

Two-Dimensional Surrogate Contact Modeling for Computationally Efficient Dynamic Simulation of Total Knee Replacements

Yi-Chung Lin

Raphael T. Haftka
Ph.D.

Department of Mechanical and Aerospace
Engineering,
University of Florida,
231 MAE-A Building,
P.O. Box 116250,
Gainesville, FL 32611-6250

Nestor V. Queipo
Ph.D.

Applied Computing Institute,
University of Zulia,
Maracaibo,
Zulia 4005, Venezuela

Benjamin J. Fregly¹
Ph.D.

Department of Mechanical and Aerospace
Engineering,
Department of Biomedical Engineering,
and Department of Orthopaedics and
Rehabilitation,
University of Florida,
231 MAE-A Building,
P.O. Box 116250,
Gainesville, FL 32611-6250
e-mail: fregly@ufl.edu

Computational speed is a major limiting factor for performing design sensitivity and optimization studies of total knee replacements. Much of this limitation arises from extensive geometry calculations required by contact analyses. This study presents a novel surrogate contact modeling approach to address this limitation. The approach involves fitting contact forces from a computationally expensive contact model (e.g., a finite element model) as a function of the relative pose between the contacting bodies. Because contact forces are much more sensitive to displacements in some directions than others, standard surrogate sampling and modeling techniques do not work well, necessitating the development of special techniques for contact problems. We present a computational evaluation and practical application of the approach using dynamic wear simulation of a total knee replacement constrained to planar motion in a Stanmore machine. The sample points needed for surrogate model fitting were generated by an elastic foundation (EF) contact model. For the computational evaluation, we performed nine different dynamic wear simulations with both the surrogate contact model and the EF contact model. In all cases, the surrogate contact model accurately reproduced the contact force, motion, and wear volume results from the EF model, with computation time being reduced from 13 min to 13 s. For the practical application, we performed a series of Monte Carlo analyses to determine the sensitivity of predicted wear volume to Stanmore machine setup issues. Wear volume was highly sensitive to small variations in motion and load inputs, especially femoral flexion angle, but not to small variations in component placements. Computational speed was reduced from an estimated 230 h to 4 h per analysis. Surrogate contact modeling can significantly improve the computational speed of dynamic contact and wear simulations of total knee replacements and is appropriate for use in design sensitivity and optimization studies. [DOI: 10.1115/1.3005152]

Introduction

For more than three decades, total knee replacement (TKR) surgery has been performed on patients with severe osteoarthritis. Though modern TKR surgery has a high success rate, some patients still experience substantial functional impairment postsurgery compared to healthy individuals [1]. Furthermore, patients today desire more than just pain relief, seeking a high level and broad variety of daily activities following TKR surgery (e.g., tennis, hiking, gardening, and even jogging) [2,3]. Thus, maximization of durability and minimization of functional limitations have become two primary design goals.

One way to pursue these goals is through the development of computational technology that permits sensitivity and optimization studies of new TKR designs. A recent step in this direction has been the development of computational models of knee wear simulator machines [4–8]. Such models permit dynamic contact and wear simulations of knee implant designs to be performed in a time- and cost-efficient manner, with model predictions being validated against experimental motion and wear measurements [4,6,8,9]. Though finite element models can be used to perform

contact calculations for these simulations, recent work has shown that elastic foundation contact models produce similar contact forces and pressures in a fraction of the CPU time [10]. Even so, CPU times on the order of 5–10 min for a one-cycle dynamic contact simulation can still be prohibitive for design sensitivity and optimization studies requiring thousands of simulations.

Computational limitations in other engineering disciplines have been overcome through the use of surrogate modeling approaches. Surrogate modeling (also known as “metamodeling” or “response surface approximation”) involves replacing a computationally costly model with a computationally cheap model constructed using data points sampled from the original model. Once the surrogate model is constructed, it is used in place of the original computationally costly model when subsequent engineering analyses are performed. A variety of surrogate model fitting methods have been proposed, including polynomial response surfaces [11–13], Kriging [14,15], radial basis functions [16,17], splines [18], and support vector machines [19,20]. While surrogate models have been utilized successfully for a variety of engineering applications [21–26], only a few studies have used surrogate models to fit input-output relationships from a computational contact model [27–29]. Furthermore, none of these studies has proposed a surrogate modeling technique to improve the computational efficiency of dynamic contact simulations.

Unlike traditional engineering applications, joint contact analyses pose unique challenges to existing surrogate modeling ap-

¹Corresponding author.

Contributed by the Bioengineering Division of ASME for publication of the JOURNAL OF BIOMECHANICAL ENGINEERING. Manuscript received November 19, 2007; final manuscript received August 22, 2008; published online February 4, 2009. Review conducted by Noshir Langrana.

proaches. Sampling data points within a hypercube defined by the upper and lower bounds of the relative pose parameters will result in few physically realistic data points for surrogate model creation. Most sample points will be either out of contact (i.e., zero contact force) or in excessive penetration (i.e., unrealistically high contact force) since contact forces are highly sensitive to the small displacements in the contact normal direction. Consequently, traditional surrogate-based modeling approaches fail to provide the accuracy needed to perform dynamic simulations. A special surrogate-based modeling approach adapted to the specific needs of joint contact analyses is therefore required.

This study proposes a novel surrogate contact modeling approach for performing computationally efficient dynamic contact and wear simulations of total knee replacements. The approach addresses the unique challenges involved in applying surrogate modeling techniques to joint contact problems. The primary objectives of the study are twofold. The first is to evaluate the computational speed and accuracy of the proposed surrogate contact modeling approach by comparing its results with those generated by an elastic foundation contact model. Sample points for the surrogate contact model are generated by the same elastic foundation contact model used in the comparative dynamic wear simulations. The second objective is to demonstrate the practical applicability of the proposed approach by analyzing the sensitivity of knee replacement wear predictions to realistic variations in input motions and loads and component placements. A Monte Carlo approach requiring thousands of dynamic wear simulations is used to perform the sensitivity study. Both objectives are pursued using a dynamic contact model of a total knee replacement constrained to planar motion in a Stanmore machine, and both demonstrate the significant computational benefits that surrogate contact modeling can provide for the analysis of TKR designs.

Methods

Surrogate Contact Model Development. We have developed a novel surrogate modeling approach to replace computationally costly contact calculations in dynamic simulations of total knee replacements. To develop the approach, we utilized a three-dimensional (3D) elastic foundation contact model [30–34] of a cruciate-retaining commercial knee implant (Osteonics 7000, Stryker Howmedica Osteonics, Inc., Allendale, NJ) constrained to sagittal plane motion in a Stanmore simulator machine [35]. The model possessed three degrees of freedom (DOFs) relative to ground: tibial anterior-posterior translation, femoral superior-inferior translation, and femoral flexion-extension. Similar to a Stanmore machine, femoral flexion-extension rotation was prescribed while the remaining two DOFs were loaded using ISO standard input curves [36]. A pair of parallel spring bumper loads was also applied to the tibia in the anterior-posterior direction to approximate the effect of ligament forces [35]. The EF contact model was treated as an applied load on the femoral component and tibial insert and used linear elastic material properties [4]. The dynamic equations for the model were derived via Kane's method using AUTOLEV symbolic manipulation software (OnLine Dynamics, Sunnyvale, CA). The equations were incorporated into a MATLAB program (The Mathworks, Natick, MA) that was used to perform forward dynamic simulations with the stiff solver ode15s.

To develop a surrogate contact model to replace the EF contact model in the larger dynamic model described above, we followed a four-step process: (1) design of experiments (DOE), (2) computational experiments, (3) surrogate model selection, and (4) surrogate model implementation (Fig. 1). The first step, DOE, is a statistical method for determining at which locations in design space computational experiments should be performed with the computationally costly model to sample the outputs of interest (Fig. 1(a)). For our sagittal plane knee implant model, locations in design space (or "sample points") were defined by the following three relative pose parameters: tibial anterior-posterior (AP) translation X , femoral superior-inferior (SI) translation Y , and femoral

flexion-extension rotation θ . The associated outputs of interest calculated by the EF contact model were the AP contact force F_x and the SI contact force F_y . Calculated contact torque was not of interest since θ was prescribed. To minimize the total number of sample points while maximizing the quality of the resulting surrogate model fit, we chose a Hammersley quasirandom (HQ) sampling method [37]. This method uses an optimal design scheme for placing sample points within a multidimensional hypercube and provides better uniformity properties than do other sampling techniques [38,39].

The second step, computational experiments, involved performing simulations with the original contact model at the sample points defined in the first step. Because contact forces are highly sensitive to some pose parameters but not others, we made one modification to the computational experiments performed at the HQ sample points (Fig. 1(b)). A traditional HQ sampling method would sample values of the three inputs X , Y , and θ within their specified bounds and use the EF contact model to calculate the two outputs F_x and F_y at each sample point. However, since SI contact force F_y is highly sensitive to small variations in SI translation Y [29,40], this approach produced a large number of sample points that were either out of contact or deeply penetrating. Inclusion of these unrealistic points in the surrogate model fitting process reduced the accuracy of the resulting surrogate contact model. To resolve this issue, we inverted the HQ sampling method for superior-inferior translation by sampling forces in the sensitive direction (i.e., F_y) and displacements in the insensitive directions (i.e., X and θ). Specifically, we used the HQ method to sample 300 pairs of X and θ values. For each of these pairs, we generated five sample points by performing five static analyses with the EF contact model using F_y values ranging from 600 N to 2400 N (approximately three times body weight) in 600 N increments plus a 10 N value for determining contact initiation. The outputs of each static analysis were F_x and Y . Evaluating 1500 sample points (X , F_y , and θ) using static analyses required approximately 3 h of CPU time on a 3.4 GHz Pentium IV PC.

The third step, surrogate model selection, involved determining the most appropriate type of surrogate model for fitting the input-output relationships defined by the computational experiments in the second step. Based on its ability to interpolate multidimensional nonuniformly spaced sample points, Kriging [15,41] was selected as the surrogate modeling method to fit the input-output relationships defined by the 1500 sample points. Originally developed for geostatistics and spatial statistics, Kriging has been widely applied in many different fields [28,42]. In mathematical terms, Kriging utilizes a combination of a polynomial trend model and a systematic departure model as follows:

$$y(\mathbf{x}) = f(\mathbf{x}) + z(\mathbf{x}) \quad (1)$$

where \mathbf{x} is a vector of input variables, $y(\mathbf{x})$ is the output to be fitted, $f(\mathbf{x})$ is a polynomial trend model, and $z(\mathbf{x})$ is a systematic departure model whose correlation structure is a function of the distances between sample points. The $f(\mathbf{x})$ term approximates the design space globally while the $z(\mathbf{x})$ term interpolates local deviations from $f(\mathbf{x})$. The necessary parameters for $z(\mathbf{x})$ were obtained through maximum likelihood estimation [15]. The MATLAB toolbox DACE [43] was used to construct all necessary Kriging models using a cubic polynomial trend function and a cubic spline correlation function based on preliminary tests.

The fourth step, surrogate model implementation, required the development of a special model inversion process to accommodate our modified HQ sampling approach (Fig. 1(c)). During a dynamic simulation, the EF contact model requires X , Y , and θ as inputs and provides F_x and F_y as outputs, whereas the surrogate contact model requires X , F_y , and θ as inputs and predicts F_x and Y as outputs. To address this inconsistency, we fitted F_y as a function of Y using the general form of a Hertzian contact model [44]

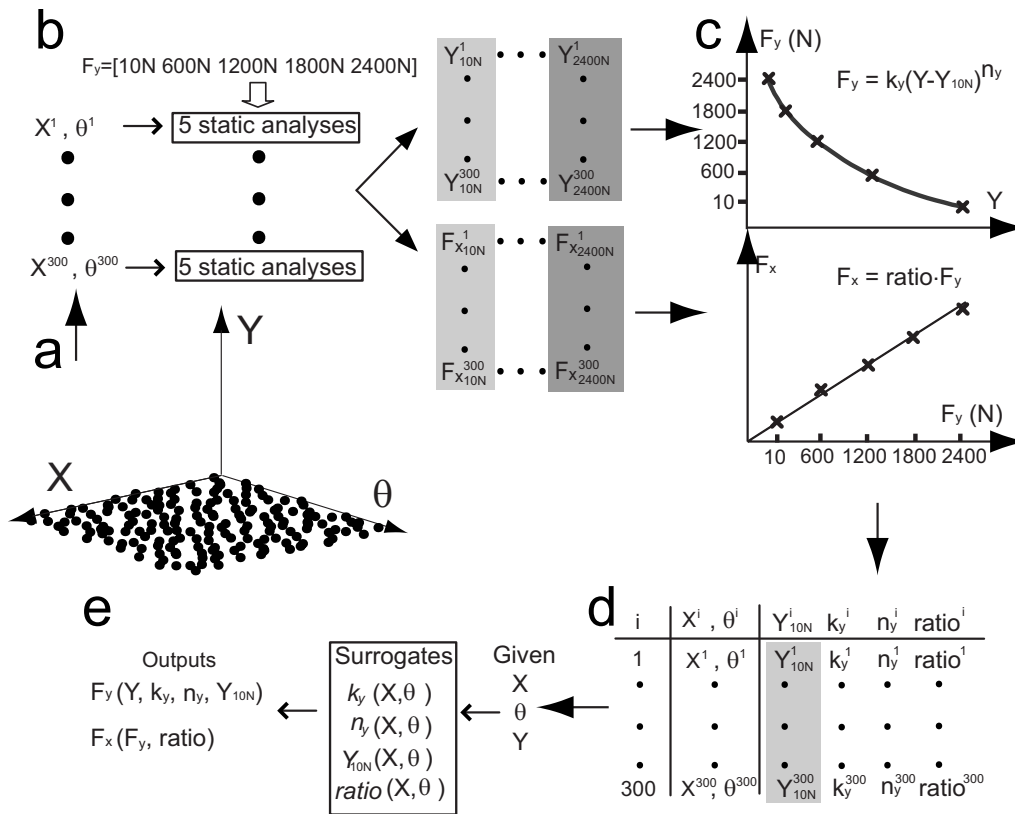


Fig. 1 Overview of surrogate contact model creation and implementation process. (a) 300 points (black dots) are sampled in the (X, θ) design space. (b) Five static analyses using five axial loads are performed at each sample point location using an elastic foundation contact model to calculate Y and F_x . (c) Generalized Hertzian contact theory is used to model F_y as a function of Y while F_x is modeled as proportional to F_y . (d) The values of Y_{10N} , k_y , n_y , and ratio from the 300 sample points are fitted as functions of X and θ using Kriging. (e) During a dynamic simulation, the surrogate contact model calculates F_y and F_x from four Kriging models given the current values of X , Y , and θ .

$$\hat{F}_y = k_y(X, \theta)[Y_{10N}(X, \theta) - Y]^{n_y(X, \theta)} \quad (2)$$

where \hat{F}_y is the predicted SI contact force, k_y is the contact stiffness, Y_{10N} is the SI translation at contact initiation, and n_y is the contact exponent. Use of Eq. (2) to calculate \hat{F}_y required fitting k_y , Y_{10N} , and n_y as a function of X and θ . Each of the 300 (X, θ) sample pairs had five sampled values of F_y (i.e., 10 N, 600 N, 1200 N, 1800 N, and 2400 N) and five associated output values of Y . Thus, Y_{10N} was already known for each (X, θ) sample pair. To find k_y and n_y for each (X, θ) sample pair, we linearized Eq. (2) by taking \log_{10} of both sides and solved for $\log_{10}(k_y)$ (and hence k_y) and n_y using linear least squares (five equations in two unknowns). This process yielded one value of k_y , Y_{10N} , and n_y for each (X, θ) sample pair.

To calculate F_x , we used the observation that F_x is close to a constant fraction of F_y at each (X, θ) sample pair as follows:

$$\hat{F}_x = \text{ratio}(X, \theta)\hat{F}_y \quad (3)$$

where \hat{F}_x is the predicted AP contact force, “ratio” is the average of the five (F_x, F_y) quotients obtained for each sample pair, and \hat{F}_y is the SI contact force predicted by Eq. (2). Thus, the final surrogate contact model utilized four separate Kriging models to fit k_y , Y_{10N} , n_y , and ratio as a function of X and θ (Fig. 1(d)), permitting the calculation of $\hat{F}_y = f(X, Y, \theta)$ and $\hat{F}_x = f(X, \theta, \hat{F}_y)$

from Eqs. (2) and (3) at any point during a dynamic simulation (Fig. 1(e)).

Surrogate Contact Model Evaluation. To evaluate the computational speed and accuracy of our surrogate contact modeling approach, we performed nine dynamic wear simulations with both the surrogate contact model and the EF contact model used to generate it. Each wear simulation utilized different nominal component placements that were variations away from the original placements (Fig. 2(a)). The specific variations simulated were all possible combinations of three SI locations of the femoral flexion axis (original ± 10 mm) and three AP slopes of the tibial insert (original ± 10 deg). To cover all possible dynamic simulation paths, we specified the ranges of the 300 HQ sampling points using the limits of motion obtained from dynamic simulations performed with the EF contact model when the SI location of the femoral flexion axis and the AP slope of the tibial insert were set to their extreme values. We quantified differences between the sets of results by calculating root mean square errors (RMSEs) in predicted contact forces, component translations, and wear volumes.

Generation of wear volume predictions with the surrogate contact model required the development of a separate surrogate modeling approach based on prediction of medial and lateral center of pressure (CoP) locations. This approach used Archard’s wear law [45] to predict wear volume for any specified number of cycles from the time history of contact forces and CoP slip velocities on each side over a one-cycle dynamic simulation. Specifically,

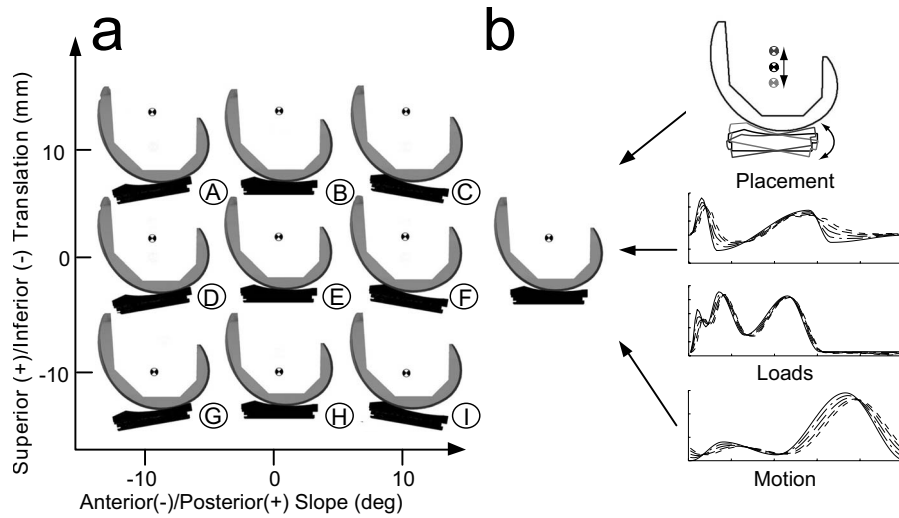


Fig. 2 Overview of global and local sensitivity analyses of predicted wear volume performed with the surrogate contact model. (a) Large variations in nominal component placements (labeled A–I) used for the global sensitivity analysis. (b) Small variations in component placements, input loads, and input motion used for local sensitivity analysis involving Monte Carlo methods applied to each nominal component placement.

$$V = N \sum_{i=1}^n \sum_{j=1}^2 k F_{ij} d_{ij} = N k \sum_{i=1}^n \sum_{j=1}^2 F_{ij} |\mathbf{v}_{ij}| \Delta t \quad (4)$$

where V is the total wear volume of the medial and lateral sides over N cycles, N is the assumed total number of cycles (5×10^6), i is a discrete time instant during the one-cycle dynamic simulation (1 through n), j is the side (medial or lateral), k is the assumed material wear rate ($1 \times 10^{-7} \text{ mm}^3/\text{N m}$ [46]), F_{ij} is the contact force magnitude at time instant i on side j , and d_{ij} is the sliding distance of the CoP at time instant i on side j . Due to symmetry of the implant geometry, the predicted value of F_{ij} on each side was taken as half the total contact force obtained from Eqs. (2) and (3). Prediction of d_{ij} on each side required determining the magnitude of the slip velocity \mathbf{v}_{ij} of the CoP and multiplying by the time increment Δt [4]. To find \mathbf{v}_{ij} , we used the kinematic concept of two points fixed on a rigid body [47]

$$\mathbf{v}_{ij} = {}^A \mathbf{v}^O + {}^A \boldsymbol{\omega}^B \times \mathbf{p}^{O\text{CoP}} \quad (5)$$

In this equation, A represents the tibial insert, B represents the femoral component, O is the femoral coordinate system origin located at the intersection of the flexion-extension and superior-inferior translation axes, CoP is the medial or lateral center of pressure location treated as a point fixed in body B , ${}^A \mathbf{v}^O$ is the velocity of point O with respect to reference frame A as determined from the relative translation of the femoral component origin, ${}^A \boldsymbol{\omega}^B$ is the angular velocity of body B with respect to reference frame A as determined from the flexion-extension of the femoral component, and $\mathbf{p}^{O\text{CoP}}$ is the position vector from point O to point CoP. While the values of ${}^A \mathbf{v}^O$ and ${}^A \boldsymbol{\omega}^B$ were obtained directly from the dynamic simulations, calculation of the CoP location on each side required six additional Kriging models (i.e., three components per side). Each Kriging model was created by fitting the average of the five CoP locations obtained from each (X, θ) sample pair. Wear volumes calculated by the surrogate contact model using the CoP-based approach were compared with those calculated by the EF contact model using both the CoP-based approach and a previously published element-based approach [4,8] to verify that the surrogate model wear volume predictions were accurate.

Surrogate Contact Model Application. As a practical application of the proposed surrogate contact modeling approach, we

performed two types of sensitivity analyses to investigate how machine setup issues affect predicted wear volume. The first was a global sensitivity analysis to investigate the effect of large variations in component placements within the ranges defined in Fig. 2(a) (i.e., $\pm 10 \text{ mm/deg}$). Specifically, we generated 400 combinations of component placements using 20 evenly spaced values of femoral flexion axis SI location and tibial insert AP slope. For each combination, the surrogate-based contact model was used to perform a dynamic simulation and predict wear volume. Percent changes in predicted wear volume were calculated relative to the maximum predicted value from all global sensitivity analysis results. The second category was a local sensitivity analysis to investigate the effect of small variations in nominal component placements and motion and load inputs by using Monte Carlo techniques (Fig. 2(b)). For the local sensitivity analysis, allowable variations in motion and load input curves were defined using experimental observations from eight different implant designs tested in a Stanmore simulator machine [36]. Deviations of the experimental input curves from the ISO standard curves were calculated and analyzed using principal component analysis [48], with the first two principal components being sufficient to account for over 98% of the variability in each type of curve. These two components were used to generate new deviation curves by selecting component weights as uniformly distributed random numbers within the bounds of the experimental curves. These deviation curves were added to the ISO standard curves to create new input motion and load curves as needed for the Monte Carlo analyses.

To distinguish between the effects of small errors in motion and load inputs and small errors in component placements, the local sensitivity analysis performed four types of Monte Carlo analyses for each of the nine nominal component placements for a total of 36 Monte Carlo analyses in all. Each type of Monte Carlo analysis varied input profiles and component placements separately or together. The first type varied motion and load profiles within 100% of their allowable ranges and also varied component placements within $\pm 1 \text{ mm/deg}$. The second type was similar but used smaller variations of only 10% of the allowable range for motion and load profiles and $\pm 0.1 \text{ mm/deg}$ for component placements. The third type varied only motion and load profiles within 100% of their allowable ranges and imposed no variations on component placements. Finally, the fourth type varied component placements within $\pm 1 \text{ mm/deg}$ but imposed no variations on motion and load

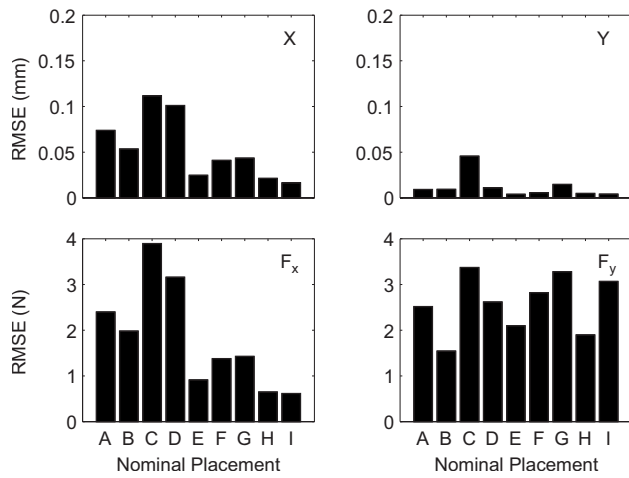


Fig. 3 Root-mean-square errors in joint motions and contact forces for each nominal component placement. Errors are computed by comparing dynamic simulation results generated with the elastic foundation contact model and the surrogate contact model.

profiles. For all types, percent changes in predicted wear volume were calculated relative to values obtained from the nominal component placements using the ISO standard inputs. Each Monte Carlo analysis was performed on a 3.4 GHz Pentium IV PC and involved at least 1000 dynamic contact simulations. The convergence criterion for each analysis was met when the mean and coefficient of variance (i.e., $100 \times$ standard deviation/mean) for the last 10% of the wear predictions were within 2% of the final mean and coefficient of variance [49].

Results

For the computational evaluation, the dynamic simulations performed with the surrogate contact model closely reproduced the planar motions, contact forces, and wear volumes predicted with the EF contact model for all nine component placements. On average, the RMSEs for planar motions and contact forces were less than 0.125 mm and 4 N, respectively (Fig. 3). While each dynamic simulation performed with the EF contact model required approximately 13 min of CPU time, each simulation performed using the surrogate contact model required 13 s. For both the surrogate contact model and the EF contact model, the CoP-based approach for predicting wear volumes reproduced the element-based EF results to within 0.5% error (Table 1).

For the practical application, the global and local sensitivity analyses revealed that large variations in component placements relative to the original locations and small variations in input motion relative to the ISO standard had a significant affect on predicted wear volume. For the global sensitivity analysis, predicted wear volume was sensitive to large variations in the SI location of the femoral flexion axis (maximum change of 26% when the AP slope of tibial insert held fixed), the AP slope of the tibial insert (maximum change of 13% when the SI location of femoral flexion axis held fixed), or both parameters simultaneously (maximum change of 33%) (Fig. 4). For the local sensitivity analysis, wear volume ranges were larger (mean of 18 mm³ compared to 2 mm³) and 50% percentile wear volume results were lower (mean of 23 mm³ compared to 31 mm³) when motion and load inputs and component placements were varied within 100% rather than 10% of their allowable ranges (Figs. 5(a) and 5(b)). When inputs and placements were varied separately within 100% of their allowable ranges, predicted wear volume was much more sensitive to motion and load inputs than to component placements (maximum change of 64% compared to 5%; Figs. 5(c) and 5(d)). Finally, when an additional global sensitivity analysis was performed with

Table 1 Comparison between wear volume predictions (mm³) made with the elastic foundation contact model and the surrogate contact model for each of the nine nominal component placements. Each prediction was made for 5×10^6 motion cycles. Both contact models can predict wear volume using a simple center of pressure-based approach, while only the elastic foundation contact model can make predictions using a more detailed element-based approach.

Nominal placement	Elastic foundation contact model		Surrogate contact model CoP-based
	Element-based	CoP-based	
A	37.07	37.05	36.98
B	34.03	34.02	33.91
C	39.04	39.03	38.87
D	32.91	32.89	32.88
E	30.78	30.76	30.77
F	34.20	34.19	34.14
G	29.37	29.36	29.32
H	26.45	26.44	26.38
I	29.37	29.35	29.31

the original component placements to separate the influence of the different inputs, predicted wear volume varied approximately linearly with small changes in each input curve but was more sensitive to changes in flexion-extension motion (maximum change of 47%; Fig. 6) than to changes in input loads (maximum change of 20% for F_x and 14% for F_y ; Figs. 6(a) and 6(b)). Each Monte Carlo analysis performed with the surrogate contact model required 4 h of CPU time compared to an estimated 230 h with the EF contact model.

Discussion

Computational speed is a major limiting factor for performing sensitivity and optimization studies of total knee replacement designs. Unlike constraint-based joint models, knee joint contact models require repeated geometry evaluations that consume the vast majority of the CPU time in a dynamic contact simulation [50]. This paper has presented a novel surrogate modeling approach for performing efficient dynamic contact simulations of human joints. The surrogate contact model was developed to replace the original EF contact model to avoid time-consuming repeated evaluation of the surface geometry during dynamic simu-

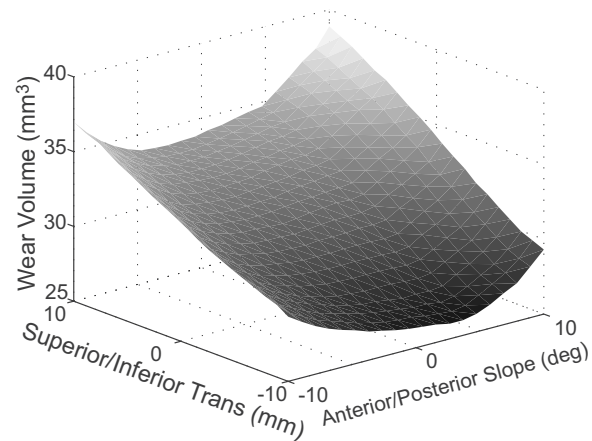


Fig. 4 Predicted wear volume as a function of large variations in component placements as calculated from 400 dynamic simulations performed with the surrogate contact model. Wear increases linearly when the femoral flexion axis is translated superiorly and quadratically when the tibial slope is increased anteriorly or posteriorly.

