹Second department of Internal Medicine, Faculty of Medicine, Tottori University, Yonago, Japan. ¹²Department of Gastroenterology, Kanazawa University Graduate School of Medicine, Kanazawa, Japan. ¹³Department of Gastroenterology and Metabolism, Ehime University Graduate School of Medicine, Ehime, Japan. ¹⁴Department of Internal Medicine, Hyogo College of Medicine, Nishinomiya, Japan. ¹⁵Central Research Laboratory, Hitachi Ltd., Kokubunji, Japan. ¹⁶Gastroenterology and Hepatology, Yamaguchi University Graduate School of Medicine, Yamaguchi, Japan. ¹⁷Research Center for Hepatitis and Immunology, International Medical Center of Japan Konodai Hospital, Ichikawa, Japan. ¹⁸These authors contributed equally to this work. Correspondence should be addressed to M.M. (mmizokami@imcj2.hosp.go.jp).

Received 29 June; accepted 21 August; published online 13 September 2009; doi:10.1038/ng.449



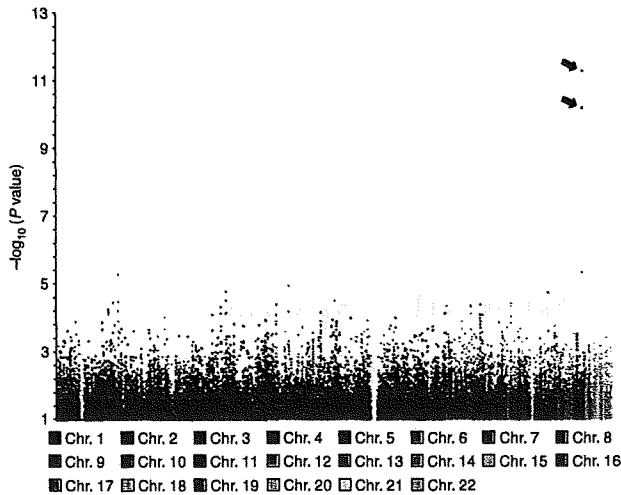


Figure 1 Genome-wide association results with PEG-IFN- α /RBV treatment in 142 Japanese patients with HCV (78 NVR and 64 VR samples). P values were calculated by using a χ^2 test for allele frequencies. The dots with arrows for chromosome 19 denote SNPs that showed significant genome-wide associations ($P < 8.05 \times 10^{-8}$) with response to PEG-IFN- α /RBV treatment.

kinetics during treatment, and amino acid pattern in the interferon sensitivity-determining region, have been reported to be significantly associated with the treatment outcome in a number of independent studies^{8–10}. Studies have also provided strong evidence that ~20% of patients with HCV genotype 1 and 5% of patients with genotype 2 or 3 have a null response to PEG-IFN- α /RBV. No definite predictor of this resistance is currently available that make it possible to bypass the initial 12–24 weeks' treatment before deciding whether treatment should be continued. If a reliable predictor of non-response were identified for use in patients before treatment initiation, then an estimated 20%, including those who have little or no chance to achieve SVR, could be spared the side effects and cost of treatment.

Host factors, including age, sex, race, liver fibrosis and obesity, have also been reported to be associated with PEG-IFN- α /RBV therapy outcome^{11,12}. However, little is known about the host genetic factors that might be associated with the response to therapy; thus far only

a few candidate genes, including those encoding type I interferon receptor-1 (*IFNAR1*) and mitogen-activated protein kinase-activated protein kinase 3 (*MAPKAPK3*), have been reported to be associated with treatment response^{13,14}. We describe here a GWAS for response to PEG-IFN- α /RBV treatment.

We conducted this GWAS to identify host genes associated with response to PEG-IFN- α /RBV treatment in 154 Japanese patients with HCV genotype 1 (82 with NVR and 72 with virologic response (VR), based on the selection criteria as described in Online Methods). We used the Affymetrix SNP 6.0 genome-wide SNP typing array for 900,000 SNPs. A total of 621,220 SNPs met the following criteria: (i) SNP call rate $\geq 95\%$, (ii) minor allele frequency (MAF) $\geq 1\%$ and (iii) deviation from Hardy-Weinberg equilibrium (HWE) $P \geq 0.001$ in VR samples. After excluding 4 NVR and 8 VR samples that showed quality control (QC) call rates of $< 95\%$, 78 NVR and 64 VR samples were included in the association analysis. **Figure 1** shows a genome-wide view of the single-point association data based on allele frequencies. Two SNPs located close to *IL28B* on chromosome 19 showed strong associations, with a minor allele dominant model (rs12980275, $P = 1.93 \times 10^{-13}$, and rs8099917, $P = 3.11 \times 10^{-15}$, respectively), with NVR to PEG-IFN- α /RBV treatment (**Table 1**). The rs8099917 lies between *IL28B* and *IL28A*, ~8 kb downstream from *IL28B* and ~16 kb upstream from *IL28A*. These associations reached genome-wide levels of significance for both SNPs in this initial GWAS cohort (Bonferroni criterion $P < 8.05 \times 10^{-8}$ (0.05/621,220)). The frequencies of minor allele-positive patients were much higher in the NVR group than in the VR group for both SNPs (74.3% in NVR, 12.5% in VR for rs12980275; 75.6% in NVR, 9.4% in VR for rs8099917). Notably, individuals homozygous for the minor allele were observed only in the NVR group. The VR group, as compared to the NVR group, showed genotype frequencies closer to those in the healthy Japanese population¹⁵, yet the minor allele frequencies were slightly higher in the transient virologic response (TVR) group (23.1%, 15.4%) than in the SVR group (9.8%, 7.8%) (**Table 1**). We applied the Cochran-Armitage test on all the SNPs and found a genetic inflation factor, λ , of 1.029 for the GWAS stage (**Supplementary Fig. 1**). We also carried out principal component analysis in 142 samples for the GWAS stage together with the HapMap samples (CEU, YRI, CHB and JPT) (**Supplementary Fig. 2**); this suggested that the effect of population stratification was negligible.

Table 1 Significant association of two SNPs (rs12980275 and rs8099917) with response to PEG-IFN- α /RBV treatment

dbSNP rsID	Nearest gene	MAF ^b (allele)	Allele (1/2)	Stage	Null responder (NVR ^a , n = 128)			Responder (VR ^a , n = 186)			Responder (SVR ^a , n = 140)			NVR vs. VR		NVR vs. SVR	
					11	12	22	11	12	22	11	12	22	OR (95% CI) ^c	P value ^d	OR (95% CI) ^c	P value ^d
rs12980275	<i>IL28B</i>	0.15 (G)	A/G	GWAS	20	54	4	56	8	0	46	5	0	20.3	1.93×10^{-13}	26.7	7.41×10^{-13}
					(25.6)	(69.2)	(5.1)	(87.5)	(12.5)	(0.0)	(90.2)	(9.8)	(0.0)	(8.3–49.9)		(9.3–76.5)	
					10	37	3	101	21	0	73	16	0	19.2	5.46×10^{-15}	18.3	8.37×10^{-13}
				Replication	(20.0)	(74.0)	(6.0)	(82.8)	(17.2)	(0.0)	(82.0)	(18.0)	(0.0)	(8.3–44.4)		(7.6–44.0)	
				Combined	30	91	7	157	29	0	119	21	0	17.7	2.84×10^{-27}	18.5	3.99×10^{-24}
					(23.4)	(71.1)	(5.5)	(84.4)	(15.6)	(0.0)	(85.0)	(15.0)	(0.0)	(10.0–31.3)		(10.0–34.4)	
rs8099917	<i>IL28B</i>	0.12 (G)	T/G	GWAS	19	56	3	58	6	0	47	4	0	30.0	3.11×10^{-15}	36.5	5.00×10^{-14}
					(24.4)	(71.8)	(3.8)	(90.6)	(9.4)	(0.0)	(92.2)	(7.8)	(0.0)	(11.2–80.5)		(11.6–114.6)	
					11	37	2	108	14	0	78	11	0	27.4	9.47×10^{-18}	25.1	1.00×10^{-14}
				Replication	(22.0)	(74.0)	(4.0)	(88.5)	(11.5)	(0.0)	(87.6)	(12.4)	(0.0)	(11.5–65.3)		(10.0–63.1)	
				Combined	30	93	5	166	20	0	125	15	0	27.1	2.68×10^{-32}	27.2	1.11×10^{-27}
					(23.4)	(72.7)	(3.9)	(89.2)	(10.8)	(0.0)	(89.3)	(10.7)	(0.0)	(14.6–50.3)		(13.9–53.4)	

^aNVR, null virologic response; VR, virologic response; SVR, sustained virologic response. The 186 VRs consisted of 46 transient virologic response (TVRs) and 140 SVRs. ^bMinor allele frequency and minor allele in 184 healthy Japanese individuals¹⁵. The MAF of the SNPs in SVR is similar to that of TVR group, whereas that of NVR is much higher (76.6%). ^cOdds ratio for the minor allele in a dominant model. ^d P value by χ^2 test for the minor allele dominant model.