

**Figure 1.** Predicted values of selected markers *vs.* radiation dose (based on individual regression models, Eq. 1). A) ROS. B) IL-6. C) CRP. D) ESR.

### Principal factor analyses of cytokines and markers

Principal factor analysis revealed that 3 factors explained different characteristics among cytokines and markers. The most important factor explained 31% of the total variation and had large coefficients for Log(ROS), Log(IL-6), Log(CRP), and Log(ESR); the second factor explained 14% of the total variation and had large coefficients for Log(TNF- $\alpha$ ), Log(ESR), and Log(Igs); and the third factor explained 13% of the total variation and had large coefficients for Log(TNF- $\alpha$ ), Log(IL-4), and Log(IL-10). This implies the presence of some latent factors possibly related to different immunological pathways or mechanisms, each involving a unique set of the cytokines and markers (that is, those with large coefficients for that particular factor). Then, the principal component analysis was conducted separately for the aforementioned cytokines and markers strongly associated with each of the first, second, and third factors to create scores quantifying the information of that latent factor. **Tables 3–5** show these 3 correlation matrices used in the principal component analyses. All correlations are positive and lie between 0.12 and 0.42, indicating a low to moderate interrelationship among the variables (37).

**TABLE 3.** Correlation matrix for logarithms of ROS, IL-6, CRP, and ESR

Variable	Log(IL-6)	Log(CRP)	Log(ESR)
Log(ROS)	0.329	0.254	0.286
Log(IL-6)	–	0.410	0.327
Log(CRP)	–	–	0.312

### Regression model for inflammation scores

Only the first principal component of each of the 3 principal component analyses contained more information than the original single cytokines and markers. We call those first principal components the 1st\_score, 2nd\_score, and 3rd\_score, which quantify the first, second, and third unobserved latent inflammation mechanisms implied by factor analysis, respectively. These scores are obtained by Eqs. 2–4:

$$\begin{aligned} \text{1st\_score} = & 0.463 * s\_Log(\text{ROS}) \\ & + 0.539 * s\_Log(\text{IL-6}) + 0.509 * s\_Log(\text{CRP}) \\ & + 0.486 * s\_Log(\text{ESR}) \quad (2) \end{aligned}$$

$$\begin{aligned} \text{2nd\_score} = & 0.527 * s\_Log(\text{TNF-}\alpha) \\ & + 0.601 * s\_Log(\text{ESR}) + 0.601 * s\_Log(\text{Igs}) \quad (3) \end{aligned}$$

$$\begin{aligned} \text{3rd\_score} = & 0.636 * s\_Log(\text{TNF-}\alpha) \\ & + 0.445 * s\_Log(\text{IL-4}) \\ & + 0.630 * s\_Log(\text{IL-10}) \quad (4) \end{aligned}$$

where the prefix *s* is used to indicate that individual (logged) variables have been standardized to have 0 means and *sd* of 1. We view these scores as single multivariate indices of inflammation, each of which incorporates the inflammation characteristics of the associated group of variables into a single variable. With

**TABLE 4.** Correlation matrix for logarithms of TNF- $\alpha$ , ESR, and Igs

Variable	Log(ESR)	Log(Igs)
Log(TNF- $\alpha$ )	0.300	0.301
Log(ESR)	–	0.421

TABLE 5. Correlation matrix for logarithms of TNF- $\alpha$ , IL-4, and IL-10

Variable	Log (IL-4)	Log (IL-10)
Log (TNF- $\alpha$ )	0.133	0.277
Log (IL-4)	–	0.125

$y$  as scores defined in Eqs. 2–4, we fit the model Eq. 1, resulting, in each case, in a prediction equation. As shown in **Table 6**, all 3 models are statistically significant (with  $P < 0.001$ ), with  $R^2 = 19.0, 19.1,$  and  $8.8\%$ , respectively: considerably better than models for most of the individual markers, thereby implying superior fits and better predictability. More important, the inflammation scores are independent of each other. The associations of the 1st\_score, 2nd\_score, and 3rd\_score with the explanatory variables; *i.e.*, radiation, gender, age, smoking, and BMI, in terms of the regression model Eq. 1, are also shown in Table 6. All the variables in the model of the 1st\_score are statistically significant, with very small  $P$  values, an exception being the BMI, for which the statistical significance is somewhat marginal. More specifically, radiation, gender, and age are statistically significant in both models for the 1st\_score and 2nd\_score. Also, radiation and age are significant for the 3rd\_score. In addition, this score may be closely associated with BMI, but not with gender and smoking. Standard error values, also presented in Table 6, of all the variables are small, thereby supporting, as desired, the good precision of the estimates of the respective slopes. **Figure 2** shows the plots of predicted values of the 1st\_score (Fig. 2A, B), 2nd\_score (Fig. 2C, D), and 3rd\_score (Fig. 2E, F) among study subjects against their given dose and age at the time of collecting plasma samples, respectively. In each of the 6 images, an increasing trend for the score is self-evident. However, in terms of the scatter in these plots, the effect of aging is more pronounced than the effect of radiation, since points are more tightly clustered around the trend in the former case.

### Inflammation scores and CD4 or naive CD4 T-cell frequencies

We also investigated the association between the inflammation scores and CD4 or naive CD4 T-cell frequencies, respectively, by adding CD4 T-cell frequencies and naive CD4 T-cell frequencies, separately, into the regression models mentioned above (**Fig. 3**). The estimates of the regression coefficient in each model were  $-0.022$  (CD4,  $P < 0.001$ ) and  $-0.016$  (naive CD4,  $P = 0.008$ ) for the 1st\_score,  $-0.017$  (CD4,  $P = 0.005$ ) and  $-0.015$  (naive CD4,  $P = 0.010$ ) for the 2nd\_score, and  $-0.006$  (CD4,  $P = 0.24$ ) and  $-0.005$  (naive CD4,  $P = 0.38$ ) for the 3rd\_score. There were negative correlations between inflammation scores and helper T-cell measurements, while only the associations of the 1st\_score and 2nd\_score were statistically significant.

### Radiation effects on age-dependent inflammation

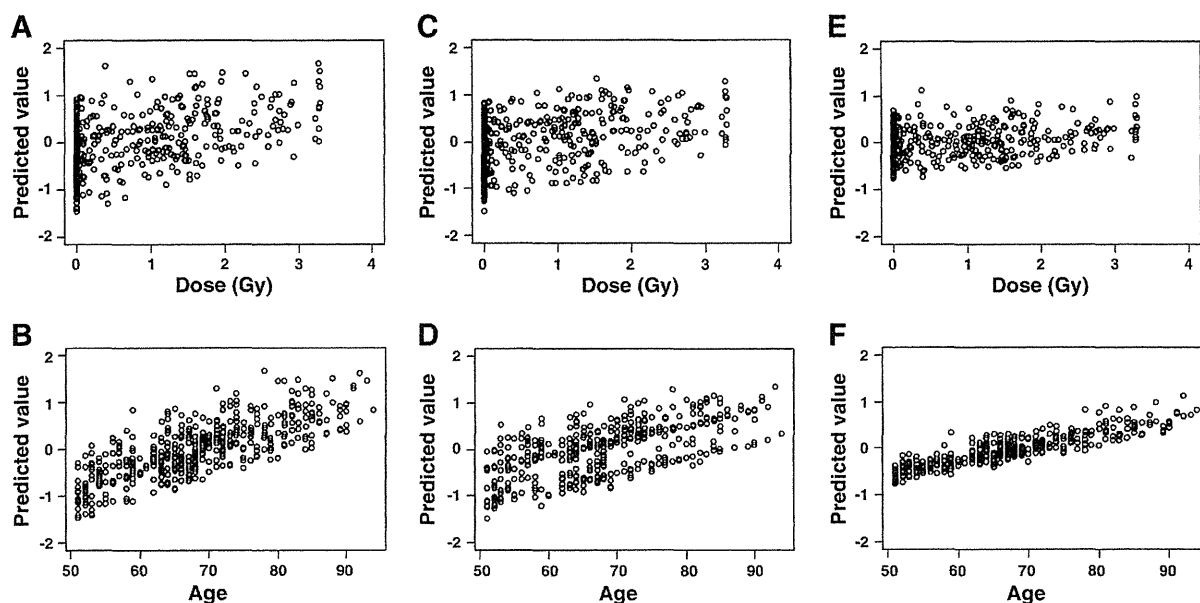
The regression slope coefficients in the fitted model shown in Table 6 can also be used to compare the effect of radiation dose with the corresponding effect of aging. Specifically, for the inflammation pathway quantified by the 1st\_score, where  $\beta_{\text{dose}} = 0.33$  and  $\beta_{\text{age}} = 0.48$ , when viewed in the context of the inflammation explained by the 4 markers in the 1st\_score, 1 Gy of radiation exposure is approximately equivalent to  $\beta_{\text{dose}}/\beta_{\text{age}} = 0.33/0.48 = 0.69$  decade of aging, or an age increase of 6.9 yr. In view of the small standard error values of both slope coefficients, this estimate is also expected to be reasonably accurate. Similarly for the inflammation pathways quantified by the 2nd\_score and 3rd\_score, 1 Gy of radiation exposure is approximately equivalent to  $0.27/0.33 = 0.82$  and  $0.12/0.30 = 0.40$  decades of aging, respectively. For the 2nd\_score and 3rd\_score, we calculated these ratios using the coefficients of models chosen with minimum Akaike's information criterion (AIC). The percentage increments in each inflammation score for the unit change in dose or age can be measured as  $100 * (\exp(\beta_{\text{dose}}) - 1)$  and  $100 * (\exp(\beta_{\text{age}}) - 1)$ , respectively. Accordingly, corresponding to 1 Gy of radiation exposure, the percentage increments in 1st\_score, 2nd\_score, and 3rd\_score are approximately 39, 31, and 13%, respectively, while for one decade of extra age, the score increments are approximated as 62, 39, and 35%, respectively.

### DISCUSSION

In this study, we measured plasma ROS levels in A-bomb survivors and used those results and previously measured results of plasma markers and cytokines for evaluation of radiation and aging effects on immune

TABLE 6. Regression model for scores of inflammation pathways (based on leading principal components as response)

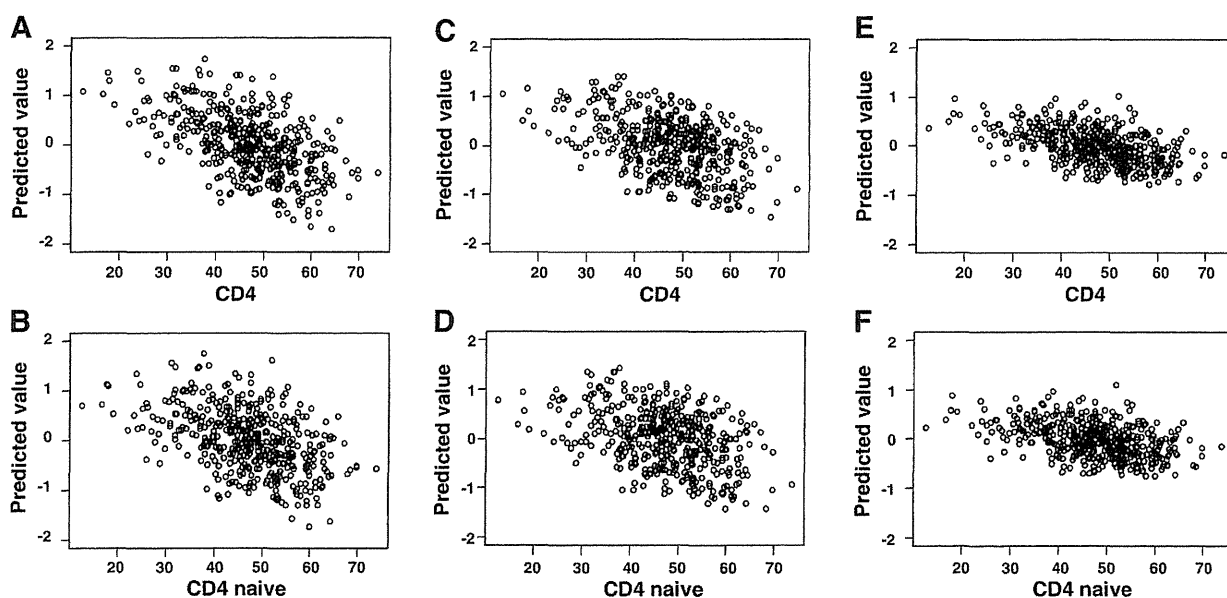
Factor	Estimate	95% CI	$P$	$R^2$
1st_score				
Dose (Gy)	0.332	(0.204, 0.461)	<0.001	0.190
Gender	0.475	(0.213, 0.737)	<0.001	
Age (10 yr)	0.482	(0.367, 0.597)	<0.001	
Smoking	0.383	(0.091, 0.676)	0.011	
BMI	0.037	(0.001, 0.073)	0.045	
2nd_score				
Dose (Gy)	0.267	(0.148, 0.386)	<0.001	0.191
Gender	0.689	(0.446, 0.932)	<0.001	
Age (10 yr)	0.304	(0.214, 0.427)	<0.001	
Smoking	0.001	(-0.406, 0.136)	0.329	
BMI	0.035	(-0.039, 0.028)	0.745	
3rd_score				
Dose (Gy)	0.124	(0.011, 0.238)	0.032	0.088
Gender	-0.025	(-0.257, 0.207)	0.836	
Age (10 yr)	0.304	(0.203, 0.406)	<.001	
Smoking	0.001	(-2.257, 0.260)	0.991	
BMI	0.035	(0.004, 0.067)	0.030	



**Figure 2.** Predicted inflammation scores as functions of radiation dose and age (as calculated for all study subjects by the 1st, 2nd, and 3rd scores, comprised of ROS, IL-6, CRP, and ESR; TNF- $\alpha$ , ESR, and Igs; and TNF- $\alpha$ , IL-4, and IL-10, respectively). *A)* 1st\_score vs. dose. *B)* 1st\_score vs. age. *C)* 2nd\_score vs. dose. *D)* 2nd\_score vs. age. *E)* 3rd\_score vs. dose. *F)* 3rd\_score vs. age.

and systemic markers of inflammation. The analyses included both inflammatory cytokines and markers (ROS, IL-6, TNF- $\alpha$ , CRP, ESR, and Igs) and antiinflammatory cytokines (IL-4 and IL-10). IL-4, IL-6, TNF- $\alpha$ , and CRP play a central role in the coordination of the inflammatory response as key proinflammatory cytokines. In the previous study, plasma levels of inflammatory cytokines and biomarkers (IL-6, TNF- $\alpha$ , CRP, and ESR) increased significantly with radiation dose or age (32, 38). We demonstrated here that 7 of these 8

cytokines and markers, including plasma ROS levels, were correlated with increasing age and radiation dose. However, the coefficient of determination in each regression was small. Then, to assess the biologically significant relationship between inflammation and radiation exposure or aging, we used multivariate statistical analyses and evaluated the systemic markers of inflammation as scores calculated by linear combinations of selected cytokines and markers. Principal factor analysis revealed that 3 different factors explained the



**Figure 3.** Association between the inflammation scores and CD4 or CD4 naive T cell frequency. *A)* 1st\_score vs. CD4 T cells. *B)* 1st\_score vs. naive CD4 T cells. *C)* 2nd\_score vs. CD4 T cells. *D)* 2nd\_score vs. naive CD4 T cells. *E)* 3rd\_score vs. CD4 T cells. *F)* 3rd\_score vs. naive CD4 T cells.

correlation structure among cytokines and markers. The first, second, and third factors had large coefficients for ROS, IL-6, CRP, and ESR; TNF- $\alpha$ , ESR, and Igs; and TNF- $\alpha$ , IL-4, and IL-10, respectively. These factors imply the presence of some latent variables possibly related to immunological pathways or mechanisms, each involving the cytokines/markers with large coefficients for that factor. In the association of the 1st\_score, 2nd\_score, and 3rd\_score with the explanatory variables, radiation, gender, and age were statistically significant in the models of the 1st\_score and 2nd\_score. These results indicate that linear combinations of ROS, IL-6, CRP, and ESR and TNF- $\alpha$ , ESR, and Igs generated such a score that was the most indicative of inflammation and revealed clear dependences on radiation dose and aging and that there may be two inflammation pathways related to radiation. One pathway is involved in an ROS-dependent pathway related to IL-6 and CRP. The other is in an ROS-independent pathway related to TNF- $\alpha$  and Igs. Our results indicated that the systemic markers of inflammation might be accelerated by these ROS-dependent and -independent pathways.

The ROS measured in the plasma at the time of sampling indicates an equilibrium level that is constitutively produced by cells and reflects the oxidative state of the body. It has been known that ROS perform essential roles in inflammation and immune responses to pathogens, including bacterial killing through the production of superoxide by reduced nicotinamide adenine dinucleotide phosphate (NADPH) oxidases during the respiratory burst in activated macrophages and neutrophils (39,40). One possibility is that the pathways related to these NADPH oxidases may be involved in ROS generation and modulate by aging and radiation. Evidence that high serum concentrations of cytokines and inflammatory proteins were associated with high levels of ROS and low levels of superoxide dismutase and glutathione peroxidase is also of particular interest (41). Circulating levels of TNF- $\alpha$  and IL-6 increased with age in the general population (42); aging was also associated with increased levels of acute-phase proteins, such as CRP and serum amyloid A (43). Moreover, whereas IFN- $\alpha$  and IFN- $\gamma$  production decreased in the elderly, IL-4 and IL-10 production increased (44), as did production of proinflammatory cytokines, such as IL-1, IL-6, IL-8, and TNF- $\alpha$  (45). Aging is also found to cause the increased levels of inflammatory cytokines and ROS, both of which have been associated with obesity, insulin resistance, and atherosclerosis (46, 47). Inflammation-induced formation of ROS is often indicative of the development of cardiovascular disease, diabetes, and cancer, and its implications for aging provide strong support for the hypothetical involvement of interrelated inflammation and oxidative stress in aging processes.

It has been reported that vascular endothelial cells might also be involved in production of many inflammation-related cytokines and ROS in plasma and that their extensive production has been implicated in

diseases such as atherosclerosis through the induction of chronic activation of the vascular endothelium and components of the immune system (48). It is likely that elevated plasma levels of CRP and IL-6 are risk factors for cardiovascular disease (49). Indeed, human CRP and complement activation are major mediators of ischemic myocardial injury (50). Experimental and clinical studies have suggested increased production of ROS both in animals and patients with acute and chronic heart failure (51–55). Our results indicate that a linear combination of ROS, IL-6, CRP, and ESR (Eq. 2) generated such a score that was the most indicative of inflammation and revealed clear dependences on radiation dose and aging. These results suggest that, collectively, radiation exposure may enhance the persistent inflammatory status of A-bomb survivors in conjunction with natural aging. Given the potential implications of our findings, a follow-up study with an increased number of subjects or retrospective study with the use of stored plasma samples, in association with the development of various inflammation-associated diseases, is warranted to confirm the clinical benefit of these scores.

It has been reported that a regression analysis indicates statistically significant association between radiation dose and leukocyte counts but not neutrophil counts (33). In addition, long-term immunological studies of A-bomb survivors have revealed that the percentages of CD4 (helper) T-cell populations, especially those of CD45RA<sup>+</sup> (naive) CD4 T cells, decreased in peripheral blood as radiation dose or age increased (34, 56). In our present studies, score values correlated negatively with the percentages of CD4 and naive CD4 T cells. We therefore suppose that elevated inflammatory score values indicate certain mechanisms which play a role in the attenuation of T-cell immunity in A-bomb survivors. Furthermore, we have previously reported that the percentages of CD4<sup>+</sup> T cells in peripheral blood lymphocytes decreased significantly in A-bomb survivors who had a history of myocardial infarction, indicating that such an immunological modification may be related to the increased risk of the disease (57). Taken together, this immunological balance impaired by aging and radiation exposure might result in acceleration of inflammatory status related to those multiple ROS-dependent and -independent inflammatory pathways.

To obtain further support of the hypothesis that A-bomb radiation has accelerated immunological aging, which, in turn, might lead to long-lasting inflammation, we are conducting a longitudinal analysis that assesses changes in immunological and clinical status with radiation exposure and aging in the A-bomb survivor population, on the basis of a number of immunological and inflammatory biomarkers interacting with each other. We have also begun to analyze genotype-phenotype associations concerning immune-related genes to take into account the genetically-regulated production and function of inflammatory cytokines and chemokines, which may in part explain

the interindividual variation in inflammatory response. Based on the studies mentioned above, it is expected that we will be able to derive the basis for the development of an integrated scoring system to verify our hypothesis of acceleration of immunosenescence due to radiation exposure. FJ

The Radiation Effects Research Foundation (RERF; Hiroshima and Nagasaki, Japan) is a private, nonprofit foundation funded by the Japanese Ministry of Health, Labor, and Welfare and the U.S. Department of Energy (DOE), the latter in part through DOE award DE-HS0000031 to the National Academy of Sciences. This was based on RERF Research Protocols 1-93, 4-02, and 3-09 and was supported in part by grants-in-aid for scientific research from the Japanese Ministry of Education, Culture, Sports Science, and Technology; the Japanese Ministry of Health, Labor, and Welfare; and the U.S. National Institute of Allergy and Infectious Diseases (NIAID; contract HHSN272200900059C). The views of the authors do not necessarily reflect those of the two governments.

## REFERENCES

- Douple, E. B., Mabuchi, K., Cullings, H. M., Preston, D. L., Kodama, K., Shimizu, Y., Fujiwara, S., and Shore, R. E. (2011) Long-term radiation-related health effects in a unique human population: lessons learned from the atomic bomb survivors of Hiroshima and Nagasaki. *Disaster Med. Public Health Prep.* 5(Suppl. 1), S122–S133
- Kodama, K., Sasaki, H., and Shimizu, Y. (1990) Trend of coronary heart disease and its relationship to risk factors in a Japanese population: a 26-year follow-up, Hiroshima/Nagasaki study. *Jpn. Circ. J.* 54, 414–421
- Pierce, D. A., Shimizu, Y., Preston, D. L., Vaeth, M., and Mabuchi, K. (1996) Studies of the mortality of atomic bomb survivors. Report 12, Part I. Cancer: 1950–1990. *Radiat. Res.* 146, 1–27
- Sharp, G. B., Mizuno, T., Cologne, J. B., Fukuhara, T., Fujiwara, S., Tokuoka, S., and Mabuchi, K. (2003) Hepatocellular carcinoma among atomic bomb survivors: significant interaction of radiation with hepatitis C virus infections. *Int. J. Cancer* 103, 531–537
- Preston, D. L., Ron, E., Tokuoka, S., Funamoto, S., Nishi, N., Soda, M., Mabuchi, K., and Kodama, K. (2007) Solid cancer incidence in atomic bomb survivors: 1958–1998. *Radiat. Res.* 168, 1–64
- Ozasa, K., Shimizu, Y., Suyama, A., Kasagi, F., Soda, M., Grant, E. J., Sakata, R., Sugiyama, H., and Kodama, K. (2012) Studies of the mortality of atomic bomb survivors, Report 14, 1950–2003: an overview of cancer and noncancer diseases. *Radiat. Res.* 177, 229–243
- Bretz, W. A., Weyant, R. J., Corby, P. M., Ren, D., Weissfeld, L., Kritchevsky, S. B., Harris, T., Kurella, M., Satterfield, S., Visser, M., and Newman, A. B. (2005) Systemic inflammatory markers, periodontal diseases, and periodontal infections in an elderly population. *J. Am. Geriatr. Soc.* 53, 1532–1537
- Ungvari, Z., Csiszar, A., and Kaley, G. (2004) Vascular inflammation in aging. *Herz.* 29, 733–740
- Krabbe, K. S., Pedersen, M., and Bruunsgaard, H. (2004) Inflammatory mediators in the elderly. *Exp. Gerontol.* 39, 687–699
- Harris, T. B., Ferrucci, L., Tracy, R. P., Corti, M. C., Wacholder, S., Ettinger, W. H., Jr., Heimovitz, H., Cohen, H. J., and Wallace, R. (1999) Associations of elevated interleukin-6 and C-reactive protein levels with mortality in the elderly. *Am. J. Med.* 106, 506–512
- Van der Meide, P. H., and Schellekens, H. (1996) Cytokines and the immune response. *Biotherapy* 8, 243–249
- Pannan, B. H., and Robotham, J. L. (1995) The acute-phase response. *New Horiz.* 3, 183–197
- Lu, Z. Y., Brailly, H., Wijdenes, J., Bataille, R., Rossi, J. F., and Klein, B. (1995) Measurement of whole body interleukin-6 (IL-6) production: prediction of the efficacy of anti-IL-6 treatments. *Blood* 86, 3123–3131
- Fey, G. H., Hattori, M., Hocke, G., Brechner, T., Baffet, G., Baumann, M., Baumann, H., and Northemann, W. (1991) Gene regulation by interleukin 6. *Biochimie (Paris)* 73, 47–50
- Tracey, K. J. (2002) The inflammatory reflex. *Nature* 420, 853–859
- Dinarello, C. A. (1991) Interleukin-1 and interleukin-1 antagonism. *Blood* 77, 1627–1652
- Papadaki, H. A., Eliopoulos, D. G., Ponticoglou, C., and Eliopoulos, G. D. (2001) Increased frequency of monoclonal gammopathy of undetermined significance in patients with nonimmune chronic idiopathic neutropenia syndrome. *Int. J. Hematol.* 73, 339–345
- Roxburgh, C. S., and McMillan, D. C. (2010) Role of systemic inflammatory response in predicting survival in patients with primary operable cancer. *Future Oncol.* 6, 149–163
- Macarthur, M., Hold, G. L., and El-Omar, E. M. (2004) Inflammation and Cancer II. Role of chronic inflammation and cytokine gene polymorphisms in the pathogenesis of gastrointestinal malignancy. *Am. J. Physiol. Gastrointest. Liver Physiol.* 286, G515–G520
- Willcox, J. K., Ash, S. L., and Catignani, G. L. (2004) Antioxidants and prevention of chronic disease. *Crit. Rev. Food Sci. Nutr.* 44, 275–295
- Peters, T., Weiss, J. M., Sindrilaru, A., Wang, H., Oreshkova, T., Wlaschek, M., Maity, P., Reimann, J., and Scharffetter-Kochanek, K. (2009) Reactive oxygen intermediate-induced pathomechanisms contribute to immunosenescence, chronic inflammation and autoimmunity. *Mech. Ageing Dev.* 130, 564–587
- Federico, A., Morgillo, F., Tuccillo, C., Ciardiello, F., and Loguercio, C. (2007) Chronic inflammation and oxidative stress in human carcinogenesis. *Int. J. Cancer.* 121, 2381–2386
- Chung, J., Lee, H. S., Chung, H. Y., Yoon, T. R., and Kim, H. K. (2008) Salicylideneamino-2-thiophenol inhibits inflammatory mediator genes (RANTES, MCP-1, IL-8 and HIF-1 $\alpha$ ) expression induced by tert-butyl hydroperoxide via MAPK pathways in rat peritoneal macrophages. *Biotechnol. Lett.* 30, 1553–1558
- Papa, S., Zazzeroni, F., Pham, C. G., Bubici, C., and Franzoso, G. (2004) Linking JNK signaling to NF-kappaB: a key to survival. *J. Cell Sci.* 117, 5197–5208
- Papa, S., Bubici, C., Zazzeroni, F., Pham, C. G., Kuntzen, C., Knabb, J. R., Dean, K., and Franzoso, G. (2006) The NF-kappaB-mediated control of the JNK cascade in the antagonism of programmed cell death in health and disease. *Cell Death Differ.* 13, 712–729
- Wullaert, A., Heyninck, K., and Beyaert, R. (2006) Mechanisms of crosstalk between TNF-induced NF-kappaB and JNK activation in hepatocytes. *Biochem. Pharmacol.* 72, 1090–1101
- Bubici, C., Papa, S., Dean, K., and Franzoso, G. (2006) Mutual cross-talk between reactive oxygen species and nuclear factor-kappa B: molecular basis and biological significance. *Oncogene* 25, 6731–6748
- Goossens, V., De Vos, K., Vercammen, D., Steemans, M., Vancompernelle, K., Fiers, W., Vandenebeele, P., and Grooten, J. (1999) Redox regulation of TNF signaling. *Biofactors* 10, 145–156
- Sone, H., Akanuma, H., and Fukuda, T. (2010) Oxygenomics in environmental stress. *Redox Rep.* 15, 98–114
- Cullings, H. M., Fujita, S., Funamoto, S., Grant, E. J., Kerr, G. D., and Preston, D. L. (2006) Dose estimation for atomic bomb survivor studies: its evolution and present status. *Radiat. Res.* 166, 219–254
- Hayashi, I., Morishita, Y., Imai, K., Nakamura, M., Nakachi, K., and Hayashi, T. (2007) High-throughput spectrophotometric assay of reactive oxygen species in serum. *Mutat. Res.* 631, 55–61
- Hayashi, T., Morishita, Y., Kubo, Y., Kusunoki, Y., Hayashi, I., Kasagi, F., Hakoda, M., Kyoizumi, S., and Nakachi, K. (2005) Long-term effects of radiation dose on inflammatory markers in atomic bomb survivors. *Am. J. Med.* 118, 83–86
- Neriishi, K., Nakashima, E., and Delongchamp, R. R. (2001) Persistent subclinical inflammation among A-bomb survivors. *Int. J. Radiat. Biol.* 77, 475–482

34. Kusunoki, Y., Kyoizumi, S., Hirai, Y., Suzuki, T., Nakashima, E., Kodama, K., and Seyama, T. (1998) Flow cytometry measurements of subsets of T, B and NK cells in peripheral blood lymphocytes of atomic bomb survivors. *Radiat. Res.* **150**, 227–236
35. Khattree, R., and Naik, D. N. (1999) *Applied Multivariate Statistics with SAS Software*, SAS Press/Wiley, Cary, NC, USA
36. Khattree, R., and Naik, D. N. (2000) *Multivariate Data Reduction and Discrimination with SAS Software*, SAS Press/Wiley, Cary, NC, USA
37. Bishop, Y. M., Fienberg, S. E., Holland, P. W., and Light, R. J. (2007) *Discrete Multivariate Analysis: Theory and Practice*, Springer, London
38. Hayashi, T., Kusunoki, Y., Hakoda, M., Morishita, Y., Kubo, Y., Maki, M., Kasagi, F., Kodama, K., Macphee, D. G., and Kyoizumi, S. (2003) Radiation dose-dependent increases in inflammatory response markers in A-bomb survivors. *Int. J. Radiat. Biol.* **79**, 129–136
39. Lambeth, J. D. (2004) NOX enzymes and the biology of reactive oxygen. *Nat. Rev. Immunol.* **4**, 181–189
40. Kanayama, A., and Miyamoto, Y. (2007) Apoptosis triggered by phagocytosis-related oxidative stress through FLIPS down-regulation and JNK activation. *J. Leukoc. Biol.* **82**, 1344–1352
41. Mantovani, G., Maccio, A., Madeddu, C., Mura, L., Gramignano, G., Lusso, M. R., Mulas, C., Mudu, M. C., Murgia, V., Camboni, P., Massa, E., Ferrelli, L., Contu, P., Rinaldi, A., Sanjust, E., Atzei, D., and Elsener, B. (2002) Quantitative evaluation of oxidative stress, chronic inflammatory indices and leptin in cancer patients: correlation with stage and performance status. *Int. J. Cancer* **98**, 84–91
42. Dobbs, R. J., Charlett, A., Purkiss, A. G., Dobbs, S. M., Weller, C., and Peterson, D. W. (1999) Association of circulating TNF-alpha and IL-6 with ageing and parkinsonism. *Acta Neurol. Scand.* **100**, 34–41
43. Ballou, S. P., Lozanski, F. B., Hodder, S., Rzewnicki, D. L., Mion, L. C., Sipe, J. D., Ford, A. B., and Kushner, I. (1996) Quantitative and qualitative alterations of acute-phase proteins in healthy elderly persons. *Age Ageing* **25**, 224–230
44. Solana, R., and Mariani, E. (2000) NK and NK/T cells in human senescence. *Vaccine* **18**, 1613–1620
45. Rink, L., Cakman, I., and Kirchner, H. (1998) Altered cytokine production in the elderly. *Mech. Ageing Dev.* **102**, 199–209
46. Pedersen, B. K., and Bruunsgaard, H. (2003) Possible beneficial role of exercise in modulating low-grade inflammation in the elderly. *Scand. J. Med. Sci. Sports* **13**, 56–62
47. Cubbon, R. M., Kahn, M. B., and Wheatcroft, S. B. (2009) Effects of insulin resistance on endothelial progenitor cells and vascular repair. *Clin. Sci. (London)* **117**, 173–190
48. Galkina, E., and Ley, K. (2009) Immune and inflammatory mechanisms of atherosclerosis (\*). *Ann. Rev. Immunol.* **27**, 165–197
49. Ridker, P. M., Rifai, N., Stampfer, M. J., and Hennekens, C. H. (2000) Plasma concentration of interleukin-6 and the risk of future myocardial infarction among apparently healthy men. *Circulation* **101**, 1767–1772
50. Griselli, M., Herbert, J., Hutchinson, W. L., Taylor, K. M., Sohail, M., Krausz, T., and Pepys, M. B. (1999) C-reactive protein and complement are important mediators of tissue damage in acute myocardial infarction. *J. Exp. Med.* **190**, 1733–1740
51. Cesselli, D., Jakoniuk, I., Barlucchi, L., Beltrami, A. P., Hintze, T. H., Nadal-Ginard, B., Kajstura, J., Leri, A., and Anversa, P. (2001) Oxidative stress-mediated cardiac cell death is a major determinant of ventricular dysfunction and failure in dog dilated cardiomyopathy. *Circ. Res.* **89**, 279–286
52. Ferrari, R., Guardigli, G., Mele, D., Percoco, G. F., Ceconi, C., and Curello, S. (2004) Oxidative stress during myocardial ischaemia and heart failure. *Curr. Pharm. Des.* **10**, 1699–1711
53. Giordano, F. J. (2005) Oxygen, oxidative stress, hypoxia, and heart failure. *J. Clin. Invest.* **115**, 500–508
54. Landmesser, U., Spiekermann, S., Dikalov, S., Tatge, H., Wilke, R., Kohler, C., Harrison, D. G., Hornig, B., and Drexler, H. (2002) Vascular oxidative stress and endothelial dysfunction in patients with chronic heart failure: role of xanthine-oxidase and extracellular superoxide dismutase. *Circulation* **106**, 3073–3078
55. Li, J. M., and Shah, A. M. (2004) Endothelial cell superoxide generation: regulation and relevance for cardiovascular pathophysiology. *Am. J. Physiol.* **287**, R1014–R1030
56. Kusunoki, Y., Yamaoka, M., Kasagi, F., Hayashi, T., Koyama, K., Kodama, K., MacPhee, D. G., and Kyoizumi, S. (2002) T cells of atomic bomb survivors respond poorly to stimulation by *Staphylococcus aureus* toxins in vitro: does this stem from their peripheral lymphocyte populations having a diminished naive CD4 T-cell content? *Radiat. Res.* **158**, 715–724
57. Kusunoki, Y., Kyoizumi, S., Yamaoka, M., Kasagi, F., Kodama, K., and Seyama, T. (1999) Decreased proportion of CD4 T cells in the blood of atomic bomb survivors with myocardial infarction. *Radiat. Res.* **152**, 539–543

Received for publication June 20, 2012.  
Accepted for publication July 24, 2012.

## METHODOLOGY

# Conventional case–cohort design and analysis for studies of interaction

John Cologne,<sup>1\*</sup> Dale L Preston,<sup>2</sup> Kazue Imai,<sup>3</sup> Munechika Misumi,<sup>1</sup> Kengo Yoshida,<sup>3</sup> Tomonori Hayashi<sup>3</sup> and Kei Nakachi<sup>4</sup>

<sup>1</sup>Department of Statistics, Radiation Effects Research Foundation, Hiroshima, Japan, <sup>2</sup>Hirosoft International Corporation, Seattle, WA, USA, <sup>3</sup>Department of Radiobiology/Molecular Epidemiology, Radiation Effects Research Foundation, Hiroshima, Japan and <sup>4</sup>Consultant, Radiation Effects Research Foundation, Hiroshima, Japan

\*Corresponding author. Department of Statistics, Radiation Effect Research Foundation, 5-2 Hijiyama Park, Minami-ku, Hiroshima 732-0815 Japan. E-mail: jcologne@rerf.jp

---

**Accepted** 1 June 2012

**Background** The case–cohort study design has received significant methodological attention in the statistical and epidemiological literature but has not been used as widely as other cohort-based sampling designs, such as the nested case–control design. Despite its efficiency and practicality for a wide range of epidemiological study purposes, researchers may not yet be aware of the fact that the design can be analysed using standard software with only minor adjustments. Furthermore, although the large number of options for design and analysis of case–cohort studies may be daunting, they can be reduced to a few simple recommendations.

**Methods** We review conventional methods for the design and analysis of case–cohort studies and describe empirical comparisons based on a study of radiation, gene polymorphisms and cancer in the Japanese atomic bomb survivor cohort.

**Results** Stratified, as opposed to simple, random subcohort selection is recommended, especially for studies of gene–environment interaction, which are notorious for lacking statistical power. Methods based on the score-unbiased exact pseudo-likelihood (or its analogue with stratified case–cohort data) are recommended for use in conjunction with the asymptotic variance estimator.

**Conclusions** We present an example of how to implement case–cohort analysis methods using SPSS, a popular statistical package that lacks some of the features necessary to directly adapt and implement published methods based on other software platforms. We also illustrate case–control analysis using Epicure, which provides greater risk-modelling flexibility than other software. Our conclusions and recommendations should help investigators to better understand and apply the case–cohort design in epidemiological research.

**Keywords** case–cohort study, pseudo-likelihood, stratification, statistical interaction, gene–environment interaction

---

## Introduction

A recent explosion in molecular genetic methods facilitates detailed investigation of mechanisms of radiation-related cancer, which is one of the most well-documented outcomes of radiation exposure in the Japanese atomic bomb survivors.<sup>1</sup> Using stored leukocytes obtained from survivors who attended clinical examinations in the Adult Health Study<sup>2</sup> conducted at the Radiation Effects Research Foundation, studies are being conducted to investigate various components of the immunogenome and cancers at several sites.<sup>3,4</sup> The standard approach to assessing risk is through the effect of radiation exposure on rates of disease or death using cohort follow-up. Analysis may be based on the proportional hazards model fit via the partial likelihood (PL).<sup>5</sup> Extensions to the proportional hazards model involving more general risk functions other than proportional hazards, such as excess relative risk (ERR) or excess absolute rates, are often used.<sup>6</sup>

Cohort studies involving costly or time-consuming methods to obtain covariate information, such as interviews, use of precious biological specimens, collecting hospital or employment records, or genotyping, can be carried out efficiently using sampling from the cohort non-cases because the cases are more influential on risk estimates and typically represent a relatively small fraction of the total cohort.<sup>7</sup> Selecting subjects from the cohort using the counter-matched nested case-control study design<sup>8</sup> can be effective, especially when studying interactions.<sup>9,10</sup> However, with separate studies for multiple outcomes (e.g. multiple sites of cancer to be studied individually) selecting control subjects separately for each outcome wastes resources and can be inefficient. An appealing alternative is the case-cohort design,<sup>11</sup> in which one selects a single subcohort from the initial cohort at its inception with a pre-specified sampling fraction, either randomly or using stratified random sampling, and adds in all cases that occur in the cohort outside the subcohort. Several discussions of the relative merits of case-cohort and nested case-control designs are available,<sup>12-15</sup> and both designs have been shown to be cost-effective for gene-disease association studies.<sup>16</sup> In this review, we address case-cohort study design and analysis focusing, in particular, on exposure-based stratified random sampling from the cohort as a means to increase statistical efficiency relative to simple random sampling.

Because the subcohort is a random or stratified random sample from the cohort, it can be analysed by itself as a cohort study using the ordinary PL, but this is clearly inefficient because of the loss of cases that occur outside the subcohort. Thus, the key feature of the case-cohort design is the inclusion of all cases that occur in the cohort regardless of whether they are in the selected subcohort. Because of such biased sampling with regard to case status, risk estimation using the ordinary PL is not appropriate, one

reason being that asymptotic distribution theory for the PL estimators breaks down because cases outside the subcohort induce non-nesting of the so-called sigma fields, or sets on which the counting process probability measure is defined.<sup>17</sup> However, that problem can be overcome using straightforward adjustments. Most importantly, naive estimates of parameter standard errors output by ordinary proportional hazards regression programmes will be incorrect because of the statistical dependencies between elements of the estimating equations that inflate the variance, but those can also be handled in a straightforward manner using influence residuals available in most regression programmes. The upshot of this is that, although the synthetic nature of a case-cohort design may seem baffling to researchers, there is no impediment to obtaining appropriate estimates of risk and standard errors as long as appropriate methods of analysis are applied, and there are methods for performing such analyses using standard software.

Conventional options for analysing randomly sampled case-cohort data include the exact pseudo-likelihood (PsL)<sup>11,18</sup> and approximate PsL,<sup>17,19</sup> with various adjustments to parameter standard error estimates. In addition, more recent approaches based on considering the case-cohort sample as a cohort study with missing data or with alternative weighting<sup>20-23</sup> or as the second phase in a two-phase design<sup>24</sup> have been proposed. A generalized approach to analysing various designs involving sampling from a cohort has also been proposed.<sup>25</sup> Although the conventional approaches have been available for some time and have been studied in some detail, the more recent approaches still require further evaluation and comparison among themselves and with the more conventional methods. Therefore, we believe it is possible to draw generalizations and conclusions regarding choice of design and method among conventional approaches, and that is the focus of this review.

The exact PsL is score unbiased and has been shown to have the best small-sample properties in simulation studies.<sup>26,27</sup> However, there was little difference among the various PsL estimators when subcohort size was not small (>15% sampling fraction relative to the full cohort).<sup>26</sup> Procedures for adjusting the standard errors of estimated parameters include the asymptotic method (based on the pseudo-score, analogous to likelihood-based methods) and the robust (empirical or 'sandwich') types of estimators. Both methods use delta-beta statistics (influence residuals) that can be obtained from most Cox regression software, and they are asymptotically equivalent. However, the robust variance estimates can be conservative (too large) in finite samples.<sup>28</sup>

If there are risk factors whose values are known in the entire cohort, an alternative approach is to construct the subcohort using stratified sampling. In that case, the methods appropriate for randomly selected subcohorts require further modification to achieve



unbiased parameter estimates and correct variance estimates. Borgan *et al.*<sup>29</sup> described three estimation methods, analogous to those for randomly sampled subcohorts. Their method III, known as the ‘swapper’ method because a case not in the subcohort swaps places with a member of the subcohort in the same sampling stratum, is the analogue of the exact PsL and is the only score-unbiased procedure among the three. Score unbiased means that the estimating equations have expected value zero; otherwise they can result in finite-sample bias even though they are statistically consistent (unbiased asymptotically). Although the computational complexity involved in analysing data from a stratified case-cohort study is not great,<sup>18</sup> it nevertheless requires some effort, and it is not clear what magnitude of efficiency gains may be achieved by stratification to guide investigators in choosing between simple and stratified subcohort selection. In addition, the score-unbiased swapper method for stratified data does not seem to be used, despite mention of a stratified design in some published works,<sup>27</sup> and it is not implemented in the case-cohort function ‘cch’<sup>30</sup> in the R package ‘survival’.<sup>31</sup>

The present report describes the conventional simple and stratified random-sampling approaches to case-cohort study design and compares them using actual data from a study of epithelial growth factor receptor (EGFR) gene repeat length polymorphism and radiation joint effects on lung cancer in atomic bomb survivors.<sup>4</sup> We provide simple recommendations as to which methods of analysis should be used, and we present examples of how to implement the methods using Epicure and SPSS, which are popular statistical packages among epidemiologists. We also provide S-plus code for implementing the swapper method for stratified case-cohort data, which makes critical use of a built-in function not available in R.

## Review of conventional case-cohort methodology

Early work focused on estimation of relative risk based on the analogy of a case-control comparison.<sup>32–34</sup> Subsequent analytical methods capitalized on the analogy with cohort follow-up and PL, although the unique sampling approach (over-selecting cases) leads to a PsL rather than the usual PL. The PsL may be used to fit the analogue of a PL analysis (e.g. the proportional hazards model). Prentice noted that parametric functions other than log-linear risk functions could be used, and the PsL can be extended to accommodate absolute rate differences and other non-proportional hazard models for estimating the hazard rate ratios for covariates.<sup>11</sup> However, with the exception of Epicure,<sup>6</sup> such models are not generally available in case-cohort analysis software.

Analyses of the case-cohort sample must adjust for bias introduced in the distributions of covariates used in calculating the denominator of the PsL, as the case-cohort sample is not population-based because of over-sampling of cases. Bias incurred by including cases outside the subcohort is corrected by not allowing those cases to contribute to risk sets other than their own (risk sets comprise cohort subjects still under observation and at risk at the times cases occur). Two approaches exist based on whether cases outside the subcohort are allowed to enter the denominators of the PsL for their own risk sets. With Prentice’s original method,<sup>11</sup> cases outside the subcohort appear in their own risk set denominators; this method is called the exact PsL. The method on which the asymptotic variance estimate was originally based<sup>17</sup> excludes cases not in the subcohort even from the denominators of their own risk sets; this method is called the approximate PsL. These inclusions or exclusions can be carried out by weighting the cases. With the exact PsL, cases outside the subcohort are given Weight 1 in their own risk set and Weight 0 otherwise; with the approximate PsL, cases outside the subcohort are given Weight 0 in all risk sets, including their own, although those cases still contribute to the numerators of the PsL in their own risk sets by virtue of being cases (defining the risk sets). In practice, this means that weighting must be accomplished by the use of offsets rather than prior weights because the latter will exclude non-subcohort cases altogether. The exact PsL is score unbiased, whereas the approximate PsL is not score unbiased and can result in some small-sample bias, particularly if the subcohort risk sets are small.<sup>18</sup>

In addition to case weighting, all subjects can be weighted using risk set weights based on follow-up time and cohort attrition: who among the subcohort subjects are still at risk as time passes. This type of weighting, proposed by Barlow, uses the inverse of the sampled fraction as a weight.<sup>35</sup> The original proposal updates the weights at each failure time (case occurrence or risk set) using time-dependent weights according to the surviving fraction in the sampled subcohort vis-à-vis the surviving fraction in the full cohort. This approximates the likelihood function denominator contribution that would pertain in the full cohort analysis and is intuitive based on counting process theory.<sup>36</sup>

Barlow’s robust variance estimate, a jackknife estimate related to the influence function, was shown to have good performance and said to be easier to calculate than either the bootstrap estimate<sup>37</sup> or the asymptotic variance estimate, the latter requiring calculation of the correlations among components of the scores (estimating equations). An advantage of the robust variance vis-à-vis the asymptotic variance is that the latter is difficult to compute in open cohorts where there is staggered entry.<sup>15</sup> Practical computation was via SAS using the PHREG procedure

with non-subcohort cases defined to enter the study just before their failure times and variances derived from the delta-beta residuals<sup>15</sup> (originally using martingale residuals<sup>35</sup>). Similar computations may be done in S-plus<sup>36</sup> or R. If the data are partitioned into risk sets (if there are time-dependent covariates, for example), the subject-specific residuals must be summed up separately for each subject over all risk sets in which the subject appears.

A simplification of the Barlow weighting approach is to use the single, original subcohort sampling fraction at all failure times (i.e. weights are fixed, not depending on risk set). This has the practical advantage that multiple records (one record for each subject for each risk set in which she or he is at risk) do not have to be constructed, as the number of records in the analysis data set can significantly increase if there are many risk sets (many cases of disease or death). As with the exact PsL estimator, the Barlow method weights each case with Weight 1 at its failure time whether or not it is in the subcohort. Therneau and Li suggested that cases in the subcohort be weighted by the appropriate risk set fraction just as the non-cases, rather than converting their weight to 1 at their time of failure, to reflect the fact that they are legitimate risk set members in the sampled data.<sup>19</sup>

Details of variance estimation are described in a number of publications.<sup>17,19,35,38</sup> The asymptotic variance consists of within-risk set variances and between-risk set covariances, the latter arising from correlations among score contributions (components of the estimating equations). Briefly, the problem is solved by adding to the usual PL parameter covariance matrix a weighted product of influence residuals for the subcohort data. One approach uses weights based on the subcohort sampling fraction.<sup>29</sup> Another uses as weights the proportion of cases actually sampled in the subcohort.<sup>19</sup> With a large amount of data, the two types of weights will be similar because the expected fraction of cases in the subcohort approaches the subcohort sampling fraction as the subcohort size increases. However, due to random sampling variation, the actual proportion of cases in the subcohort may deviate from the expected value. Essentially all of the so-called 'robust' variance estimators that have been proposed can be expressed as special cases of a general influence function variance estimator that is robust to misspecification of the model relating the risk parameters to the underlying distribution function.<sup>38</sup>

All of the aforementioned methods consistently estimate the risk and variance asymptotically (with large studies), but they differ in their efficiency (estimates of the parameter variances) and small- or finite-sample bias (in the case of the approximate PsL). Therefore, a major concern is how well they perform in studies of the size typically encountered in epidemiological practice. According to Onland-Moret *et al.*,<sup>26</sup> the 'most simple approach'—Prentice's

original, exact PsL method—provides estimates and standard errors closest to the full-cohort results (i.e. using all subjects with no case-cohort sampling) when the cohort is small (about 1000 subjects or less) or when the subcohort is small (sampling fraction 0.15 or less with a cohort size of 1000 subjects). Those authors also demonstrated that the approximate PsL is inferior to the exact PsL method. However, they did not implement Barlow's method as originally proposed: they used the initial sampling fraction rather than time-dependent weights. Langholz and Jiao<sup>18</sup> reported that the approximate PsL estimator does not perform as well as the exact estimator when risk set sizes are small (e.g. when there are large numbers of cases and a small subcohort sampling fraction). The approximate PsL estimator is also not score unbiased in finite samples,<sup>29</sup> which means that it can result in small-sample bias.

From a practical standpoint, either the exact or approximate PsL involves about the same amount of data manipulation, and both require constructing duplicate data records for subcohort cases to accommodate the case weighting described earlier. The approximate method requires one set of records for all subcohort members (including subcohort cases) having Weights 1, and a second set of records for all cases (including both subcohort and non-subcohort cases) having Weights 0. Weight 0 is achieved by including a large, negative offset,  $-100$  say, which when included in the proportional hazards model becomes  $\exp\{-100\}$  or a number essentially equal to zero to the level of machine precision. A weight of one is achieved by using an offset of zero ( $\exp\{0\} = 1$ ). This means that cases not in the subcohort will define risk sets but will not enter into the denominators of the PsL calculations, whereas subcohort cases are included in all risk set denominators up to (in time) and including their own. This data setup was illustrated by Therneau and Li, who also described how to calculate the asymptotic variance estimates, using SAS and S-plus.<sup>19</sup> Their S-plus routine calculates the variances directly, whereas their SAS routine requires the IML procedure and an additional calculation by hand.

With the exact PsL, non-subcohort cases should be included in the denominators of their own risk sets. Instead of calculating offsets to down weight the cases, it is easier to specify entry and exit times in such a way as to prevent non-subcohort cases from contributing to the overall subcohort follow-up while allowing them to contribute to their own risk sets. One record is constructed for each case with entry time set to a small increment before its onset time. The increment should be much smaller than the differences between successive failure times so as to avoid overlap of risk sets. Another record is constructed for each subcohort subject (including subcohort cases) with exit time set to a small increment before its onset or censoring time. This ensures that

subcohort cases contribute to all risk sets before their time of onset. Langholz and Jiao<sup>18</sup> described SAS commands to set up the data for the exact PsL. They also described fitting the proportional hazards model and computing asymptotic and robust variance estimators in SAS using the same data setup, although (as with the Therneau and Li<sup>19</sup> SAS variance estimate) the asymptotic variance estimate requires an additional calculation by hand.

The exact PsL with asymptotic variance estimate is also available in the Epicure survival analysis module 'Peanuts'.<sup>6</sup> In Epicure, the user only need specify, in addition to the usual PL variables (entry time, exit time and failure indicator), an indicator of which cases are outside the subcohort (CSCHRT) and the subcohort sampling fraction; thus, substantially less setup is required than with SAS or S-plus. Furthermore, Epicure allows fitting general risk models, including absolute rates; therefore, it is our method of choice for unstratified case-cohort analyses.

Up to this point, we have focused on case-cohort studies involving simple random selection of the subcohort. Stratification (stratified random sampling) can be used to increase efficiency ('exposure stratification') and to deal with confounding ('confounder stratification').<sup>18</sup> In the present work, we are concerned with efficiency of risk estimation when exposure is known in the entire cohort and interest lies in making inference about interactions with exposure (exposure risk modification). Thus, we focus on exposure stratification. Stratified sampling in case-cohort studies merely involves randomly selecting subcohort subjects separately within exposure strata. As with all stratified sampling scenarios, the optimal number of strata and optimal numbers of persons to sample from each stratum are issues worth consideration,<sup>29</sup> but we do not pursue them here as they are difficult to generalize.

Borgan *et al.*<sup>29</sup> defined three methods for handling stratified case-cohort data. The first of their estimators (method I) is analogous to the approximate PsL estimator in the unstratified situation in that a case outside the subcohort does not enter into the denominator, even in its own risk set. The second estimator (method II) is somewhat analogous to the exact PsL, but non-cases are weighted according to their respective stratum proportions. The third estimator (method III) weights all subjects in the subcohort according to their respective stratum proportions and further deals with cases not in the subcohort by randomly selecting one member of the subcohort in the same stratum to be removed or 'swapped'. The reason for this approach is that the exact PsL for stratified data remains score unbiased if non-subcohort cases are excluded from the denominators of their likelihood terms, but such exclusion is inefficient.<sup>29</sup> Swappers may be chosen on a risk set basis if the analysis includes time-dependent covariates; otherwise, the swapper is selected from the initial subcohort strata, and if the

randomly chosen swapper has been censored as of the time of the current risk set, swapping for the corresponding non-subcohort case is simply ignored.

The swapper method (method III) is our recommended method for stratified case-cohort data, as it is the only score-unbiased method. Unfortunately, it is not included among the procedures for stratified case-cohort data analysis in the *cch* function of the R survival package. We therefore adapted to S-plus the SAS macros of Langholz and Jiao.<sup>18</sup> The script is provided in the Supplementary Appendix A1 available at *IJE* online, but cannot be directly implemented in R due to the lack of a necessary function—`match.data.frame`—that is included with S-plus, but not with R.

Mark and Katki<sup>38</sup> stated that variance estimation with the stratified case-cohort design might be handled using their influence function variance estimator starting with a stratified proportional hazards model, but did not provide details except to suggest that one approach is to treat cases with missing information as if they occurred outside the subcohort. Additive models may also be entertained in the stratified case-cohort design,<sup>39</sup> as may frailty models for cluster data.<sup>40</sup> Many other options and proposals exist in the literature, but they are too numerous to describe in detail. Table 1 presents an overview of the evolution of conventional methods for case-cohort study design and analysis. We conclude that score-unbiased methods with asymptotic variance are preferable based on both theoretical and small-sample properties. In particular, the robust variance cannot be implemented with stratified data, and risk-set (time-dependent) weights do not seem to be necessary with PsL estimators (although they may be related to efficiency with other types of estimators<sup>22</sup>). We therefore recommend the use of the exact PsL (using its analogue, the swapper method, in the stratified case) and the asymptotic variance estimates. We restrict our attention to that approach in the remainder of this work.

### Empirical investigations

The initial design of the immunogenome and cancer case-cohort study included two subcohorts: one using simple random sampling and one using stratified random sampling with stratification on radiation dose. Five strata were defined approximately as quintiles, but with round numbers as the cutpoints (Table 2). Both subcohorts were constructed using a sampling fraction of 0.5 to reduce genotyping effort by about one-half while retaining most of the full-cohort power.<sup>3</sup> To minimize the effort involved in molecular analyses and to facilitate direct comparison between the two sampling procedures, the two subcohorts were selected with the greatest amount of overlap possible between them. This was done by modifying the simple random subcohort, randomly adding to or excluding from each stratum the

**Table 1** Overview of developments in case-cohort study design and analysis

Reference	General approach	Special characteristics
Prentice <sup>11</sup>	Initial proposal of case-cohort design based on time-to-event analysis <sup>a</sup>	Defined exact pseudolikelihood (PsL) analogous to partial likelihood (PL)
Self and Prentice <sup>17</sup>	Asymptotic distribution of approximate PsL	Asymptotic normality of pseudolikelihood estimate and consistency of Prentice's variance
Wacholder et al <sup>37</sup>	Variance estimators derived under superpopulation (sampling without replacement) model or using the bootstrap	Superpopulation variance valid under $\beta=0$ Bootstrap method non-standard due to lack of knowledge of covariate information for all cohort members
Barlow <sup>35</sup>	Robust variance estimator based on influence functions Weights based on subcohort membership within risk sets	Analogous to Lin and Wei <sup>41</sup> robust variance for the full-cohort proportional hazards (PH) model Preserves correct expectation of PsL denominator in each risk set
Lin and Ying <sup>20</sup>	Missing-data approach to estimation with proportional hazards model	Variance estimators more readily derived with large data than those of Self/Prentice or Wacholder et al.
Barlow et al <sup>15</sup>	Jackknife variance estimator computed via delta-beta residuals	SAS macros available through statlib
Therneau and Li <sup>19</sup>	Intuitive description of fitting PH model with approximate pseudolikelihood and asymptotic variance estimation	SAS and S-plus examples provided
Chen and Lo <sup>21</sup>	More efficient use of case and cohort information in the estimating functions	Derived three methods depending on extent of knowledge of case proportion in the population
Borgan et al <sup>29</sup>	PsL estimators for stratified case-cohort data	Score-unbiased 'swapper' method analogous to unstratified exact PsL
Mark and Katki <sup>38</sup>	Generalized influence function approach to variance computation	Dealing with missing data among subcohort or case subjects
Scheike and Martinussen <sup>23</sup>	Maximum likelihood via EM algorithm (including non-sampled cohort members who do not fail)	Efficiency gain related to disease intensity Variance estimated by profile likelihood
Kulich and Lin <sup>22</sup>	Doubly weighted estimation utilizing covariate information available on the full cohort	Efficiency gain greatest for continuous covariates, less for binary covariates
Nan <sup>42</sup>	Estimation via efficient scores	Relaxes some assumptions underlying the proportional hazards model
Langholz and Jiao <sup>18</sup>	Computational methods for random and stratified designs using exact PsL and asymptotic variance	SAS code provided for data construction, proportional hazards model fitting and variance estimation
Samuelsen et al <sup>28</sup>	Intuitive methods analogous to those of Therneau and Li <sup>19</sup> for the asymptotic variance	Script provided for asymptotic variance estimation using S-plus/R Compared asymptotic and robust variance estimators
Kulathinal et al <sup>27</sup>	Design options and overview of analysis methods	SAS and R command examples
Moger et al <sup>40</sup>	Case-cohort methods for cluster-based sampling	Extended pseudolikelihood approach to gamma frailty and copula failure models

<sup>a</sup>Prentice<sup>11</sup> attributes the earliest ideas similar to the case-cohort design to Kupper et al<sup>43</sup> and Miettinen<sup>32</sup>, but seems to be the first to describe the link to partial likelihood and use of continuous explanatory factors

**Table 2** Numbers of subjects in the immunogenome study cohort and case-cohort subcohorts

	Radiation dose stratum (whole-body dose, weighted milliGray)					Total
	0 ≤ 1	1 ≤ 5	5 ≤ 100	500 ≤ 1500	1500 <	
Number of cohort subjects	1433	516	1110	1005	618	4682
Number of lung cancer cases	29	12	28	27	27	123
Random subcohort						
Total subcohort subjects	655	238	502	454	277	2126
Subcohort lung cases	17	6	10	16	13	62
Non-subcohort lung cases	12	6	18	11	14	61
Stratified subcohort						
Total subcohort subjects	428	421	425	423	428	2125
Subcohort lung cases	11	10	7	16	20	64
Non-subcohort lung cases	18	2	21	11	7	59

necessary number of individuals to obtain a stratified sample.

The first study performed using the immunogenome and cancer case-cohort design, a study of lung cancer and *EGFR* gene, was originally conducted using the simple random case-cohort sample.<sup>4</sup> However, the *EGFR* gene was ultimately assessed on the full cohort (this would not generally be the case with a case-cohort study). Hence, the two designs (simple and stratified) can be directly compared and various sampling fractions (i.e. <50%) can be evaluated via Monte Carlo simulation drawing from the full cohort. We performed such simulations analysing the data using the exact PsL (or its analogue with stratified data, the swapper method) and asymptotic variance adjustment, the methods of choice noted earlier.

Relative efficiency (efficiency of the case-cohort sample relative to that of the full cohort) was calculated as the inverse ratio of standard errors. Analyses were conducted using S-plus (version 6.2 for Windows) or R (version 2.10.0 for Windows). Analyses of simple case-cohort data were made using the `cch` function in the R survival package with the default 'Prentice' option (exact PsL). Analyses of stratified case-cohort data were performed using method I or II in the R survival package `cch` function. For method III, we adapted the SAS programs described by Langholz and Jiao<sup>18</sup> in S-plus using the `coxph` function (Supplementary Appendix A1 available at *IJE* online). Asymptotic standard errors were based on Cox regression delta-beta statistics [`S` function `residuals.coxph(type="dfbeta")`]. Residuals for the non-failures were used and summed over risk sets within individuals using the `aggregate` function in S-plus. Age was used as the primary time scale in proportional hazards regression with left truncation at entry age.<sup>44</sup> Covariates were as follows: city of residence, gender, year of birth (centred at the median, 1927, and divided by 5 years), smoking frequency (number of cigarettes smoked per day divided by 10), whole-body radiation dose (wGy: weighted Gray skin

dose, using weights 10 for neutron dose and 1 for gamma dose), an indicator of *EGFR* gene CA repeat length < 38 (short-repeat genotype) and the product of radiation dose and *EGFR* CA repeat length indicator. Confidence intervals (95%) and tests were based on Wald statistics. Results obtained using SPSS and Epicure are provided along with relevant commands in the Supplementary Appendix A1 available at *IJE* online.

## Results

Table 3 shows results of analyses of the actually selected simple and stratified 50% case-cohort designs. There was little practical difference between the two designs in terms of estimated parameters and standard errors, consistent with the expectation that a large case-cohort sample should provide close to full-cohort efficiency. Thus, the original analysis,<sup>4</sup> which was based on the simple random subcohort with 50% sampling fraction, should be reasonably efficient. The exact PsL method with asymptotic variance for the random case-cohort sample was programmed in S-plus and gave results identical (to three decimal places) to those obtained with the R `cch` function (data not shown). The three methods for a stratified case-cohort sample produced nearly identical results. How the methods compare with smaller sample sizes is addressed below.

Proportional hazards imply multiplicative effects; therefore, the negative interaction suggests a submultiplicative joint effect of radiation and CA repeat length polymorphism. What this implies about biological interaction, however, cannot be assessed without examining other scales of the joint effect, such as a purely additive statistical model.<sup>45</sup> An additive ERR model for the joint effects of radiation and repeat length polymorphism (which can only be fitted using Epicure with the simple case-cohort sample; see Supplementary Appendix A1 available at *IJE* online) produced virtually the same result as the multiplicative

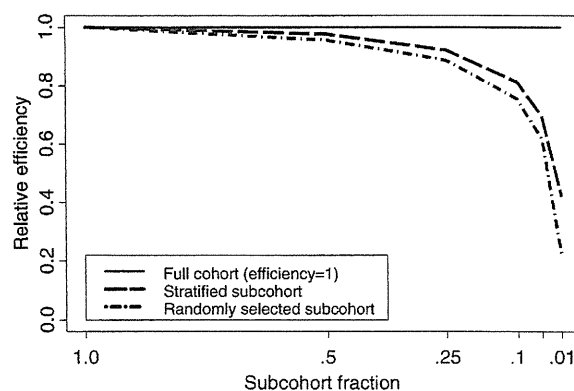
**Table 3** Proportional-hazards model fits to random and stratified subcohorts with 50% sampling fraction [coefficient (SE), Wald *P*-value]

Parameter	Case-cohort sampling method			
	Random (exact method)	Stratified		
		Method I	Method II	Method III (swapper)
Year of birth (5-year difference; centred at 1927)	0.61 (0.12)	0.62 (0.12)	0.62 (0.12)	0.62 (0.12)
<i>P</i> -value	<0.001	0.001	0.001	0.001
City (Nagasaki: Hiroshima) <sup>a</sup>	0.01 (0.21)	-0.10 (0.21)	-0.11 (0.21)	-0.10 (0.21)
<i>P</i> -value	0.96	0.65	0.69	0.63
Gender (female: male)	-0.61 (0.23)	-0.62 (0.22)	-0.63 (0.22)	-0.62 (0.22)
<i>P</i> -value	0.008	0.005	0.005	0.005
Number of cigarettes smoked per day (increment of 10)	0.45 (0.08)	0.42 (0.07)	0.42 (0.07)	0.42 (0.08)
<i>P</i> -value	<0.001	<0.001	<0.001	<0.001
Radiation dose (Gray)	0.40 (0.10)	0.37 (0.10)	0.38 (0.10)	0.37 (0.10)
<i>P</i> -value	<0.001	<0.001	<0.001	<0.001
EGFR repeat polymorphism (indicator of length <38)	0.60 (0.24)	0.52 (0.24)	0.51 (0.24)	0.52 (0.24)
<i>P</i> -value	0.012	0.029	0.034	0.027
Interaction between radiation and EGFR polymorphism	-0.40 (0.19)	-0.37 (0.18)	-0.37 (0.18)	-0.37 (0.18)
<i>P</i> -value	0.036	0.041	0.041	0.041

<sup>a</sup>Although not statistically significant, city of residence was included to avoid possible confounding of the gene-radiation interaction.

proportional hazards model, except that the radiation dose response is linear rather than log linear. The ERR for radiation was 1.00/wGy (RR=2.00 at 1 wGy) compared with the log-linear relative risk model  $RR = \exp\{0.396\} = 1.49$  at 1 wGy. With either the multiplicative relative risk or the additive ERR model, the effect of the short CA repeat length polymorphism is to approximately double the cancer risk: the ERR for the short repeat length polymorphism was 1.06 (RR=2.06) compared with the log-linear model  $RR = \exp\{0.601\} = 1.82$ . With both models, statistical interaction between radiation dose and repeat length polymorphism essentially negated the radiation dose-response main effect. Interaction on the additive ERR scale was -1.09, which effectively cancels the radiation ERR of 1.00. On the multiplicative scale, interaction parameters were close to equal but opposite in sign compared with the radiation main effect (Table 3). Thus, there is apparently no practical radiation effect among individuals with the short CA repeat length polymorphism.<sup>4</sup>

Results of Monte Carlo sampling from the full cohort with sampling fractions varying from 0.5 down to 0.01 are shown in Figure 1 and Table 4. Relative efficiency decreased with decreasing sampling fraction. Relative efficiency of the simple case-cohort sample declined more dramatically than that of the exposure-stratified sample—that is, the stratified design becomes relatively more efficient with smaller subcohort fractions (Figure 1).



**Figure 1** Relative efficiency of case-cohort analysis of gene-radiation interaction as a function of subcohort fraction: simple vs stratified sampling

To compare by simulation the three methods for stratified data, we ran the S-plus swapper function in R by importing without modification the S-plus match.data.frame function. All three methods produced similar results with subcohort fractions down to ~10%, but with sampling fraction of 5% or 1% method I had greater bias and larger standard error than methods II and III (Table 4). Overall, method III (the swapper method) performed the best, having the least bias and smallest standard errors, although differences between methods II and III were trivial.

**Table 4** Mean and average standard error (SE) of estimates of the radiation–gene interaction with stratified case–cohort data using simulation from the full cohort<sup>a</sup>

Subcohort fraction <sup>b</sup> (%)	Method I Mean (SE)	Method II Mean (SE)	Method III (swapper) Mean (SE)
50	−0.43 (0.18)	−0.43 (0.18)	−0.43 (0.18)
25	−0.43 (0.19)	−0.43 (0.19)	−0.43 (0.19)
10	−0.45 (0.22)	−0.45 (0.21)	−0.45 (0.21)
5	−0.47 (0.26)	−0.46 (0.25)	−0.45 (0.25)
1	−0.72 (0.64)	−0.48 (0.44)	−0.43 (0.41)

<sup>a</sup>Results are based on 2000 simulations at each subcohort fraction.

<sup>b</sup>For comparison, the estimate and standard error with the full cohort data were −0.42 and 0.17, respectively.

## Discussion and conclusions

The case–cohort design is not used as frequently as it might be given the advantages it offers epidemiologists, perhaps because the design is misunderstood given its synthetic nature. It is hoped that the present review will serve to heighten investigators' awareness of the definition and advantages of the case–cohort design and provide guidance to facilitate its implementation and analysis. Issues that can be resolved in a fairly general manner include: whether to stratify, whether to perform exact or approximate PsL analysis, whether to use asymptotic or robust variance estimates and whether/how to weight the observations in the analysis.

Our results suggest that selecting the subcohort using stratified, rather than simple, random sampling can be beneficial in terms of statistical efficiency, especially for studying interaction, when the sampling fraction is substantially less than one-half. Although relatively small subcohort fractions are typical, subcohorts as large as 50% (such as in the immunogenome and cancer study described here) are not beyond imagination. Even such a large study reduces by almost half the amount of covariate collection effort yet retains efficiency close to that of the full cohort. If an investigator is able to conduct a study with large subcohort fraction, there may be little benefit from the stratified sampling. However, there is little additional effort involved in selecting and analysing a stratified subcohort; therefore, if there are factors for which it makes sense to stratify, we recommend doing so, regardless of the sampling fraction.

Simulation studies conducted by others have demonstrated that the exact PsL is preferable to the approximate PsL in small samples, even though the two are asymptotically equivalent. Because there is little difference in the data preparation required for analysis, we recommend that the exact PsL be used. Borgan *et al.*<sup>29</sup> compared methods I, II and III for stratified case–cohort data, using simulations with a

10% subcohort fraction and did not note any major differences among them. However, from our Table 4, it is apparent that method I can incur more small-sample bias and lose efficiency more rapidly than the other methods with decreasing subcohort sampling fractions <10%. Method III produced the best results overall. These new findings support our conclusion that the swapper method (method III) is the method of choice for stratified case–cohort data.

Other authors have noted that the robust variance estimate can perform poorly in small samples, and the robust variance estimate cannot be used with stratified data. We therefore recommend use of the asymptotic variance estimate with both the simple (unstratified) and stratified designs.

Because risk set weighting is related to the counting-process theory underlying PL analysis, it is expected to perform well in theory. The results of Onland-Moret *et al.*<sup>26</sup> did not resolve its performance vis-à-vis other methods because they used the simplified version with constant weight.<sup>26</sup> In the simulations of Borgan *et al.*,<sup>29</sup> time-dependent weights only slightly improved efficiency of the stratified design. We conclude that gains from time-dependent weighting probably are not sufficient to justify the additional data construction effort with unstratified data. With stratified data, the data construction facilitates time-dependent weighting, but a variance estimator is not yet available. It has been noted that choice of appropriate weights might be affected when specimen storage duration or analytical batch effects are a concern (i.e. for cases outside the subcohort in a prospectively conducted case–cohort study<sup>46</sup>); this would not be a problem with retrospectively conducted case–cohort studies. In fact, although we are not aware that it has been proposed before, it should be possible to perform risk set subsampling of nested control subjects from the subcohort, if the timing of case ascertainment could lead to biases due to storage or batch effects.

Missing data are an important consideration in observational studies. Regardless of design, measurements made on biological specimens can be missing for a number of reasons. Mark and Katki<sup>38</sup> suggested that, as long as missingness is not related to covariate values, case–cohort analyses without time-varying weights should remain valid if cases with missing covariate information are simply ignored.

One advantage of the nested case–control and case–cohort designs as compared with the traditional case–control study design is their ability to estimate absolute risk difference (as opposed to relative risk) and survival probability. However, when there is sampling of cases or missing case data, estimators of those quantities require modification to avoid bias.<sup>47</sup> Another advantage of the cohort-based sampling designs is that the reference group can be used to study population characteristics, such as allele frequencies or Hardy–Weinberg equilibrium; the control subjects



for a traditional case-control study are not population based, being left over after cases have been removed from the population.

The ability to study multiple outcomes individually without having to obtain separate sets of comparison subjects is often cited as one of the benefits of the case-cohort design, but this aspect is seldom illustrated or examined in the literature. In particular, the use of the same comparison subjects induces correlations among results for multiple outcomes. Sørensen and Andersen<sup>48</sup> investigated this via asymptotic theory and simulations and, using competing risks and martingale theory, derived a consistent estimate of the asymptotic covariance matrix. They found that correlations increase with smaller subcohort sampling fractions or, equivalently, with larger cohorts and a fixed subcohort size. For example, with a cohort of size 500 and subcohort sampling fraction of 10%, correlations were substantially larger than 0.5 for the scenarios they studied. With a subcohort sampling fraction of ~50%, they estimated correlations of 0.1–0.2 for their scenarios. In the immunogenome and cancer study described here, where the subcohort sampling fraction was 50%, correlations between effect estimates for multiple cancer outcomes should therefore not be a problem (studies involving cancers at sites other than the lung are ongoing). However, such correlations could be a problem with more typical designs (25 or 10%), which may be used if more expensive procedures, such as larger numbers of loci or whole-genome sequencing, are used.

Another interesting aspect of the case-cohort design is that outcomes occurring more frequently have smaller effective control:case ratios given that the number of comparison subjects is fixed. Kulathinal *et al.*<sup>27</sup> noted that choosing a subcohort size suitable for more common outcomes would provide flexibility for including additional, less common outcomes after the fact of subcohort selection. With common outcomes, sampling of cases may be preferable to using all cases in terms of reducing effort without seriously affecting statistical efficiency.<sup>49</sup>

Even though ascertainment of some covariates, such as genotypes, may be conducted only on the case-cohort sample, efficiency may be lost because analysis based on the case-cohort sample alone wastes other information available in the entire cohort. For example, if modelling adjustment is made for important factors or potential confounders that are known in all cohort subjects, that adjustment is inefficient if restricted only to the case-cohort sample. Further efficiency gains may be obtained by considering the case-cohort sample as the second phase in a two-phase study design by weighting the analysis of the case-cohort sample through post-stratification or calibration.<sup>24,50</sup> A similar strategy can be applied to nested case-control studies,<sup>51</sup> and more general multi-phase methods can be considered.<sup>52</sup>

Alternatively, efficiency might be improved by considering the case-cohort study as a full cohort with missing data using, for example, inverse probability of sampling weighted methods.<sup>53</sup> We do not pursue those methods in this work, as they are beyond the scope of traditional case-cohort analysis, and we believe that more comparative work needs to be done to fully understand the relative merits of these new approaches vis-à-vis the conventional design and analysis of case-cohort studies. We are currently exploring the two-phase approach in comparison with the conventional methods recommended in the present review.

The typical case-cohort analysis is based on a restricted model for the hazard rate (incidence or mortality rate) as a log-linear function of explanatory covariates. Kong and Cai<sup>54</sup> described how to fit accelerated failure time models to case-cohort data. The accelerated-failure-time model does not require the assumption of proportional hazards, although that assumption can also be relaxed in standard proportional hazards models by including interactions of covariates with the basic time variable. Apart from Langholz and Jiao<sup>18</sup> describing how to estimate absolute risks, a frustrating aspect of many current software implementations of case-cohort analysis methods is the restriction to standard proportional hazards models. In radiation risk assessment, more general models are routinely used.<sup>1</sup> A Bayesian full-likelihood approach<sup>55</sup> may allow more general risk models. At present, Epicure only allows fitting general risk models with unstratified case-cohort data, but a procedure to accommodate stratified data should become available in the near future.

Because we have focused on gene-environment interaction, a few words about such tests are appropriate. Failure to adequately model the main effect of an environmental risk factor can lead to inflated type I errors (false positives) in tests of gene-environment interaction involving that factor.<sup>56</sup> This result emphasizes the need for use of general risk models, such as the ERR model used in our example. In our opinion, models for main effects of genotype have often not been given due consideration. Perhaps this is due to the popularity of the Cochran-Armitage trend test, which is applied under the assumption of a linear (co-dominant) genomic model. It is unfortunate that one of the major software packages for genomic analysis (R haplo.stats<sup>57</sup>) does not include the arbitrary two-parameter genomic model which we deem an important starting point (e.g. for testing homogeneity with biallelic loci). The case-only design can be powerful when gene and environmental factor are independent, and hybrid methods address the fact that such independence often cannot be verified.<sup>58</sup> However, the case-only design is limited to a multiplicative test of statistical interaction, whose use has been over-emphasized.<sup>59</sup> Furthermore, the subcohort in a



case-cohort design is suitable for cross-sectional analyses,<sup>46</sup> such as assessing allele frequencies and Hardy-Weinberg equilibrium (an important but not essential assumption for association studies<sup>60,61</sup>). Finally, we note that if the only impact of a gene is to modify the effect of environmental exposure, a main effect of genotype/haplotype may not be required in the model.<sup>62</sup> Alternatively, it is conceivable that a gene evidencing both a direct effect on outcome and modification of risk of an environmental factor might not necessarily exert both effects according to the same genomic model if they involve separate mechanisms. A number of further issues are discussed in detail in two recent publications.<sup>63,64</sup>

In conclusion, for studying interactions with an exposure variable, we recommend that case-cohort studies be conducted based on exposure-stratified random sampling when exposure can be ascertained in the entire cohort. When stratification is not necessary, the exact PsL with asymptotic variance estimates may be used and the most flexible implementation is to be found in the Epicure software. With stratified case-cohort data, we recommend the score-unbiased method III (the swapper method) of Borgan *et al.*,<sup>29</sup> with variances estimated using the asymptotic method. Stratified data may be analysed using the SAS macros of Langholz and Jiao<sup>18</sup> or the S-plus function provided here (Supplementary Appendix A1 available at *IJE* online). Further work is needed to extend existing software for stratified case-cohort data to allow fitting of risk models that are more general than the standard proportional hazards model.

#### KEY MESSAGE

- We summarize the literature on case-cohort design and analysis and make recommendations regarding its implementation for studies of gene-environment interaction. Stratified, rather than simple, random sampling should be preferred when exposure information exists for the entire cohort. The current state of relevant software is reviewed, and needs are identified for further development of software tools that allow fitting general, flexible risk models. Given that the subcohort can be used for assessment of population parameters, such as allele frequencies and Hardy-Weinberg equilibrium, the case-cohort design is well suited to molecular epidemiological research.

## References

- <sup>1</sup> Preston DL, Ron E, Tokuoka S *et al.* Solid cancer incidence in atomic bomb survivors: 1958-1998. *Radiat Res* 2007;**168**:1-64.
- <sup>2</sup> Yamada M, Wong FL, Fujiwara S, Akahoshi M, Suzuki G. Noncancer disease incidence in atomic bomb survivors, 1958-1998. *Radiat Res* 2004;**161**:622-32.
- <sup>3</sup> Hayashi T, Kusunoki Y, Kyoizumi S *et al.* Relationship between cancer development and genetic polymorphisms among A-bomb survivors, focusing on immune-related genes. Research Protocol 4-04. Hiroshima, Japan: Radiation Effects Research Foundation, 2004; [http://www.rerf.jp/programs/rparchiv\\_c/rp04-04.htm](http://www.rerf.jp/programs/rparchiv_c/rp04-04.htm) (5 July 2012, date last accessed).
- <sup>4</sup> Yoshida K, Nakachi K, Imai K *et al.* Lung cancer susceptibility among atomic bomb survivors in relation to CA repeat number polymorphism of epidermal growth factor receptor gene and radiation dose. *Carcinogenesis* 2009;**30**: 2037-41.
- <sup>5</sup> Cox DR. Regression models and life tables. *J R Stat Soc Series B* 1972;**34**:187-202.
- <sup>6</sup> Preston DL, Lubin JH, Pierce DA, McConney ME. *Epicure Users Guide*. Seattle: Hirosoft International Corporation, 1993.

## Supplementary Data

Supplementary Data are available at *IJE* online.

## Funding

This work was supported by the Radiation Effects Research Foundation (RERF), Hiroshima and Nagasaki, Japan, a private, non-profit foundation funded by the Japanese Ministry of Health, Labour and Welfare (MHLW) and the U.S. Department of Energy (DOE), the latter in part through the National Academy of Sciences [RERF Research Protocol RP 4-04]; the Ministry of Education, Culture, Sports, Science and Technology of Japan [Grant-in-Aid for Scientific Research (B) 21390199]; and the Ministry of Health, Labour and Welfare of Japan [Grant-in-Aid for Cancer Research 22090501].

## Acknowledgements

The authors express sincere gratitude to Professor Bryan Langholz for assistance with the swapper method for stratified case-cohort data and to Dr Robert Abbott for guidance on the logistic regression approximation to Cox regression (both have provided their consent to be so acknowledged).

**Conflict of interest:** D.L.P. was the principal developer of the Epicure software, which is now owned by Risk Sciences International (RSI). Although he no longer owns the software, there is a possibility he might receive royalties on some future sales of the software.

- <sup>7</sup> Breslow NE, Day NE. *Statistical Methods in Cancer Research, Volume II – The Design and Analysis of Cohort Studies*. Lyon: International Agency for Research on Cancer, 1987.
- <sup>8</sup> Langholz B, Borgan Ø. Counter-matching: a stratified nested case-control sampling method. *Biometrika* 1995; **82**:69–79.
- <sup>9</sup> Cologne J, Langholz B. Selecting controls for assessing interaction in nested case-control studies. *J Epidemiol* 2003; **13**:193–202.
- <sup>10</sup> Cologne JB, Sharp GB, Neriishi K, Verkasalo PK, Land CE, Nakachi K. Improving the efficiency of nested case-control studies of interaction by selecting controls using counter matching on exposure. *Int J Epidemiol* 2004; **33**:485–92.
- <sup>11</sup> Prentice RL. A case-cohort design for epidemiologic cohort studies and disease prevention trials. *Biometrika* 1986; **73**:1–11.
- <sup>12</sup> Langholz B, Thomas DC. Nested case-control and case-cohort methods of sampling from a cohort: a critical comparison. *Am J Epidemiol* 1990; **131**:169–76.
- <sup>13</sup> Wacholder S. Practical considerations in choosing between the case-cohort and nested case-control designs. *Epidemiology* 1991; **2**:155–58.
- <sup>14</sup> Volovics A, van den Brandt PA. Methods for the analysis of case-cohort studies. *Biom J* 1997; **39**:195–214.
- <sup>15</sup> Barlow WE, Ichikawa L, Rosner D, Izumi S. Analysis of case-cohort designs. *J Clin Epidemiol* 1999; **52**:1165–72.
- <sup>16</sup> Zeng D, Lin DY, Avery CL, North KE, Bray MS. Efficient semiparametric estimation of haplotype-disease associations in case-cohort and nested case-control studies. *Biostatistics* 2006; **7**:486–502.
- <sup>17</sup> Self SG, Prentice RL. Asymptotic distribution theory and efficiency results for case-cohort studies. *Ann Stat* 1988; **16**:64–81.
- <sup>18</sup> Langholz B, Jiao J. Computational methods for case-cohort studies. *Comp Statist Data Anal* 2007; **51**:3737–48.
- <sup>19</sup> Therneau TM, Li H. Computing the Cox model for case cohort designs. *Lifetime Data Anal* 1999; **5**:99–112.
- <sup>20</sup> Lin DY, Ying Z. Cox regression with incomplete covariate measurements. *J Am Stat Assoc* 1993; **88**:1341–49.
- <sup>21</sup> Chen K, Lo SH. Case-cohort and case-control analysis with Cox's model. *Biometrika* 1999; **86**:755–64.
- <sup>22</sup> Kulich M, Lin DY. Improving the efficiency of relative-risk estimation in case-cohort studies. *J Am Stat Assoc* 2004; **99**:832–44.
- <sup>23</sup> Scheike TH, Martinussen T. Maximum likelihood estimation for Cox's regression model under case-cohort sampling. *Scand J Stat* 2004; **31**:283–93.
- <sup>24</sup> Breslow NE, Lumley T, Ballantyne CM, Chambless LE, Kulich M. Using the whole cohort in the analysis of case-cohort data. *Am J Epidemiol* 2009; **169**:1398–405.
- <sup>25</sup> Chen K. Generalized case-cohort sampling. *J R Stat Soc Series B* 2001; **63**:791–809.
- <sup>26</sup> Onland-Moret NC, van der ADL, van der Schouw YT *et al*. Analysis of case-cohort data: a comparison of different methods. *J Clin Epidemiol* 2007; **60**:350–55.
- <sup>27</sup> Kulathinal S, Karvanen J, Saarela O, Kuulasmaa K. Case-cohort design in practice – experiences from the MORGAM Project. *Epidemiol Perspect Innov* 2007; **4**:15.
- <sup>28</sup> Samuelsen SO, Ånestad H, Skrandal A. Stratified case-cohort analysis of general cohort sampling designs. *Scand J Stat* 2007; **34**:103–19.
- <sup>29</sup> Borgan Ø, Langholz B, Samuelsen SO, Goldstein L, Pogoda J. Exposure stratified case-cohort designs. *Lifetime Data Anal* 2000; **6**:39–58.
- <sup>30</sup> Breslow N. "cch: Fits proportional hazards regression model to case-cohort data". Function in Therneau T: "Package 'survival': Survival analysis, including penalised likelihood". *R package*; version 2.36-12, <http://cran.r-project.org/web/packages/survival/survival.pdf>, February 15 2012 (accessed 15 March 2012).
- <sup>31</sup> Therneau T. Package 'survival': Survival analysis, including penalised likelihood. *R package*; version 2.36-12, February 15, 2012, <http://cran.r-project.org/web/packages/survival/survival.pdf> (accessed 15 March 2012).
- <sup>32</sup> Miettinen S. Design options in epidemiologic research: an update. *Scand J Work Environ Health* 1982; **8**(Suppl 1): 7–14.
- <sup>33</sup> Greenland S. Adjustment of risk ratios in case-base studies (hybrid epidemiologic designs). *Stat Med* 1986; **5**: 579–84.
- <sup>34</sup> Sato T. Estimation of a common risk ratio in stratified case-cohort studies. *Stat Med* 1992; **11**:1599–605.
- <sup>35</sup> Barlow WE. Robust variance estimation for the case-cohort design. *Biometrics* 1994; **50**:1064–72.
- <sup>36</sup> Therneau TM, Grambsch PM. *Modeling Survival Data: Extending the Cox Model*. New York: Springer-Verlag, 2000.
- <sup>37</sup> Wacholder S, Gail MH, Pee D, Brookmeyer R. Alternative variance and efficiency calculations for the case-cohort design. *Biometrika* 1989; **76**:117–23.
- <sup>38</sup> Mark SD, Katki H. Influence function based variance estimation and missing data issues in case-cohort studies. *Lifetime Data Anal* 2001; **7**:331–44.
- <sup>39</sup> Kulich M, Lin DY. Additive hazards regression for case-cohort studies. *Biometrika* 2000; **87**:73–87.
- <sup>40</sup> Moger TA, Pawitan Y, Borgan Ø. Case-cohort methods for survival data on families from routine registries. *Stat Med* 2008; **27**:1062–74.
- <sup>41</sup> Lin DY, Wei LJ. The robust inference for the Cox proportional hazards model. *J Am Stat Assoc* 1989; **84**:1074–78.
- <sup>42</sup> Nan B. Efficient estimation for case-cohort studies. *Can J Stat* 2004; **32**:403–19.
- <sup>43</sup> Kupper LL, McMichael AJ, Spirtas R. A hybrid epidemiologic study design useful in estimating relative risk. *J Am Stat Assoc* 1975; **70**:524–28.
- <sup>44</sup> Cologne J, Hsu W-L, Abbott RD *et al*. Proportional hazards regression in epidemiologic follow-up studies: an intuitive consideration of primary time scale. *Epidemiology* 2012; **23**:565–73.
- <sup>45</sup> Rothman KJ, Greenland S, Lash TL. *Modern Epidemiology*. 3rd edn. Philadelphia: Lippincott Williams & Wilkins, 2008.
- <sup>46</sup> Rundle AG, Vineis P, Ahsan H. Design options for molecular epidemiology research within cohort studies. *Cancer Epidemiol Biomarkers Prev* 2005; **14**:1899–907.
- <sup>47</sup> Mark SD, Katki HA. Specifying and implementing non-parametric and semiparametric survival estimators in two-stage (nested) cohort studies with missing case data. *J Am Stat Assoc* 2006; **101**:460–71.
- <sup>48</sup> Sørensen P, Andersen PK. Competing risks analysis of the case-cohort design. *Biometrika* 2000; **87**:49–59.
- <sup>49</sup> Cai J, Zeng D. Power calculation for case-cohort studies with nonrare events. *Biometrics* 2007; **63**:1288–95.
- <sup>50</sup> Breslow NE, Lumley T, Ballantyne CM, Chambless LE, Kulich M. Improved Horvitz-Thompson estimation of model parameters from two-phase stratified samples: applications in epidemiology. *Stat Biosci* 2009; **1**:32–49.

- <sup>51</sup> Breslow NE, Chatterjee N. Design and analysis of two-phase studies with binary outcome applied to Wilms tumour prognosis. *Appl Stat* 1999;**48**:457–68.
- <sup>52</sup> Whittemore AS. Multistage sampling designs and estimating equations. *J R Stat Soc Series B* 1997;**59**:589–602.
- <sup>53</sup> Robins JM, Rotnitzky A, Zhao LP. Estimation of regression coefficients when some regressors are not always observed. *J Am Stat Assoc* 1994;**89**:846–66.
- <sup>54</sup> Kong L, Cai J. Case-cohort analysis with accelerated failure time model. *Biometrics* 2009;**65**:135–42.
- <sup>55</sup> Kulathinal S, Arjas E. Bayesian inference from case-cohort data with multiple end-points. *Scand J Stat* 2006;**33**:25–36.
- <sup>56</sup> Cornelis MC, Tchetgen EJT, Liang L *et al*. Gene-environment interactions in genome-wide association studies: a comparative study of tests applied to empirical studies of type 2 diabetes. *Am J Epidemiol* 2012;**175**:191–202.
- <sup>57</sup> Sinnwell JP, Schaid DJ. Package haplo.stats': Statistical analysis of haplotypes with traits and covariates when linkage phase is ambiguous. *R package*; version 1.5.5, February 25, 2012. <http://cran.r-project.org/web/packages/haplo.stats/haplo.stats.pdf> (accessed 15 March 2012).
- <sup>58</sup> Mukherjee B, Ahn J, Gruber SB, Chatterjee N. Testing gene-environment interaction in large-scale case-control association studies: possible choices and comparisons. *Am J Epidemiol* 2012;**175**:177–90.
- <sup>59</sup> Mukherjee B, Ahn J, Gruber SB, Chatterjee N. Response to invited commentary "GE-whiz! Ratcheting up gene-environment studies". *Am J Epidemiol* 2012;**175**:208–9.
- <sup>60</sup> Trikalinos TA, Salanti G, Khoury MJ, Ioannidis JPA. Impact of violations and deviations in Hardy-Weinberg equilibrium on postulated gene-disease associations. *Am J Epidemiol* 2006;**163**:300–9.
- <sup>61</sup> Lin DY, Zeng D, Millikan R. Maximum likelihood estimation of haplotype effects and haplotype-environment interactions in association studies. *Genet Epidemiol* 2005;**29**:299–312.
- <sup>62</sup> Thomas DC, Lewinger JP, Murcray CE, Gauderman WJ. Invited commentary: GE-whiz! Ratcheting gene-environment studies up to the whole genome and the whole exposome. *Am J Epidemiol* 2012;**175**:203–7.
- <sup>63</sup> Khoury MJ, Wacholder S. Invited commentary: from genome-wide association studies to gene-environment-wide interaction studies—challenges and opportunities. *Am J Epidemiol* 2009;**169**:227–30.
- <sup>64</sup> Thomas D. Gene-environment-wide association studies: emerging approaches. *Nat Rev Genet* 2010;**11**:259–72.
- <sup>65</sup> Samuelsen SO. *Case-cohort studies: Pre-course 13. Norwegian Epidemiology Conference Tromsø 23–24*. Presentation slides, November, 2005 <http://folk.uio.no/osamuels/cascohort4.pdf> (accessed 15 March 2012).

## Compromised hematopoiesis and increased DNA damage following non-lethal ionizing radiation of a human hematopoietic system reconstituted in immunodeficient mice

Changshan Wang<sup>1</sup>, Shunsuke Nakamura<sup>1</sup>, Motohiko Oshima<sup>1</sup>, Makiko Mochizuki-Kashio<sup>1</sup>, Yaeko Nakajima-Takagi<sup>1</sup>, Mitsujiro Osawa<sup>1</sup>, Yoichiro Kusunoki<sup>2</sup>, Seishi Kyoizumi<sup>2</sup>, Kazue Imai<sup>2</sup>, Kei Nakachi<sup>2</sup> & Atsushi Iwama<sup>1</sup>

<sup>1</sup>Department of Cellular and Molecular Medicine, Graduate School of Medicine, Chiba University, Chiba, and <sup>2</sup>Department of Radiobiology/Molecular Epidemiology, Radiation Effects Research Foundation, Hiroshima, Japan

### Abstract

**Purpose:** Precise understanding of radiation effects is critical to development of new modalities for the prevention and treatment of radiation-induced damage. In this study, we evaluated the effects of non-lethal doses of X-ray irradiation on human hematopoietic stem and progenitor cells (HSPC) reconstituted in NOD/Shi-*scid*, IL2R $\gamma^{\text{null}}$  (NOG) immunodeficient mice.

**Materials and methods:** We transplanted cord blood CD34<sup>+</sup> HSPC into NOG mice irradiated with 2.0 Gy via tail veins. At the 12th week after transplantation, the NOG mice were irradiated with 0, 0.5, 1.0, 2.0, or 4.0 Gy, and the radiation effects on human HSPC *in vivo* were evaluated.

**Results:** Although a majority of the mice irradiated with 2.0 Gy or more died in 12 weeks after irradiation, the mice that were exposed to 0.5 or 1.0 Gy of irradiation survived and were subjected to analysis. The chimerism of human CD45<sup>+</sup> hematopoietic cells in peripheral blood and bone marrow (BM) of the recipient mice was reduced in an X-ray dose-dependent manner after irradiation. Percentages of human CD34<sup>+</sup> HSPC as well as human CD34<sup>+</sup>CD38<sup>-</sup> HSC in BM similarly declined. CD34<sup>+</sup>CD38<sup>-</sup> HSC purified from the humanized mice at the 12th week after irradiation showed significantly increased numbers of phosphorylated H2AX ( $\gamma$ H2AX) foci, a marker of DNA breaks, in an X-ray dose-dependent manner. Expression of p16<sup>INK4A</sup>, a hallmark of aging of HSC, was also detected only in HSPC from irradiated mice.

**Conclusions:** With further refinement, the humanized mouse model might be effectively used to study the biological effects of non-lethal radiation *in vivo*.

**Keywords:** Human HSPC; radiation effect; DNA damage; NOG immunodeficient mice

### Introduction

Effects of ionizing radiation on human hematopoietic stem and progenitor cells (HSPC) have not been well characterized *in vivo* due to difficulties in bioassays using bone

marrow (BM) samples. In mice, sublethal doses of total-body irradiation (TBI) have been reported to induce hematopoietic stem cells (HSC) senescence accompanied with impaired self-renewal capacity of HSC and upregulated biomarkers of cellular senescence, including p16<sup>INK4a</sup> (Wang et al. 2006). Furthermore, radiation-induced genomic instability in normal hematopoietic cells has been demonstrated to persist for a prolonged period *in vivo* (Watson et al. 2001). We have also previously reported that genomic effects of irradiation on the murine hematopoietic system persisted *in vivo* for long periods in terms of an increased number of micronucleated reticulocytes in peripheral blood (PB) as an indicator of genomic instability (Hamasaki et al. 2007). In contrast, radiation effects on human HSPC have been evaluated mainly *in vitro* (Rübe et al. 2011). Although human hematopoietic progenitor cells have been evaluated *in vivo* using the human fetus bone transplanted into *scid/scid* (SCID) mice (Kyoizumi et al. 1994), more precise understanding of radiation effects *in vivo* on HSPC is critical to develop new modalities for the prevention and treatment of radiation-induced damage.

In this study, we demonstrate that 0.5 and 1.0 Gy TBI affected the function of human HSC reconstituted in NOG immunodeficient mice. We propose that the humanized mouse model is useful for the evaluation of the radiation effects on human hematopoiesis, particularly in the dose region where long-term effects on the hematopoietic system as well as other organs have been observed among atomic-bomb survivors (Nakachi et al. 2008).

### Materials and methods

#### Transplantation

NOG immunodeficient mice were supplied by the Central Institute for Experimental Animals (Kawasaki, Kanagawa, Japan) and maintained in the Animal Research Facility of the Graduate School of Medicine, Chiba University, in

Correspondence: Prof. Atsushi Iwama, MD, Department of Cellular and Molecular Medicine, Graduate School of Medicine, Chiba University, 1-8-1 Inohana, Chuo-ku, Chiba 260-8670, Japan. Tel: +81 43 226 2187. Fax: +81 43 226 2191. E-mail: aiwama@faculty.chiba-u.jp

(Received 21 June 2012; revised 27 July 2012; accepted 2 September 2012)