

# Ball bearing diagnosis using a homogenous hybrid database in a supervised machine learning

Souleymane Sow<sup>1,2</sup>, Xavier Chimentin<sup>1</sup>, Lanto Rasolofondraibe<sup>2</sup>, Olivier Cousinard<sup>1</sup>

<sup>1</sup>Institute of Thermic, Mechanics, Material (IThMM), University of Reims Champagne-Ardennes, France

<sup>2</sup>Center of Research for Information and Communication Science and Technology (CReSTIC), University of Reims Champagne-Ardennes, France

[souleymane.sow@univ-reims.fr](mailto:souleymane.sow@univ-reims.fr)

## Abstract

Digital twin (DT) are often described as a virtual and dynamic representation of a system. They guarantee interaction between physical and virtual spaces. In the context of maintenance 4.0, the lack of historical data can be caused by an impossible instrumentation for complex systems. To face it, DT offers the possibility to simulate several operating modes which can serve for a diagnostic. This operation can be made by using machine learning algorithm (MLA) through a diagnosis by classification. But the challenge is to identify the best use of both data historical and simulated on a hybridisation database to make the most reliable diagnosis. In this paper, a digital twin combining a discrete element model (DEM) and a finite element model (FEM) is developed to generate data with an outer race default signature. These generated data with five sizes of defaults are also measured on the test bench. According to a percentage, historical data are used to build the homogenous hybrid database. Two MLAs (Support Vector Machines and K-Nearest Neighbours) are used to perform a classification by training the homogenous hybrid database and the test is realised by using the rest of historical data. The results of this approach show a better reliability than existing methods on the tested datasets also it's allowed to evaluate the contribution of historical data in homogenous hybridisation process.

**Keywords:** digital twin, bearings flexibility, hybrid model, homogeneous hybrid database, fault diagnosis, classification, artificial intelligence, SVM, KNN, supervised machine learning.

## 1 Introduction

Digital twins have emerged as a powerful tool for representing complex systems virtually and dynamically. They enable interaction between physical and virtual spaces and have found particular application in the field of maintenance 4.0. However, the lack of historical data due to the inability to instrument complex systems has been a significant challenge. To address this issue, digital twins offer the possibility of simulating several operating modes to aid in diagnostic procedures. Developing a digital twin requires reference data acquired on the test bench, which will subsequently be used to update the numerical model [1].

Ball bearings are often seen as the most critical component of rotating machinery. They are the object of a number of studies and many numerical models have been developed. In the paper [2], the authors presented various models of digital rolling elaborated in the literature. Machine learning algorithms (MLAs) have been used for classification-based diagnosis [3]–[5], but effectively using both historical and simulated data in a hybridisation database remains a challenge.

This paper presents a digital twin that combines a discrete element model (DEM) and a finite element model (FEM) to generate data with an outer race default signature. Data with five different sizes of defects are also measured on a test bench called “historical data”. According to a percentage, they are used to build a homogenous hybrid database with the entire generated data, and two MLAs (Support Vector Machines [6] and K-Nearest Neighbours [7]) are trained by the homogenous hybrid database. The test is carried out using

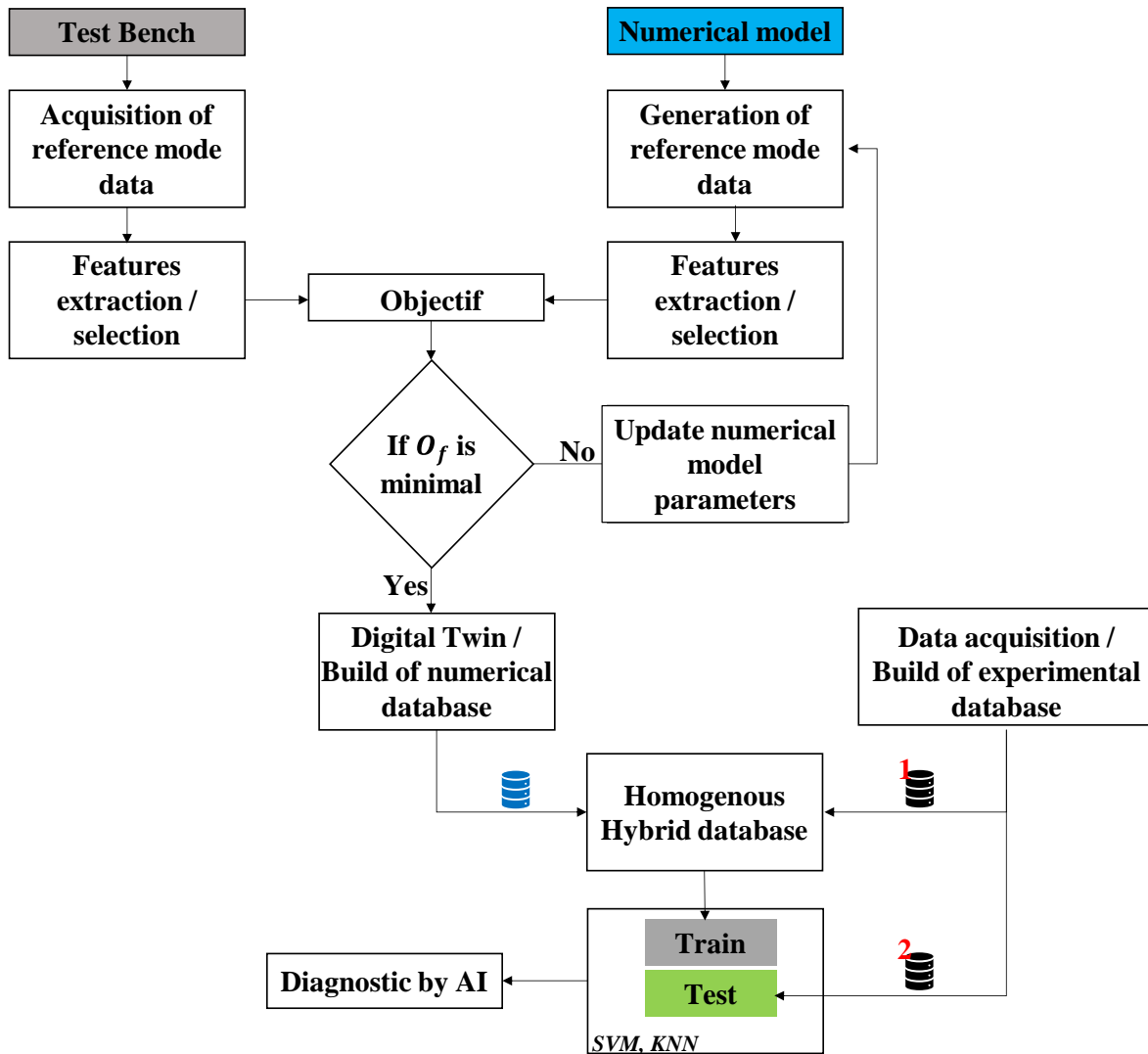
the rest of the historical data, and the contribution of historical data in the homogenous hybridisation process is evaluated. The results show better reliability than existing methods on the tested datasets.

The paper is divided into two sections in addition to the introduction and conclusion. Section 2 presents the methods used in this study. In section 3, the results of the application of the method are presented.

## 2 Methods

### 2.1 Global methodology

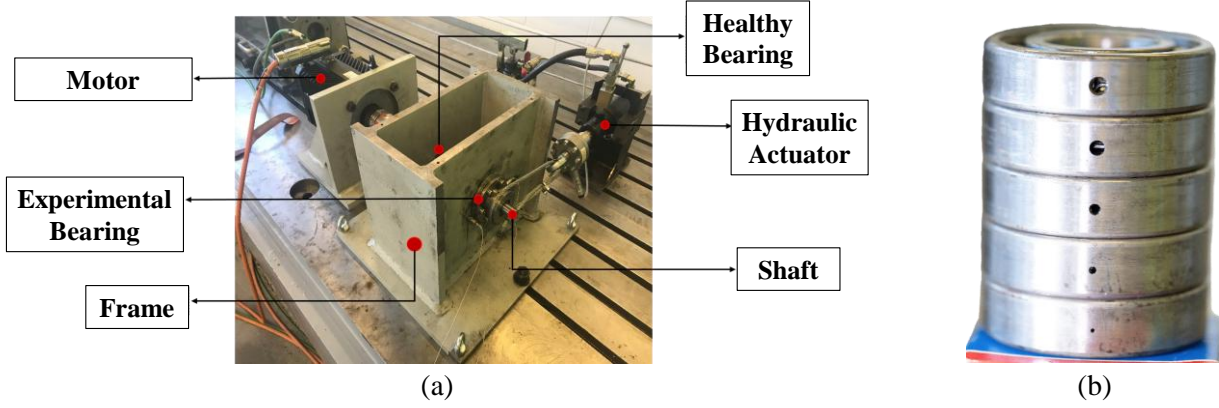
As shown in **Figure 1**, the approach employed in this study revolves around acquiring data from reference operating mode. Temporal and frequency indicators are then extracted from this signal. Subsequently, the numerical model of the test bench is developed to generate data from the reference operational mode as well. The numerical model combine a discrete element model (DEM) and a finite element model (FEM). These two models communicate with each other to form a single concentrated parameter model. By using the initially selected indicators and the experimental data, the numerical model is updated by minimizing the error defined within an objective function. Once the numerical model is updated, the remaining operating modes are generated and measured to build the respective numerical and experimental databases. The experimental database is divided into two. The numerical data is combined with the first numerical database to form a homogeneous hybrid database. To perform diagnostic classification, the hybrid database is used to train two machine learning algorithms (SVM and KNN). The other second of experimental data, which was not used in the hybridisation process, is employed for prediction testing.



**Figure 1** : Methodology based on training a hybrid database in a diagnosis by classification

## 2.2 Test bench and data acquisition

The test bench is a bearing module diagnosis as shown in **Figure 2a**. It consists of a frame and a steel shaft equipped with two SKF 6206 BC 32x64x18 ball bearings. The shaft is powered by a synchronous electric motor with a maximum output of 10 kW and is pulled by a hydraulic actuator connected to a steel cable. A variable speed drive, controlled by a PLC, regulates the rotation frequency of the shaft. The radial load can be manually adjusted, and a 300 daN force are applied to the rotating shaft at 1000 rpm. The signals obtained from the test bench are related to defects on the outer ring, ranging from 1 mm to 5 mm, as shown in **Figure 2b**. **Table 1** provides detailed specifications of the studied test bench.



**Figure 2** : Description: (a) Bearing module test bench; (b) Outer race fault severities from 1 mm to 5 mm

Parameters	Motor	Frame	Shaft	Bearing
Mass (kg)	8		2.5	
Moment of inertia (kg m <sup>2</sup> )	0.01125		0.0108	
Young's modulus (N m <sup>-2</sup> )		$2 * 10^{11}$	$2 * 10^{11}$	
Poisson ratio		0.3	0.3	
Density (kg m <sup>-3</sup> )		7900	7900	
Number of balls				9
Inner race radius (mm)				30
Outer race radius (mm)				62
Ball diameter (mm)				9.52
Pitch diameter				44
Backlash $\gamma$ ( $\mu$ m)				18
Bearing stiffness k (N m <sup>-1</sup> )				$8.5 * 10^9$

**Table 1** : Test bench characteristics

Experimental data are collected by using two (2) uniaxial piezoelectric accelerometers DJB A/120/VT with a sensitivity of 10 ( $\pm 10\%$ ) mv/g covering a frequency range of 34 kHz. A data collector type "OROS 36" is set at a sampling frequency of 51.2 kHz. Data are collected in 3.2 seconds and save into a historical database. For each data collected, a cut-off is done in 10 sub signals with 16384 points signals. Each acquisition provides 20 signals for each functional mode. That makes a total of 100 signals for the experimental database. For a signal having the signature of a defect on the outer race has a theoretical defect frequency as described in **equation 1** at 59.37 Hz. Where  $N_b$  is the number of balls,  $d$  is the ball diameter,  $D$  is the pitch diameter,  $\theta$  radial clearance and  $f_r$  the cage speed.

$$BPFO = \frac{N_b}{2} \left[ 1 - \left( \frac{d}{D} * \cos \theta \right) \right] f_r \quad (1)$$

As disciplines such as statistics, signal processing, and computer science continue to advance, techniques for extracting vibration features are constantly improving. One approach consists in working with the raw signal, without any filtering, and calculating statistical moments to detect faults in the system. Features are extracted in both time and frequency domains. In the time domain, the features extracted are: RMS, Peak, Kurtosis, Crest Factor, Skewness, Impulse Factor, Shape Factor, Average, Standard deviation, Talaf and the THIKAT. In the frequency domain, features such as: Mean frequency, Root mean square frequency, Frequency centre, Standard deviation frequency are extracted from the signals.

The Sequential Backward Selection (SBS) algorithm is used to reduce the features to the most relevant ones [8]. It allow to select the combination of features that provides the best separation data of different operating modes. In order to optimise time calculation and to represent each signal by a point in a 3D Space, the number of features to be selected by the SBS is limited to three (3).

### 2.3 Numerical model

The design is based on the following hypotheses: (i) the bearing elements are assumed to move in one plane; (ii) the angles between the balls are considered constant; (iii) the outer race is assumed to be fixed to the frame and the inner race to be fixed with a rigid contact to the transmission shaft; (iv) taking the flexibility of the bearings into consideration a finite element model of the frame is developed. The final model is a constitution of a DEM and FEM models. The model has nine (9) nodes each having three degrees of freedom (DOF), which are: flexion in the  $\vec{v}$  direction; flexion along the  $\vec{w}$  direction and a torsion in  $\vec{u}$  direction. The DEM model was proposed by Farhat *et al.* [9] based on Patil *et al.* and Harsha *et al.* works [10], [11]. As shown in **Figure 3a**, from the DEM are modelled the shaft, the motor, the ball bearings with the different operating modes. The shaft is modelled as a Timoshenko beam with nine (9) different nodes. The global shaft mass matrix  $[M_S]$  and the stiffness matrix  $[K_S]$  are extracted from shaft modelling. Associated to the node 1, the motor is modelled by its mass  $m_e$  and its polar inertia  $I_e$  in a global motor mass matrix  $[M_M]$ . Only the excitations forces due to the crushing of the ball to the races are considered in the equation of motion. As given in **equation 2** and **3**, the contact force depend on the angular position  $\theta_i$  for each ball, the radial backlash  $\gamma$ , the displacement on the  $\vec{v}$  and  $\vec{w}$  directions and the Hertz contact coefficient  $k$ . For a default located on the outer race of the bearing with a depth  $h$  and from an angle  $\varphi_d$ , a clearance is added when each ball pass on the default [12]. The added clearance  $Delta \varphi_i$  presents a periodicity at the Ball Pass Frequency Outer race (BPFO). To constitute the global forces vector  $\vec{F}$ , the motor torque  $C_m(t)$ , gravity forces are calculated in function of the shaft mass  $m$  and the gravity acceleration  $g$ , the radial load  $\vec{F}$  applied on the node 8 in the  $\vec{v}$  direction are included to vector containing the excitation forces in the nodes 3 and 6

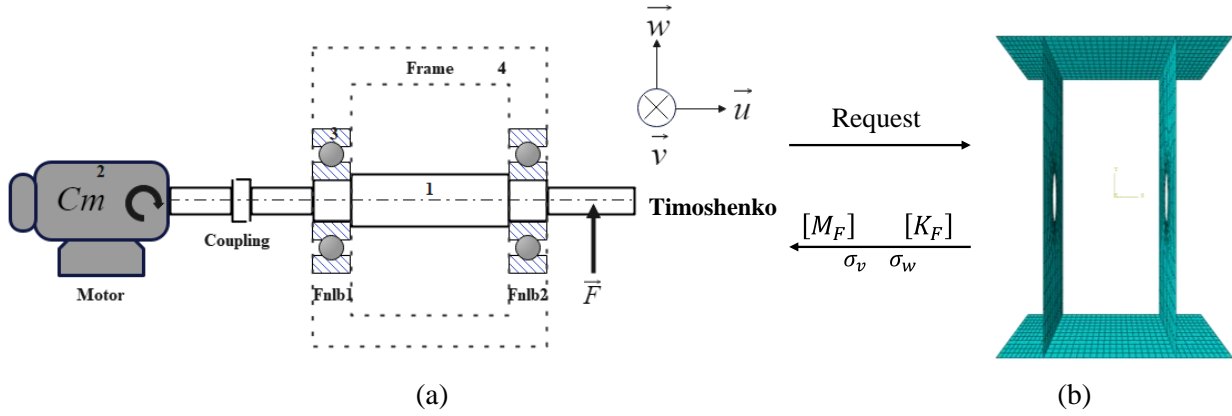
$$F_V(t) = \sum_{i=0}^{Nb} k[v(t) \cos(\theta_i(t)) + w(t) \sin(\theta_i(t)) - (\gamma + Delta \varphi_i(t))]^{\frac{3}{2}} + \cos(\theta_i(t)) \quad (2)$$

$$F_w(t) = \sum_{i=0}^{Nb} k[v(t) \cos(\theta_i(t)) + w(t) \sin(\theta_i(t)) - (\gamma + Delta \varphi_i(t))]^{\frac{3}{2}} + \sin(\theta_i(t)) \quad (3)$$

To highlight the contribution of modelling bearing flexibility, the frame is modelled on a finite element calculation software ABAQUS. Four-node shell elements (S4R) is used for the meshing type with a medial-axis to take consider the geometries on the lateral faces. From this mesh, 1534 elements are obtained. By fixing the base of the frame, the mass  $[M_F]$  and stiffness  $[K_F]$  matrices of the frame are extracted, as well as the stresses on the rings caused by the load applied to the shaft radially along the  $\vec{v}$  direction.

To reflect the system's dynamics, the equation of motion as given in **equation 4**, is reconstituted under MATLAB. Where  $[M]$  is the global system mass matrix with the assembled different mass  $[M_S]$ ,  $[M_M]$  and  $[M_F]$ . The global system stiffness  $[K]$  is a sum of different stiffnesses  $[K_S]$ ,  $[K_F]$ .  $[C]$  is the global damping matrix, calculated proportionally to the mass matrix  $[M]$  and the average value of the stiffness matrix  $[K]$  as given in **equation 5** [13].  $\alpha$  and  $\beta$  two real Rayleigh coefficients calculated to make sure that, the damping is

viscous with  $\alpha = 0.6$  and  $\beta = 6 * 10^{-4}$ .  $q$  is the generalized vector of coordinates that are defined by the degrees of freedom of each node  $q = \{v_n, w_n, \theta_n\}, n \in (1..9)$ . The flow chart presented in Patil *et al.* work [10] is used to solve the equation of motion.



**Figure 3 :** Numerical model : (a) overview; (b) finite element model

$$[M]\ddot{q} + [C]\dot{q} + [K]q = F(t) \quad (4)$$

$$[C] = \alpha[M] + \beta[K] \quad (5)$$

As the aim is to make a diagnosis by hybridising the data, so the characteristics of the generated data should be at least in the same order as those of the experimental data. The digital twin should reflect the actual behaviour of the physical machine, which means that it is necessary to update the model parameters. Additionally, to ensure that the numerical model accurately reflects the behavior of the physical machine, it is imperative to update the model parameters. By referring to the work of Wang *et al.*[14], the numerical model is update by calculating the objective function. The error between  $sf_i$  and  $mf_i$  the relevant  $i^{\text{th}}$  features extracted respectively from the simulated and measured signals given in **equation 6**, where  $p = \{E, \nu, \alpha, \beta\}$  represents the vector of the parameters to be optimised in the model (model parameters). Choosing these parameters is a crucial step that greatly depends on the results of the update. Updated values are required to be physically acceptable. The permissible upper and lower limits, UB and LL are defined as  $\pm 5\%$  of the initial parameters. The signal acquired for the case of 4 mm of a severity fault is chosen as the reference signal. It is used to recalibrate the numerical model.

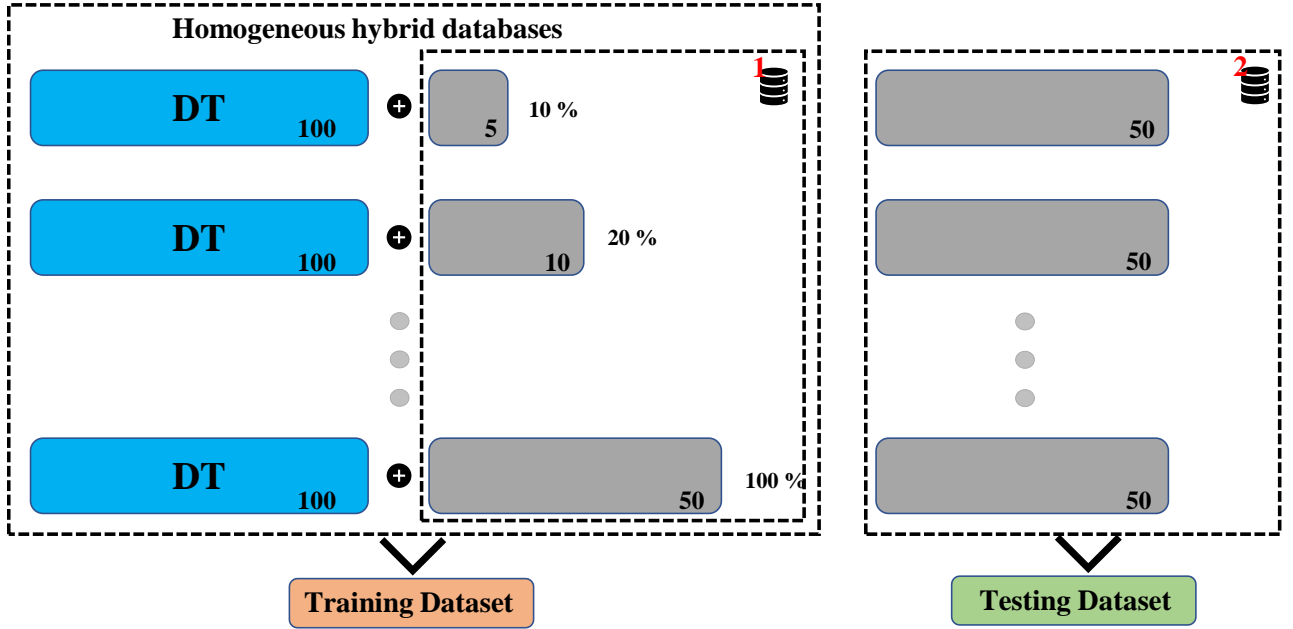
$$O_f = \min ||\{sf_i\} - \{mf_i\}||_2^2 \quad (6)$$

## 2.4 Homogenous hybridisation model

In the presented study, a single hybrid model is created by merging experimental data with data generated by a numerical model. The homogeneity of this hybrid model relies on the similarity of the fused data from the experimental and numerical databases, which should correspond to the same operating modes. The homogeneous hybridisation model, as illustrated in **Figure 4**, is carefully selected.

During each iteration, the experimental database is initially divided into two parts. The first part, containing 10% to 100% of the data for each severity, is used to construct the hybrid model. This portion is then combined with the complete numerical database, and both are employed to train the two MLAs, namely SVM (Support Vector Machines) [15] and KNN (K-Nearest Neighbours) [16].

The second experimental database, composed of the same operating modes but excluded from the hybridisation process, is exclusively dedicated to conducting prediction tests. It serves as an independent dataset for evaluating the performance of the MLAs.

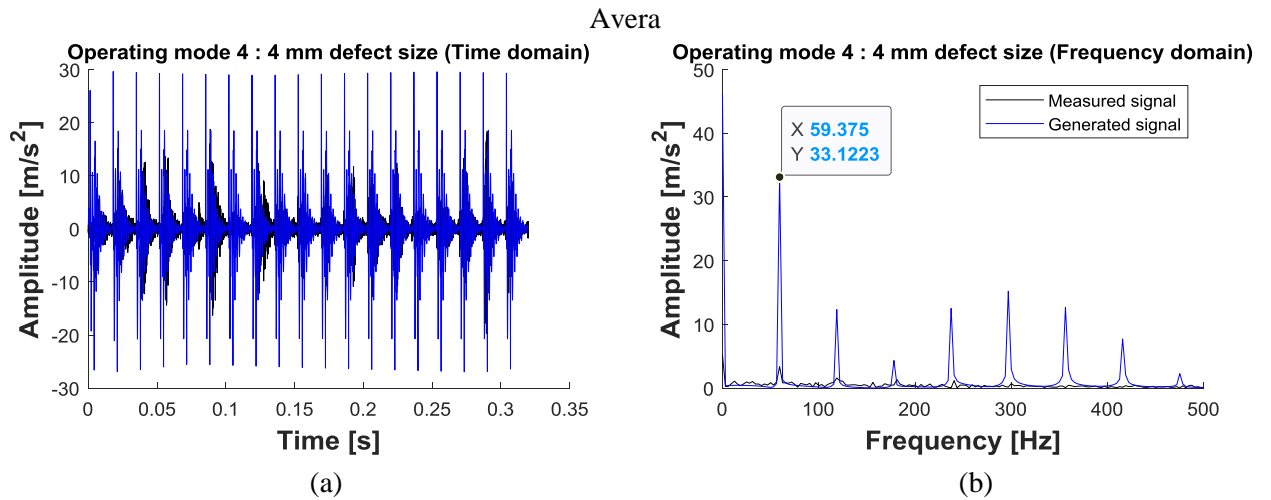


**Figure 4 :** Construction of the homogenous hybrid database methodology

### 3 Results

The selected indicators, SK, SF, and Talaf, have been found to be crucial for the analysis of the system's behavior. However, when generating data using a non-updated numerical model, significant differences are observed between the generated data and the measured data on the test bench (**Figure 5**). The discrepancies are mainly evident in the time domain, with the amplitude of the generated signal being 1.5 times higher than the measured signal. In the frequency domain, a power spectrum envelop is examined in the range of [0-500] Hz for spectral analysis. The outer race defect's frequency is present at 59.375 Hz in both signals, but the magnitude and spectral distribution show variations. The observed difference between the generated and historical signals justifies the need to update the numerical model. To achieve this, the error is minimised as shown in **equation 7** and the optimization problem is solved through the use of the lsqnonlin MATLAB function. The method is particularly valuable for tackling nonlinear problems characterized by an objective function featuring numerous local minima or strong parameter correlation.

$$\text{Er}(\{E, v, \alpha, \beta\}) = \left\{ \begin{matrix} SK_s \\ SF_s \\ Talaf_s \end{matrix} \right\} - \left\{ \begin{matrix} SK_m \\ SF_m \\ Talaf_m \end{matrix} \right\} \quad (7)$$

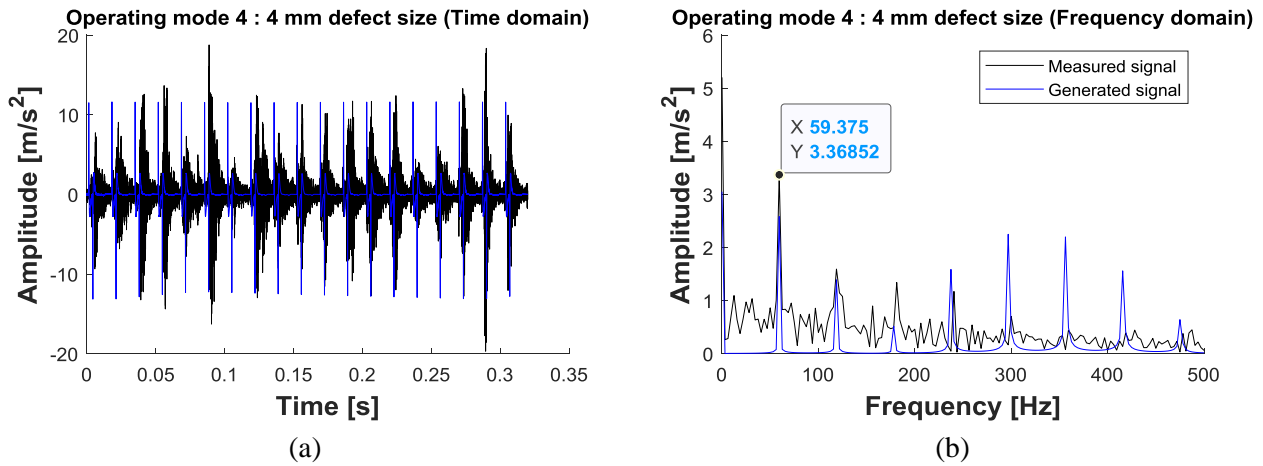


**Figure 5 :** Data generated with a non-updated model for operating mode 4 describing a defect size of 4 mm: (a) data representation in time domain; (b) data representation in frequency domain with a BPFO at 59.375 Hz

By generating fresh parameter values listed in **Table 2**, it allows the creation of new data sets across various operating modes. In fact, the updated signal of operating mode 4, is shown in both the time and frequency domains in **Figure 6**. The signals are now in the same order.

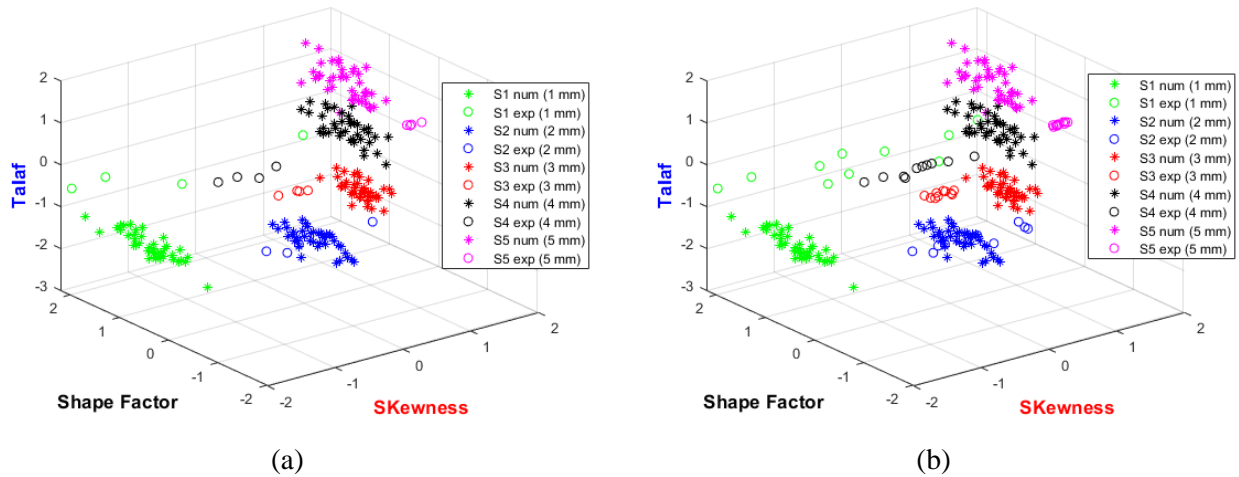
Updated parameters	Symbols	Before update	After update
Young's modulus ( $\text{N m}^{-2}$ )	$E$	$2.5 * 10^{11}$	$2.38 * 10^{11}$
Poisson ratio	$\nu$	0.3	0.25
Rayleigh coefficient	$\alpha$	0.5	0.54
Rayleigh coefficient	$\beta$	$6 * 10^{-4}$	$5.4 * 10^{-4}$

**Table 2** : The numerical model optimized parameters



**Figure 6** : Data generated with an updated model for operating mode 4 describing a defect size of 4 mm: (a) data representation in time domain; (b) data representation in frequency domain with a BPFO at 59.375 Hz

In **Figure 7a**, a homogeneous hybridisation is shown with a 10 % contribution from the experimental data of each operating mode. In **Figure 7b**, a hybridisation with a 60% contribution from experimental data is displayed.



**Figure 7** : Spatial representation of training hybridisation database: (a) 10 % experimental data contributions in homogenous hybrid database; (b) 60 % experimental data contributions in homogenous hybrid database

The accuracy of the test is depicted in **Figure 8**, demonstrating the relationship between the rate of added experimental data to the numerical database and the achieved accuracy using the hybridisation method. As the contribution of experimental data increases, the diagnostic accuracy improves accordingly. A comparison



is presented between the non-updated and updated models, where data from both models are utilized for hybridisation with experimental data. The results demonstrate a progressive increase in SVM accuracy as the percentage of experimental data used in hybridisation rises. Starting with 0% experimental data, the accuracy is 90%. However, as the percentage of introduced experimental data increases, the SVM accuracy improves accordingly. At 60% experimental data, the accuracy reaches 99%, and it reaches 100% when 70% to 100% of the experimental data is used. Similarly, the results reveal that KNN accuracy gradually improves with an increasing percentage of experimental data used in hybridisation. Initially, with 0% experimental data, the accuracy is 88%. As the percentage of introduced experimental data increases, the KNN accuracy also improves. There is an increase up to 97% with 60% experimental data, followed by a stabilization at 100% accuracy when 70% to 100% of the experimental data is used. These outcomes underscore the positive impact of experimental data on the performance of the KNN model, even though its influence appears slightly lower compared to that observed for the SVM.

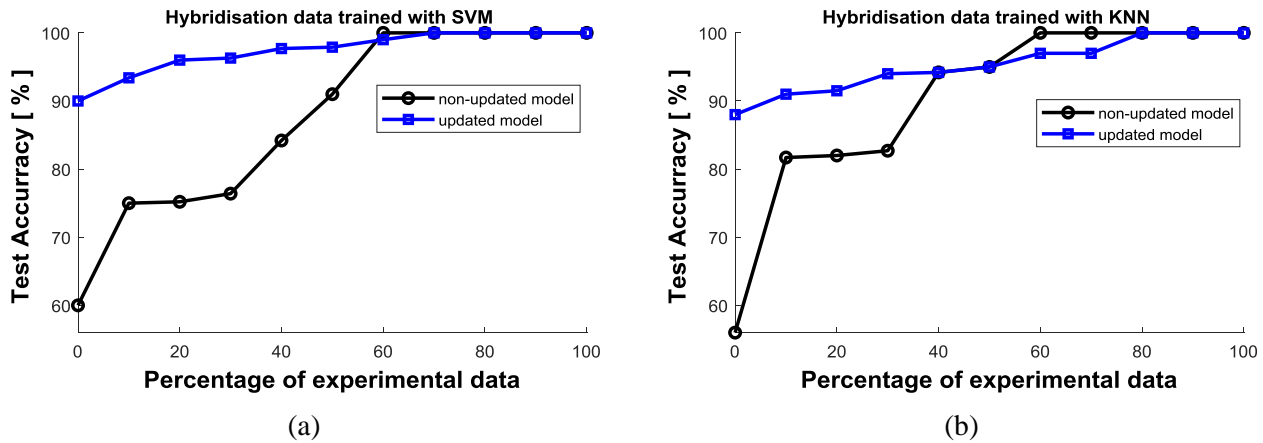


Figure 8 : Classification accuracy: (a) SVM; (b) KNN

## 4 Conclusion

This study highlights the significance of digital twin technology and machine learning algorithms in addressing the challenges associated with fault diagnosis in complex systems. By leveraging hybrid databases and integrating various modeling techniques, this research opens avenues for further advancements in the field and encourages continued exploration of innovative methods for accurate and efficient fault detection and prognosis.

The work proposes a numerical model of a test bench composed of bearings, created by combining finite element and discrete element models. The model was used to diagnose bearing faults of varying severities on the outer ring. The approach involved training on a hybrid database and testing on the remaining experimental database. The results show the reliability of using a hybrid database in diagnostic classification, with an average accuracy of 97.3% for the Support Vector Machine (SVM) and 95.2% for the K-Nearest Neighbours (KNN) algorithm. The high accuracy achieved by the SVM and KNN algorithms demonstrates their effectiveness in accurately identifying and classifying different fault conditions in the bearings. These results further validate the use of digital twin technology and machine learning algorithms for fault diagnosis in complex systems.

Moving forward, future research could focus on developing methods for constructing and optimizing hybrid databases for different types of systems, such as gearbox modules, and faults, including inner race defects. Additionally, exploring advanced feature selection techniques and incorporating other machine learning algorithms may contribute to improving diagnostic accuracy even further.

## Acknowledgment

The University of Reims Champagne-Ardenne and the Grand-Est region are acknowledged for their support in funding this study, and the authors wish to express their sincere gratitude towards them.



## References

- [1] Q. Qi and F. Tao, 'Digital Twin and Big Data Towards Smart Manufacturing and Industry 4.0: 360 Degree Comparison', *IEEE Access*, vol. 6, pp. 3585–3593, 2018, doi: 10.1109/ACCESS.2018.2793265.
- [2] A. Das and S. Ray, 'A Review on Diagnostic Techniques of Bearing Fault and its modeling in Induction Motor', in *2020 IEEE Calcutta Conference (CALCON)*, Feb. 2020, pp. 502–505. doi: 10.1109/CALCON49167.2020.9106511.
- [3] S. Tyagi and S. K. Panigrahi, 'A DWT and SVM based method for rolling element bearing fault diagnosis and its comparison with Artificial Neural Networks', *JACM*, vol. 3, no. 1, Apr. 2017, doi: 10.22055/jacm.2017.21576.1108.
- [4] E. Lahner *et al.*, 'Artificial neural networks in the recognition of the presence of thyroid disease in patients with atrophic body gastritis', *World J Gastroenterol*, vol. 14, no. 4, pp. 563–568, Jan. 2008, doi: 10.3748/wjg.14.563.
- [5] P. Tamilselvan and P. Wang, 'Failure diagnosis using deep belief learning based health state classification', *Reliability Engineering & System Safety*, vol. 115, pp. 124–135, Jul. 2013, doi: 10.1016/j.res.2013.02.022.
- [6] K.-B. Duan, J. C. Rajapakse, and M. N. Nguyen, 'One-Versus-One and One-Versus-All Multiclass SVM-RFE for Gene Selection in Cancer Classification', in *Evolutionary Computation, Machine Learning and Data Mining in Bioinformatics*, E. Marchiori, J. H. Moore, and J. C. Rajapakse, Eds., in Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2007, pp. 47–56. doi: 10.1007/978-3-540-71783-6\_5.
- [7] G. Guo, H. Wang, D. Bell, Y. Bi, and K. Greer, 'KNN Model-Based Approach in Classification', in *On The Move to Meaningful Internet Systems 2003: CoopIS, DOA, and ODBASE*, R. Meersman, Z. Tari, and D. C. Schmidt, Eds., in Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2003, pp. 986–996. doi: 10.1007/978-3-540-39964-3\_62.
- [8] S. J. Reeves and Z. Zhe, 'Sequential algorithms for observation selection', *IEEE Transactions on Signal Processing*, vol. 47, no. 1, pp. 123–132, Jan. 1999, doi: 10.1109/78.738245.
- [9] M. H. Farhat, F. Chaari, X. Chiementin, F. Bolaers, and M. Haddar, 'Dynamic Remaining Useful Life Estimation for a Shaft Bearings System', in *Smart Monitoring of Rotating Machinery for Industry 4.0*, F. Chaari, X. Chiementin, R. Zimroz, F. Bolaers, and M. Haddar, Eds., in Applied Condition Monitoring. Cham: Springer International Publishing, 2022, pp. 169–178. doi: 10.1007/978-3-030-79519-1\_11.
- [10] M. S. Patil, J. Mathew, P. K. Rajendrakumar, and S. Desai, 'A theoretical model to predict the effect of localized defect on vibrations associated with ball bearing', *International Journal of Mechanical Sciences*, vol. 52, no. 9, pp. 1193–1201, Sep. 2010, doi: 10.1016/j.ijmecsci.2010.05.005.
- [11] S. P. Harsha, K. Sandeep, and R. Prakash, 'The effect of speed of balanced rotor on nonlinear vibrations associated with ball bearings', *International Journal of Mechanical Sciences*, vol. 45, no. 4, pp. 725–740, Apr. 2003, doi: 10.1016/S0020-7403(03)00064-X.
- [12] M. H. Farhat, X. Chiementin, F. Chaari, F. Bolaers, and M. Haddar, 'Digital twin-driven machine learning: ball bearings fault severity classification', *Meas. Sci. Technol.*, vol. 32, no. 4, p. 044006, Feb. 2021, doi: 10.1088/1361-6501/abd280.
- [13] R. G. Parker, S. M. Vijayakar, and T. Imajo, 'Non-linear dynamic response of a spur gear pair: modelling and experimental comparisons', *Journal of Sound and Vibration*, vol. 237, no. 3, pp. 435–455, Oct. 2000, doi: 10.1006/jsvi.2000.3067.
- [14] J. Wang, L. Ye, R. X. Gao, C. Li, and L. Zhang, 'Digital Twin for rotating machinery fault diagnosis in smart manufacturing', *International Journal of Production Research*, vol. 57, no. 12, pp. 3920–3934, Jun. 2019, doi: 10.1080/00207543.2018.1552032.
- [15] J. Huang, X. Hu, and X. Geng, 'An intelligent fault diagnosis method of high voltage circuit breaker based on improved EMD energy entropy and multi-class support vector machine',

*Electric Power Systems Research*, vol. 81, no. 2, pp. 400–407, Feb. 2011, doi: 10.1016/j.epsr.2010.10.029.

- [16] Q. Wang, Y. B. Liu, X. He, S. Y. Liu, and J. H. Liu, ‘Fault Diagnosis of Bearing Based on KPCA and KNN Method’, *Advanced Materials Research*, vol. 986–987, pp. 1491–1496, 2014, doi: 10.4028/www.scientific.net/AMR.986-987.1491.