# Estimation of varying coefficients for a growth curve model

**Kenichi Satoh[1],* and Hirokazu Yanagihara[2],***

[1] *Department of Environmetrics and Biometrics*

*Research Institute for Radiation Biology and Medicine, Hiroshima University*

*1-2-3 Kasumi, Minami-ku, Hiroshima 734-8553, JAPAN*

[2] *Department of Mathematics, Graduate School of Science, Hiroshima University*

*1-3-1 Kagamiyama, Higashi-Hiroshima 739-8626, JAPAN*

*email:* ksatoh@hiroshima-u.ac.jp

**email:* yanagi@math.sci.hiroshima-u.ac.jp

SUMMARY. In this paper, a new approach to a growth curve model is developed which uses time-varying coefficients. Since the mean structure of the growth curve model has many unknown parameters depending on both covariates and time trend designs, it can be difficult to understand and interpret. Using varying coefficient functions, the effects of covariates can be evaluated and visualized more easily. The functional confidence intervals are derived theoretically and a procedure is proposed to test whether the effects of covariates are significant.

KEY WORDS: Confidence interval; Growth curve model; Hypothesis testing; Longitudinal data; Repeated measurements; Varying coefficient.

## 1. Introduction

Let $y_i(t)$ $(i = 1, \ldots, n)$ be an observation of the $i$th subject at time $t$, and $\varepsilon_i(t)$ $(i = 1, \ldots, n)$ be a error term with $E[\varepsilon_i(t)] = 0$. We suppose that observations for different subjects are independent, and we let $\boldsymbol{x}(t)$ be a $q$-dimensional known design vector, where $n$ is the sample size. We write $\boldsymbol{y}_i = (y_i(t_{i1}), \ldots, y_i(t_{ip_i}))'$, $\boldsymbol{\varepsilon}_i = (\varepsilon_i(t_{i1}), \ldots, \varepsilon_i(t_{ip_i}))'$ and $\boldsymbol{X}_i = (\boldsymbol{x}(t_{i1}), \ldots, \boldsymbol{x}(t_{ip_i}))'$. Suppose that $\boldsymbol{\varepsilon}_1, \ldots, \boldsymbol{\varepsilon}_n$ are independent and let the covariance matrix of $\boldsymbol{\varepsilon}_i$ be denoted by $\mathrm{Cov}[\boldsymbol{\varepsilon}_i] = \boldsymbol{\Sigma}_i$. Then, the traditional growth curve model (GCM) is defined by

$$\boldsymbol{y}_i = X_i \boldsymbol{\alpha}_i + \boldsymbol{\varepsilon}_i, \quad (i = 1, \ldots, n), \tag{1}$$

where $\boldsymbol{\alpha}_i$ is a $q$-dimensional unknown parameter vector. It is well known that the GCM is a useful model for describing how growth varies over time. For example, if we fit a polynomial curve of degree $q - 1$ to the trend over time, the design vector is taken to be $\boldsymbol{x}(t) = (1, t, \ldots, t^{q-1})'$. If a more complicated trend is required, nonparametric curves can be applied, taking basis functions to be $B$-splines, or utilizing a Gaussian kernel. Note that when $p_1 = \cdots = p_n = p$ and $t_{1j} = \cdots = t_{nj} = t_j$ for $j = 1, \ldots, p$, we have a so-called balanced design; otherwise we have an unbalanced design.

Historically, in the GCM, explanations of the variation between individual subjects have been attempted in terms of known covariates $a_{ij}$ $(i = 1, \ldots, n; j = 1, \ldots, k)$. This means that $\boldsymbol{\alpha}_i$ is regressed by $\boldsymbol{a}_i = (a_{i1}, \ldots, a_{ik})'$, i.e., $\boldsymbol{\alpha}_i = \boldsymbol{\Theta}' \boldsymbol{a}_i$, where $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_k)'$ is a $k \times q$ unknown parameter matrix (see Potthoff and Roy (1964), Laird and Ware (1982), Vonesh and Carter (1987)). For example, if we express a time trend using a polynomial curve, the coefficients of the polynomial function are regressed by $\boldsymbol{a}_i$. However, it can

be difficult to understand the effects of each covariate, because $\boldsymbol{\Theta}$ has many unknown parameters, the interpretations of which depend on both design vectors $\boldsymbol{x}(t)$ and $\boldsymbol{a}_i$, particularly when the design is based on smooth curves, for example, basis functions for nonparametric regression. Furthermore, we usually consider confidence intervals and hypothesis testing involving only single elements of $\boldsymbol{\Theta}$. Since an element of $\boldsymbol{\Theta}$ is essentially a coefficient of a polynomial expressing a time trend, interval estimation and testing of a single element are meaningless in most cases. We are not concerned with the effect of a covariate on the coefficient of a polynomial function; however we are interested in the effects of covariates on the amount of growth. Therefore, we consider a new approach to the growth curve model in order to determine the effect of each covariate.

By using $\boldsymbol{\Theta}$, the GCM in (1) can be rewritten as

$$y_i(t) = \boldsymbol{a}_i' \boldsymbol{\Theta} \boldsymbol{x}(t) + \varepsilon_i(t), \quad (t = t_{i1}, \ldots, t_{ip_i}; i = 1, \ldots, n). \tag{2}$$

In the ordinary GCM, $\boldsymbol{\Theta}$ and $\boldsymbol{a}_i$ are described together as $\boldsymbol{\alpha}_i$. On the other hand, we discuss here the possibility of pairing $\boldsymbol{\Theta}$ and $\boldsymbol{x}(t)$ as a function of $t$. A varying coefficients model (see West *et al.* (1985) and Hastie and Tibshirani (1993)) is known to be an important model for analyzing longitudinal data. By defining $\boldsymbol{\beta}(t) = (\beta_1(t), \ldots, \beta_k(t))' = \boldsymbol{\Theta} \boldsymbol{x}(t)$, we can regard the growth curve model in (2) as a varying coefficients model. The varying coefficient functions are semiparametric curves, the shapes of which are controlled by $\boldsymbol{x}(t)$. The mean structure of the observations can be rewritten as

$$E[y_i(t)] = \boldsymbol{a}_i' \boldsymbol{\Theta} \boldsymbol{x}(t) = \boldsymbol{a}_i' \boldsymbol{\beta}(t) = \sum_{j=1}^{k} a_{ij} \beta_j(t). \tag{3}$$

3

By considering $\beta_j(t)$, we can easily explain the effects of covariates. Thus, meaningful interval estimation and hypothesis testing can be carried out, e.g., a functional confidence interval for $\beta_j(t)$ can be determined and a hypothesis test of $\beta_j(t) \equiv 0$ can be carried out.

This paper is organized as follows: In §2 an estimate for the varying coefficient is discussed and its confidence interval is derived theoretically. Moreover, we propose a procedure to test for the significance of the effects of a covariate. Two numerical examples are presented in §3. §4 contains a discussion and our conclusions.

## 2. Estimation of varying coefficients and their evaluation

Now we are interested in estimating $\boldsymbol{\beta}(t)$ in (3). In fact, this is equivalent to estimating $\boldsymbol{\Theta}$ in (2). The estimator $\hat{\boldsymbol{\Theta}} = (\hat{\boldsymbol{\theta}}_1, \ldots, \hat{\boldsymbol{\theta}}_k)'$ for $\boldsymbol{\Theta}$ can be obtained from the results by Potthoff and Roy (1964) and Rao (1965) for the balanced case or Laird and Ware (1982), Vonesh and Carter (1987) for the unbalanced case. Let the asymptotic variance of $\hat{\boldsymbol{\theta}}_j$ $(j = 1, \ldots, k)$ be denoted by $\boldsymbol{\Omega}_j$, where $\boldsymbol{\Omega}_j = O(n^{-1})$. For both the balanced and unbalanced cases, the estimator satisfies $\sqrt{n}(\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j) \xrightarrow{d} N_q(\mathbf{0}, n\boldsymbol{\Omega}_j)$ $(n \to \infty)$ and an estimator $\hat{\boldsymbol{\Omega}}_j$ can also be derived for $\boldsymbol{\Omega}_j$. Let $\omega_j(t) = \boldsymbol{x}(t)'\boldsymbol{\Omega}_j\boldsymbol{x}(t)$ and $\hat{\omega}_j(t) = \boldsymbol{x}(t)'\hat{\boldsymbol{\Omega}}_j\boldsymbol{x}(t)$. From the result of Rao (1973, p.60), for any time point $t \in \mathbb{R}$, we have

$$\frac{\{\hat{\beta}_j(t) - \beta_j(t)\}^2}{\hat{\omega}_j(t)} = \frac{\{(\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j)'\boldsymbol{x}(t)\}^2}{\boldsymbol{x}(t)'\hat{\boldsymbol{\Omega}}_j\boldsymbol{x}(t)} \leq (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j)'\hat{\boldsymbol{\Omega}}_j^{-1}(\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j). \quad (4)$$

Note that $n\hat{\boldsymbol{\Omega}}_j$ is a consistent estimator of $n\boldsymbol{\Omega}_j$. Therefore, the right side of (3) converges to $\chi_q^2$ when $n$ tends to infinity because $\hat{\boldsymbol{\Omega}}_j^{-1/2}(\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j) \xrightarrow{d}$

$N_q(\mathbf{0}, \boldsymbol{I}_q)$ $(n \to \infty)$. Hence, we obtain the following equation:

$$\max_{t \in \mathbb{R}} \frac{\{\hat{\beta}_j(t) - \beta_j(t)\}^2}{\hat{\omega}_j(t)} \leq (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j)' \hat{\boldsymbol{\Omega}}_j^{-1} (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j) \xrightarrow{d} \chi_q^2. \tag{5}$$

Since the right side of (3) does not depend on $t$, we obtain a functional confidence interval for $\beta_j(t)$ as follows: Let $c_{q,\alpha}$ be the upper $\alpha \times 100$ percentage point of $\chi_q^2$, i.e., $\Pr(\chi_q^2 \geq c_{q,\alpha}) = \alpha$. A confidence interval for the varying coefficient $\beta_j(t)$ for $j = 1, \ldots, k$ can be obtained by

$$\mathcal{I}_{j,\alpha}(t) = \left[ \hat{\beta}_j(t) - \sqrt{\hat{\omega}_j(t) c_{q,\alpha}}, \quad \hat{\beta}_j(t) + \sqrt{\hat{\omega}_j(t) c_{q,\alpha}} \right]. \tag{6}$$

For any time point $t$, the coverage probability of the interval $\mathcal{I}_{j,\alpha}(t)$ is asymptotically given by $\Pr(\beta_j(t) \in \mathcal{I}_{\alpha,j}(t) : {}^\forall t \in \mathbb{R}) \geq 1 - \alpha$. Note that the confidence interval at a fixed time point $t$ can also be obtained from $\{\hat{\beta}_j(t) - \beta_j(t)\}^2 / \omega_j(t) \xrightarrow{d} \chi_1^2$, and this interval is narrower than the functional confidence interval given in (6).

Directly, we have a test statistic for the null hypothesis that a varying coefficient is zero for any time point, i.e.,

$$H_0 : \ \beta_j(t) = 0 \quad ({}^\forall t \in \mathbb{R}). \tag{7}$$

This implies the corresponding covariate has no effect on observations. Let us define $T_j(t) = \hat{\beta}_j(t)^2 / \omega_j(t)$ $(j = 1, \ldots, t)$. In order to test the hypothesis $H_0$ in (7), we have to consider the probability $\Pr(T_j(t) \leq x : {}^\forall t \in \mathbb{R})$. Note that

$$\Pr(T_j(t) \leq x : {}^\forall t \in \mathbb{R}) = \Pr\left( \max_{t \in \mathbb{R}} T_j(t) \leq x \right).$$

Now we define $T_j = \hat{\boldsymbol{\theta}}_j' \hat{\boldsymbol{\Omega}}_j^{-1} \hat{\boldsymbol{\theta}}_j$. The result in (4) implies that $\max_{t \in \mathbb{R}} T_j(t) \leq T_j$. Hence, $T_j$ is the test statistic for the null hypothesis $H_0$. Recall that $T_j$

is asymptotically null distributed as $\chi_q^2$. Therefore, the null hypothesis $H_0$ is rejected when $T_j > c_{q,\alpha}$, or the $p$-value can be calculated from the chi-square approximation $\Pr(\chi_q^2 > T_j)$.

Although $\hat{\beta}_j(t)$ can be obtained easily, $\hat{\Omega}_j$ is not always clearly defined. Therefore, before concluding this section, we describe $\hat{\Omega}_j$ for the two growth curve models.

BALANCED CASE: Note that $X_i$ can be rewritten as $X$ in the balanced case. The GCM proposed by Potthoff and Roy (1964) is defined by

$$y_i = X\Theta'a_i + \varepsilon_i, \quad (i = 1, \ldots, n),$$

where $\varepsilon_1, \ldots, \varepsilon_n \sim i.i.d. \ N_p(\mathbf{0}, \Sigma)$. Let $\hat{\Theta}$ and $\hat{\Sigma}$ be estimators of $\Theta$ and $\Sigma$, respectively. By using $\hat{\Theta}$ and $\hat{\Sigma}$, an estimator of $\Omega_j$ is given by

$$\hat{\Omega}_j = \frac{n-1}{n-(p-q)-1} e_j'(A'A)^{-1} e_j \otimes (X'\hat{\Sigma}^{-1}X)^{-1},$$

where $e_j$ is a $k$-dimensional basis vector with $j$th entry equal to 1 and the others equal to 0.

UNBALANCED CASE: The GCM proposed by Vonesh and Carter (1987) is defined by

$$y_i = X_i\Theta'a_i + \varepsilon_i, \quad (i = 1, \ldots, n),$$

where $\varepsilon_1, \ldots, \varepsilon_n$ are independent random vectors and $\varepsilon_i \sim N_{p_i}(\mathbf{0}, \sigma^2 I_{p_i} + X_i\Delta X_i')$. Let $\hat{\Theta}$, $\hat{\Delta}$ and $\hat{\sigma}^2$ be estimators of $\Theta$, $\Delta$ and $\sigma^2$, respectively. By using $\hat{\Theta}$, $\hat{\Delta}$ and $\hat{\sigma}^2$, an estimator of $\Omega_j$ is given by

$$\hat{\Omega}_j = (e_j' \otimes I_q) \left\{ \sum_{i=1}^n a_i a_i' \otimes (\hat{\Delta} + \hat{\sigma}^2 X_i'X_i)^{-1} \right\}^{-1} (e_j \otimes I_q).$$

## 3. Numerical Examples

3.1 *Male and Female Data*

We apply our estimation method to data considered by Potthoff and Roy (1964), consisting of measurements of the distance $(mm)$ from the center of the pituitary to the pteryomaxillary fissure for 11 girls and 16 boys at 4 different ages (8, 10, 12, 14 years old). Here we fit the model (1) to the data by letting $\boldsymbol{a}_i = (1,0)'$ for the girls and $(1,1)'$ for the boys and taking $\boldsymbol{x}(t) = (1, t, t^2)'$ for $t \in \{8, 10, 12, 14\}$. The first covariate expresses the distance for girls and the second covariate expresses the sex effect, which is determined by the additional distance for the boys beyond that of the girls. Since the design is balanced with respect to the time points, we have estimators of $\boldsymbol{\Theta}$, $\boldsymbol{\Omega}_1$ and $\boldsymbol{\Omega}_2$ from the result of Potthoff and Roy (1964), which are given by

$$\hat{\boldsymbol{\Theta}} = \begin{pmatrix} 17.096 & 0.537 & -0.003 \\ 4.946 & -0.852 & 0.053 \end{pmatrix},$$

$$\hat{\boldsymbol{\Omega}}_1 = \begin{pmatrix} 26.065 & -4.634 & 0.199 \\ -4.634 & 0.843 & -0.037 \\ 0.199 & -0.037 & 0.002 \end{pmatrix},$$

$$\hat{\boldsymbol{\Omega}}_2 = \begin{pmatrix} 43.985 & -7.820 & 0.335 \\ -7.820 & 1.423 & -0.062 \\ 0.335 & -0.062 & 0.003 \end{pmatrix},$$

respectively. Here we focus on the varying coefficient function for the sex effect. It is indicated in Figure 1, together with the 95% confidence interval curves. The significance level was calculated as $p = 0.002$.

3.2 *Mice body weight data*

Here we apply the proposed method to data about the body weights of mice (Watanabe *et al.*, 1996). One of our main concerns in these experiments was to determine whether neutron-induced genetic damage in parental germline

7

cells can affect the $F_1$ offspring growth or their body weight, which is an important index for measuring the health of mice. The parental mice received a single whole-body exposure to $^{252}$Cf neutrons at doses of 0, 50, 100, or 200 cGy. There were 1,232 body-weight observations of 124 mice. Body weights were measured at the same calendar time, but the birth dates for each mouse varied over two weeks. The observation ages for the mice are therefore not balanced, with age values varying between 2.03 and 13.55. The number of pairs of (male, female) mice having doses of 0, 50, 100, and 200 cGy were (4,4), (19,17), (22,25), and (19,14), respectively.

Because the variances of male body weights were greater than those of female body weights, the original weights were transformed by common logarithms as indicated in Figure 2. We thus take $y_{ij}$ to denote the logarithmic body weight $(g)$ of the $j$th measurement of the $i$th individual at an age of $t_{ij}$ months old for $i = 1, \cdots, 124$. The covariates of the $i$th individual are expressed as $\boldsymbol{a}_i = (a_{i1}, a_{i2}, a_{i3})'$ where $a_{i1} \equiv 1$, $a_{i2} = 1$ for a male $i$th mouse, $a_{i2} = 0$ for a female mouse, $a_{i3}$ is the parental dose of the $i$th mouse, and and $\boldsymbol{x}(t) = (1, t, t^2)'$ for $t_{ij}$.

Applying the estimation method proposed by Vonesh and Carter (1987), estimators of $\boldsymbol{\Theta}$, $\boldsymbol{\Omega}_j$, $j = 1, 2$ and 3 are obtained. The fitted curves and estimated time-varying coefficients are shown in Fig 2. The sex effect was significantly positive over the whole observation period, and it increased nearly monotonically. On the other hand, the dose effect was almost negative and seemed to be constant with respect to age. The significance level was calculated as $p = 0.004$. Therefore, the fitted curves in Fig 2(a) seem to be parallel with respect to dose level and are clustered by sex.

8

## 4. Conclusion and Discussion

In this paper, we have considered a new view of GCM, regarding $\boldsymbol{\theta}_j' \boldsymbol{x}(t)$ as a varying coefficient $\beta_j(t)$. Since elements of $\boldsymbol{\Theta}$ in the GCM are essentially coefficients of a polynomial function expressing the time trend, it is meaningless to focus only on a single element of $\boldsymbol{\Theta}$. Nevertheless, in the ordinary GCM, interval estimation and testing a hypothesis about only a single element of $\boldsymbol{\Theta}$ are usually considered, i.e., interval estimation for coefficients of polynomial function is carried out, or a test is performed to determine whether a coefficient of the polynomial function is zero. In real data analysis, we are interested in the effects of covariates on the amount of growth. Even if we find that, up to statistical significance, a coefficient of a polynomial function is 0, then, needless to say, it is not necessarily the case that the effect of the covariate is 0. Utilizing this new approach, that is, by varying the coefficient $\beta_j(t)$ over time, we can examine the effect over time of the covariate. Such an effect can be easily visualized by the proposed functional confidence interval for $\beta_j(t)$. Furthermore, by means of the proposed test procedure, we can determine the significance or otherwise of the effect of the covariate. Although it can be difficult to understand and interpret the results that are obtained, there is no doubt that GCM is a very useful model for describing a time trend, and our method may help in understanding results obtained from an actual GCM data analysis.

Our claim is that, by modifying GCM in this way, we are able to propose a new estimation procedure for a varying coefficient model in the case that the covariate is independent of time. It is well known that estimating a varying coefficient model is very complicated, and hence, many computational tasks

are required to obtain estimates. In particular, it has been impossible to construct a confidence interval for the entire varying coefficient function. On the other hand, our estimation method is very easy, i.e., estimates can be obtained from the results of the GCM, and so it is unnecessary to execute a special computer program to compute the functional confidence interval and the $p$-value of the test statistic. Therefore, our result may also be valuable in the sense of providing an easy estimation method for a varying coefficient model.

## References

HASTIE, T. AND TIBSHIRANI, R. (1993). Varying-coefficient models, *J. Roy. Statist. Soc. Ser.* **B**, **55**, 757–796.

LAIRD, N. M. AND WARE, J. H. (1982). Random-effects models for longitudinal data, *Biometrika*, **38**, 963–974.

POTTHOFF, R. F. AND ROY, S. N. (1964) A generalized multivariate analysis of variance model useful especially for growth curve problems, *Biometrika*, **51**, 313–326.

Rao, C. R. (1965). The theory of least squares when the parameters are stochastic and its application to analysis of growth curves, *Biometrika*, **52**, 447–458.

Rao, C. R. (1973). *Linear Statistical Inference and Its Applications* (2nd ed.), John Wiley, New York.

Vonesh, E. F. and Carter, R. L. (1987). Efficients inference for random-coefficient growth curve models with unbalanced data, *Biometrics*, **43**, 617–628.

Watanabe, H., Takahashi, T., Lee, J., Ohtaki, M., Roy, G., Ando, Y., Yamada, K., Gotoh, T., Kurisu, K., Fujimoto, N., Satow, Y. and Ito, A. (1996). Influence of paternal $^{252}Cf$ Neutron exposure on abnormal Sperm, embryonal Lethality, and Liver tumorigenesis in $F_1$ offspring of mice, *Japan. J. Cancer Research*, **87**, 51–57.

West, M., Harison, P. J. and Migon, H. S. (1985). Dynamic generalized linear models and Bayesian forecasting (with discussion), *J. Amer. Statist. Assoc.*, **80**, 73–97.

**Figure 1.** Estimated varying coefficient curve $\hat{\beta}_2(t)$ for male and female data. The broken lines represent 95% confidence intervals derived using Theorem 1. (a) The jointed gray lines show individual observations; the solid lines are fitted curves. (b) Sex effects: $\boldsymbol{\beta}_2(t)$ versus time.

**Figure 2.** Estimated varying coefficient curves $\hat{\beta}_j(t)$, $j = 1, 2$ and $3$ for mice body weight data. The dotted curves represent the $\pm 1.96$ standard error bands. The broken lines represent 95% confidence intervals derived using Theorem 1. (a) The jointed gray lines show individual observations; the solid lines are fitted curves. (b) Time effects: $\boldsymbol{\beta}_1(t)$ versus time. (c) Sex effects: $\boldsymbol{\beta}_2(t)$ versus time. (d) Dose effects: $\boldsymbol{\beta}_3(t)$ versus time.
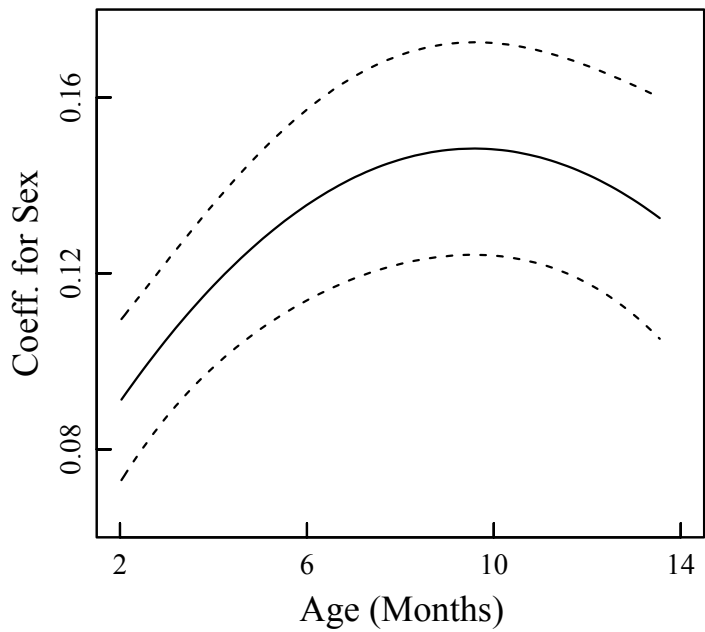
(a) Observation and Fitted Curves

(b) Sex Effect

Figure 1.

Figure 2.